



**HAL**  
open science

# Reconfigurable Logic Architectures based on Disruptive Technologies

Pierre-Emmanuel Gaillardon

► **To cite this version:**

Pierre-Emmanuel Gaillardon. Reconfigurable Logic Architectures based on Disruptive Technologies. Other. Ecole Centrale de Lyon, 2011. English. NNT : 2011ECDL0027 . tel-00674438

**HAL Id: tel-00674438**

**<https://theses.hal.science/tel-00674438>**

Submitted on 27 Feb 2012

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## *Thèse de l'Université de Lyon*

*Délivrée par l'Ecole Centrale de Lyon  
Spécialité : Dispositifs de l'Electronique Intégrée  
Soutenue publiquement le 15 Septembre 2011*

*par*

*M. Pierre-Emmanuel Gaillardon*

*Ingénieur CPE-Lyon*

*Master Recherche GEGP « Dispositifs de l'Electronique Intégrée »*

*Préparée au Laboratoire d'Intégration Silicium des Architectures Numériques  
(CEA-LETI-DACLE)*

# ***Reconfigurable Logic Architectures based on Disruptive Technologies***

*Ecole Doctorale Electronique, Electrotechnique, Automatique*

### *Composition du jury :*

<i>M. A. Souifi,</i>	<i>Professeur, INSA Lyon,</i>	<i>en qualité de Président</i>
<i>M. S. Mitra,</i>	<i>Associate Professor, Stanford University,</i>	<i>en qualité de Rapporteur</i>
<i>M. L. Torres,</i>	<i>Professeur, Université Montpellier II,</i>	<i>en qualité de Rapporteur</i>
<i>M. J.O. Klein,</i>	<i>Professeur, Université Paris Sud,</i>	<i>en qualité d'Examineur</i>
<i>M. I. O'Connor,</i>	<i>Professeur, Ecole Centrale de Lyon,</i>	<i>en qualité de Directeur</i>
<i>M. F. Clermidy,</i>	<i>HDR, CEA-LETI,</i>	<i>en qualité d'Encadrant</i>



# *Acknowledgments*

---

While this thesis is written in English, I feel nice to write most of these acknowledgments in French, as it will enforce my thanks.

Afin de débiter ces remerciements, je vous propose une petite anecdote. En fin d'année 2007, la veille de la fête des lumières, je me trouvais en entretien à l'Ecole Centrale de Lyon en vue d'une bourse de thèse. Je rencontrais ainsi deux personnes qui présentèrent chacun leurs travaux pendant 30 minutes. Autant dire qu'après cette heure de présentation (par deux ténors de la communication scientifique), je fis bien pâle figure lorsque je dus me présenter. A cette époque, je n'avais que pour seul arme la motivation. Près de 4 ans après, j'espère que ces deux personnes sont satisfaites du travail accompli. Pour ma part je le suis et ce fût un véritable honneur que d'avoir travaillé avec le **Professeur Ian O'Connor** et le **Docteur Fabien Clermidy**. Que dire à part vous remercier bien chaleureusement ? Il y aurait évidemment trop à narrer mais je peux tout de même insister sur la patience de Ian pour ses nombreuses corrections orthographiques, la pleine maîtrise du marketing de titre de Fabien et la pertinence de vos encadrements laissant une grande part à l'autonomie. Je ne parle même pas de vos qualités humaines tant de fois montrées jusqu'aux tréfonds de Gaslamp quarter.

Ces deux personnes sont les premières garantes de la pertinence de ce travail, mais une thèse ne serait rien sans un jury de grande qualité. Je tiens tout d'abord à remercier le **Professeur Lionel Torres** et le **Professeur Subhasish Mitra** pour avoir accepté de rapporter ce travail. Il en va de même pour le **Professeur Abdelkader Souifi** et le **Professeur Jacques-Olivier Klein** pour avoir accepté de l'examiner. It is my pleasure to have this thesis reviewed by such a significant board. Thank you very much for having accepted this charge.

Avant de remercier mes collaborateurs, je souhaiterais dédier un paragraphe à **Haykel Ben Jamaa**. Ce fût un réel plaisir que de travailler avec lui pendant la durée de son post-doc et après son embauche. Haykel est une personne de très grande qualité sur le plan professionnel et humain. J'ai énormément appris à ses côtés du point de vue de la technique, de la méthodologie (Haykel a toujours une idée en stock à creuser et toujours une solution à un problème !) et de la communication (Haykel est le roi de la publication !).

Cette thèse, de part les diverses technologies utilisées, a été le fruit de nombreuses collaborations internes comme externes. Sur le plan interne, j'ai eu grand plaisir à travailler avec les différentes équipes du LETI. Je remercierai ainsi **Marina Reyboz** pour sa gentillesse et son expertise dans les modèles compacts, **Giovanni Betti Beneventi** et **Luca Perniola** pour leur collaboration sur les mémoires à changement de phase, **Perrine Batude** pour son prometteur procédé 3-D monolithique, sa sympathie et sa salsa, **Thierry Poiroux** pour sa connaissance sur le graphène et sa gestion du projet SNOW, **Paul-Henry Morel** et **Thomas Ernst** pour leurs procédés à base de nanofils et leur ouverture d'esprit lors de chaque collaboration concepteurs/technologues, **Gérald Cibrario** et **Marjorie Gary** pour leur superbes DKs (Vive la PDK Team !), **Karim Azizi Mourier** le champion du terminal sans qui il ne serait pas possible d'avoir un tel environnement de CAD, **Sigrid Thomas** et **Véronique Robert**, nos chères ingénieurs brevets, toujours disponibles pour entendre une idée. Un grand merci également à **Thomas Gauthier** et à **Houcine Oucheikh**, les deux stagiaires d'école d'ingénieur que j'ai eu le plaisir de co-encadrer et qui ont réalisé un travail



appréciable. Dans le cadre du projet Nanograin, des collaborations ont été entretenues avec des partenaires tels que l'Institut des Nanotechnologies de Lyon et le Laboratoire de l'Intégration du Matériau au Système à Bordeaux. J'eus ainsi le grand plaisir de travailler avec **Nataliya Yakymets, Kotb Jabeur, Sébastien Le Beux, Sébastien Frégonèse, Cristell Maneux** et **Thomas Zimmer**. Je remercie également le Professeur **Giovanni De Micheli**, de l'Ecole Polytechnique Fédérale de Lausanne, pour m'avoir invité à donner un séminaire au sein de son groupe et proposé un poste de collaborateur scientifique. Je suis très heureux de pouvoir travailler au sein du LSI à compter de l'automne.

Je remercie également toutes les personnes que j'ai eu l'occasion de cotoyer pendant ces années de thèse, et tout particulièrement les membres de mon laboratoire : notre chef **François Bertrand**, les permanents **Didier V., Christian, Denis, Frédéric, Jérôme, Olivier, Alexandre, Edith** (A quand l'organisation d'une « special session » et d'un tutorial Salsa ?), **Yvain, Pascal, Ivan, Sébastien, Romain, Didier L., Bastien** et non permanents qui se sont succédés au sein du laboratoire, tout particulièrement mes amis **Riadh, Imen, Pallavi, Santhosh, Rodolphe, Antoine, Bilel, Jean-Philippe** et **Maxime** ; les membres du laboratoire L3I et particulièrement **Anais, Bertrand** et **Rémi** ; et du laboratoire LE2TH, particulièrement **Cyril, Patrick, Stéphanie, Frédéric, Grégory** et **Hubert**.

Toute cette belle organisation ne serait rien sans le devouement de l'administration. J'ai été heureux de mener ces travaux au sein du département DACLE sous la direction de **Jean-René Lequepeys** puis de **Thierry Collette**, mais également de profiter des conseils et de l'expérience de **Marc Belleville** et d'**Ahmed Jerraya**. Je remercie tout particulièrement **Armelle** et **Caroline**, nos adjointes de direction qui sont de vraies fées.

Durant ces années de thèse, j'ai eu l'opportunité d'effectuer des enseignements au sein de l'école d'ingénieur CPE Lyon, mais également au sein des classes préparatoires CPE. Je remercie ainsi **Nacer Abouchi, Evelyne Steffen, François Bois, Thierry Tixier, François Joly** et **Renaud Daviot** pour m'avoir proposé ces enseignements. Je remercie également **Tahar Limane** et **Jean-Marie Becker** toujours disponibles pour répondre à des questions théoriques.

Sortant à présent du cadre purement institutionnel, je voudrais faire une place de choix à **Sophie** pour avoir relu une grande partie de ce manuscrit. Merci Doudou ! Je crois que l'un comme l'autre, nous souhaitons également en profiter pour remercier et féliciter les Le Bars pour leur mariage (☺).

J'aimerai également remercier mes amis plongeurs Grenoblois. Arrivé en thèse comme simple plongeur, j'ai mis à profit ces quelques années pour évoluer jusqu'au monitorat. Ce fût très agréable d'appartenir au « Subatomic » où je salue mes amis **Jean-Benoit, Frédéric, Jean-Louis, Sylvain, Hervé, Sandrine, Laetitia, Jean-Christophe** et tous les autres. Je remercie également tous mes MF2 préférés pour m'avoir mené jusqu'au MF1 : **Corinne, Sébastien, Fabien** et **Pascal**. C'est toujours un grand plaisir de partager quelques moments sous l'eau avec vous tous.

L'occasion est bonne également pour saluer mes amis qui subissent toujours indirectement les affres de la thèse. Je remercie tout particulièrement la Ouff Team pour avoir assuré le robot en plein dans ma période de rédaction. Merci donc à **Cyril, Philippe, Anne-Sophie** et **Rémi**. Merci tout particulièrement à **Cyril** car lui plus que les autres a subi un grand nombre de fois mes retards pour cause de publication ou de manuscrit à terminer.

Enfin, je remercie ma famille pour tout le soutien qu'ils m'apportent. Toujours présents, j'espère leur apporter un peu de fierté (La palette de Champagne parle d'elle-même !) en compensation de tout ce qu'ils font pour moi. Merci **Maman, Papa, Pépé, Tonton F, Tatie D, Tonton J, Tatie V** et **Zouzine**. Je pense également fort à mes deux grand-mères au moment d'écrire ces lignes.

# *Table of Contents*

---

<b>CHAPTER 1</b>	Introduction.....	1
1.1	The Scaling Era .....	1
1.2	New Challenges.....	3
1.2.1	Process Integration and Devices .....	3
1.2.2	Systems and Tools.....	4
1.3	Architectural and Technological Opportunities .....	4
1.3.1	Reconfigurable Architectures.....	5
1.3.2	Technological Evolutions.....	5
1.3.2.1	MOS Extensions .....	6
1.3.2.2	Emerging Devices and Memories.....	6
1.4	The Nanoscale Chicken and Egg Dilemma.....	6
1.5	Work Methodology.....	7
1.6	Research Contributions.....	8
1.6.1	Methodology and Tools .....	8
1.6.2	Memories and Routing Resources for FPGA.....	9
1.6.3	Logic Blocks for FPGA.....	9
1.6.4	Nanoscale Architectures.....	10
1.7	Organization of the Thesis.....	10
<b>CHAPTER 2</b>	Background and Motivation .....	13
2.1	Conventional Reconfigurable Architecture Overview .....	14
2.1.1	The FPGA Architecture .....	14
2.1.1.1	Generalities .....	14
2.1.1.2	General Architectural organization.....	15
2.1.1.3	Logic Block Architecture.....	16
2.1.1.4	Routing Architecture.....	18
2.1.1.5	Architectural Parameterization and Optimum .....	19
2.1.2	Market Trends .....	20
2.1.2.1	Increase of the CLBs Complexity.....	20
2.1.2.2	Transition from Homogeneity to Heterogeneity.....	21
2.1.2.3	Non-Volatility Features .....	22

2.1.3	Limitations of FPGAs .....	23
2.2	Emerging Reconfigurable Architectures Overview .....	24
2.2.2	Morphic Approach .....	24
2.2.3	Heterogeneous Approach .....	25
2.2.3.1	Regular Architectures .....	25
2.2.3.2	Enhanced Reconfigurable Architectures.....	27
2.2.4	Global Comparisons and Discussions .....	28
2.3	Architectural Formalization and Template.....	29
2.3.1	Global Statement .....	29
2.3.2	Generic Architectural Template .....	30
2.4	Conclusion and Work Position .....	31
<b>CHAPTER 3</b>	<b>Innovative Structures for Routing and Configuration.....</b>	<b>33</b>
3.1	Context and Objectives.....	34
3.1.1	Context Position .....	34
3.1.2	Objectives.....	35
3.2	Proposal 1: Resistive memory technologies.....	35
3.2.1	Introduction .....	35
3.2.2	Phase Change Memory Properties and Technological Assumptions.....	36
3.2.2.1	Physical Phenomena .....	36
3.2.2.2	Technological Assumptions.....	37
3.2.2.3	Opportunities.....	37
3.2.3	Elementary Memory Node .....	38
3.2.3.1	Concept .....	38
3.2.3.2	Programming Circuitry .....	39
3.2.3.3	N-ary Logic Generalization .....	40
3.2.3.4	Technological Integration Opportunities .....	40
3.2.4	Routing Elements .....	41
3.2.4.1	Concept .....	42
3.2.4.2	Programming Circuitry .....	42
3.2.5	Functional validation and performance characterization .....	43
3.2.5.1	Methodology .....	43
3.2.5.2	Compact Modeling.....	43
3.2.5.3	Transient Simulations .....	44
3.2.5.4	Performance Estimation.....	45
3.2.5.5	Discussion .....	47
3.3	Proposal 2: Monolithic 3-D Integration Process .....	48
3.3.1	Introduction .....	48
3.3.2	Technological Assumptions .....	49

3.3.3	Elementary Blocks .....	50
3.3.3.1	3-D Memory.....	51
3.3.3.2	LUT Impact.....	51
3.3.3.3	Routing Structure Impact.....	51
3.3.4	Performance Characterization .....	54
3.3.4.1	Methodology .....	54
3.3.4.2	SRAM Configuration Performance .....	54
3.3.4.3	LUT Performance .....	55
3.3.4.4	Cross point Performance.....	55
3.3.4.5	Comments .....	56
3.4	Proposal 3: Vertical Silicon Nanowire FET Process.....	56
3.4.1	Introduction .....	56
3.4.2	Technological Assumptions .....	57
3.4.3	Vertical 3-D Logic .....	58
3.4.3.1	Smart Vias.....	58
3.4.3.2	Logic Gates .....	58
3.4.3.3	SRAMs.....	59
3.4.4	Performance Characterization .....	60
3.4.4.1	Methodology .....	60
3.4.4.2	Device Characterization.....	61
3.4.4.3	NOT Gate Characterization .....	62
3.4.4.4	Discussion .....	63
3.5	Global Comparisons and Discussions .....	63
3.6	Conclusion .....	64
<b>CHAPTER 4</b>	<b>Architectural Impact of 3-D Configuration and Routing Schemes.....</b>	<b>67</b>
4.1	Motivation and Global Methodology .....	68
4.2	Benchmarking Tool for FPGA-like Architectures .....	68
4.3	Resistive Memory-based FPGA Performance.....	69
4.3.1	PCM-FPGA Architecture.....	69
4.3.2	Methodology .....	70
4.3.3	Simulation of Large Circuits .....	70
4.3.3.1	Impact on the Routing Structures .....	71
4.3.3.2	Impact on the Configuration Memories.....	71
4.3.4	Impact of Technologies.....	71
4.3.5	Discussion .....	72
4.4	3-D Monolithic Integrated FPGA Performances.....	73
4.4.1	Overall FPGA Architectural View.....	73
4.4.2	Methodology .....	73

4.4.3	Simulation of Large Circuits .....	74
4.4.3.1	Monolithic 3-D Integration Impact on Global Area .....	74
4.4.3.2	Monolithic 3-D Integration Impact on Critical Path Delay .....	74
4.4.4	Discussion .....	75
4.5	Vertical NWFETs-based FPGA Performances .....	75
4.5.1	Real 3-D Routing for FPGAs .....	75
4.5.2	Methodology .....	76
4.5.3	Simulation of Large Circuits .....	76
4.5.4	Discussions .....	77
4.6	Discussion .....	77
4.7	Conclusion .....	78
<b>CHAPTER 5</b>	<b>Disruptive Logic Blocks .....</b>	<b>81</b>
5.1	Context and Objectives .....	82
5.2	Proposal 1: Ambipolar Carbon Electronics .....	83
5.2.1	Introduction .....	83
5.2.2	CNT-based Devices Properties and Opportunities .....	84
5.2.3	Technological Assumptions and Device Modeling .....	86
5.2.3.1	Process Flow and Layout Proposal .....	86
5.2.3.2	DG-CNFET Compact Model .....	87
5.2.3.3	Process Tuning .....	87
5.2.4	Functionality improvement : In-Field Reconfigurability .....	88
5.2.5	Performance Evaluation .....	89
5.2.5.1	DG-CNFET Evaluation Methodology .....	89
5.2.5.2	Silicon CMOS Performance Evaluation Methodology .....	90
5.2.5.3	Simulation Results .....	91
5.2.6	Dynamic Logic Circuit Improvement .....	93
5.2.6.1	Concept .....	93
5.2.6.2	Buffer Cell Illustration .....	93
5.2.6.3	Density Improvement Example .....	94
5.2.7	Discussion .....	94
5.3	Proposal 2: 1-D Silicon Crossbars .....	95
5.3.1	Introduction .....	95
5.3.2	Technological Assumptions .....	96
5.3.3	SiNWFET Configurable Logic Cell .....	96
5.3.3.1	Nanowire Crossbar Logic .....	96
5.3.3.2	Multiplexer Design Methodology .....	97
5.3.3.3	Layout Design .....	98
5.3.4	Performance Evaluations .....	99

5.3.4.1	Methodology .....	99
5.3.4.2	Electrical Simulations .....	99
5.3.4.3	Discussion .....	100
5.4	Proposal 3: Lithographic Crossbars.....	101
5.4.1	Introduction .....	101
5.4.2	Technological Assumptions .....	101
5.4.3	Logic Design Methodology.....	102
5.4.3.1	Inverter construction .....	102
5.4.3.2	Generalization to complex logic circuits .....	102
5.4.4	Performance Evaluation .....	103
5.4.4.1	Methodology .....	103
5.4.4.2	Performance Estimations with Parasitic-Free Circuits .....	104
5.4.4.3	Impact of the Width of Horizontal Wires Including Parasitics.....	106
5.4.4.4	Impact of the Parasitic Capacitances .....	106
5.4.4.5	Impact of Scaling Including Parasitics .....	107
5.4.4.6	Performance Estimations including Parasitics.....	108
5.4.4.7	Discussion and Guidelines.....	108
5.5	Global Comparisons .....	110
5.6	Conclusion.....	110
<b>CHAPTER 6</b>	<b>Disruptive Architectural Proposals and Performance Analysis.....</b>	<b>113</b>
6.1	Introduction .....	114
6.2	Architectural Proposals.....	115
6.2.1	Introduction .....	115
6.2.2	MCluster Organization.....	115
6.2.2.1	MCluster: Adding a “Logic Layer” .....	115
6.2.2.2	MCluster: Simplifying the Interconnect Overhead.....	116
6.2.3	BLE and CLB organization.....	117
6.2.4	FPGA Final Organization.....	117
6.3	Benchmarking Tool .....	118
6.3.1	General Overview of the Flow .....	119
6.3.2	MPACK: Matrix PACKer .....	120
6.3.2.1	Concept .....	121
6.3.2.2	Mapping Algorithm : Architectural Optimization .....	121
6.3.2.3	Mapping Algorithm : Brute-Force Mapper.....	124
6.3.2.4	Clustering Algorithm .....	126
6.4	Evaluation of Fixed Interconnect Topologies .....	127
6.4.1	Methodology .....	128
6.4.2	Packing Success Rate .....	128

6.4.3	Fault Tolerance.....	129
6.4.3.1	Physical Origin of Defects at Nanoscale .....	129
6.4.3.2	Analysis Protocol and Results .....	129
6.4.4	Average Interconnect Length .....	130
6.4.5	Discussion .....	131
6.5	Global Architecture Evaluations .....	131
6.5.1	Methodology .....	131
6.5.2	Evaluation Results.....	131
6.5.2.1	Impact of the Granularity.....	131
6.5.2.2	Impact of Latch Depopulation .....	132
6.5.2.3	Performance Comparison with CMOS .....	133
6.5.3	Discussion .....	134
6.6	Conclusion.....	137
<b>CHAPTER 7</b>	<b>Conclusions and Future Work .....</b>	<b>139</b>
7.1	Global Comparison.....	139
7.2	Conclusions and Contributions.....	140
7.2.1	Global Conclusions .....	140
7.2.2	Contributions to Methodology and Tools .....	141
7.2.3	Contributions to Memories and Routing Resources for FPGA.....	141
7.2.4	Contributions to Logic Blocks for FPGA .....	142
7.2.5	Contributions to Architectures .....	143
7.3	Future Work.....	143
7.3.1	Design Generalization and Technology Heterogeneity .....	143
7.3.2	Fault Tolerance.....	144
References	.....	i
List of Publications	.....	xv
List of Acronyms	.....	xvii
<b>APPENDIX A</b>	<b>Memento of Presented Technologies.....</b>	<b>xxi</b>
<b>APPENDIX B</b>	<b>Résumé Etendu .....</b>	<b>xxxiii</b>

# List of Figures

---

Figure 1.	Processor transistor count and Moore's Law since 1971 [4] .....	2
Figure 2.	Projected evolution of dimensions and mask cost [6].....	3
Figure 3.	Spectrum of implementation solutions for data processing algorithms .....	5
Figure 4.	Electronic systems hierarchical organization and opportunities [7] .....	6
Figure 5.	The Chicken and Egg Dilemma of the disruptive technology .....	7
Figure 6.	Template based Fast Evaluation methodology.....	8
Figure 7.	Island-style Field Programmable Gate Array organization.....	16
Figure 8.	Transistor-level programmability for FPGA [13] .....	16
Figure 9.	Simple gate FPGA logic block [14] .....	17
Figure 10.	Complex gate FPGA logic blocks – a) Actel Act-1 LB [15] b) Quicklogic pASIC1 LB [16] .....	17
Figure 11.	Configurable Logic Block architecture [17] .....	18
Figure 12.	Basic Logic Element architecture [17].....	18
Figure 13.	Details on the routing organization in an island-style FPGA [17].....	19
Figure 14.	Virtex-6 logic block [20].....	21
Figure 15.	Xilinx Spartan-3A architecture organization [31].....	22
Figure 16.	Schematic of a merged programming-data path floating gate transistor [35]....	23
Figure 17.	Magnetic Tunnel Junction unbalanced flip flop (a) and associated above IC structure (b) [37].....	23
Figure 18.	Field Programmable Gate Arrays area repartition per block [38].....	24
Figure 19.	NanoPLA architecture [154] .....	25
Figure 20.	Dynamic logic implementation of AND, NAND, OR and NOR functions [65] .....	26
Figure 21.	NASIC structure (implementing a 1-bit adder) [65] .....	26
Figure 22.	Structure of the programmable NWFET (a), associated characterization (b) and nanowire-nanowire coupled multigate device (c) [50].....	27
Figure 23.	XOR function example built with four ambipolar transistors [52] .....	28
Figure 24.	a) Transistor level schematic and b) configuration table for CNT reconfigurable cell [122] .....	28
Figure 25.	Layered organization of the generic template architecture .....	31
Figure 26.	Schematic of PCM device in Logic 0 (named RESET) and Logic 1 (named SET) configurations and of the programming pulses suitable to obtain the states. ....	37
Figure 27.	Cross sectional schematic showing a PCM device integration. ....	38
Figure 28.	Logic-in-PCM elementary memory node .....	38
Figure 29.	Node in read configuration.....	39
Figure 30.	Node in write configuration .....	39
Figure 31.	Line sharing illustration in standalone-memory-like architecture .....	40



Figure 32.	Ternary logic node based on PCMs .....	40
Figure 33.	Cross-sectional view of the memory node using standard process integration.....	41
Figure 34.	Full Back-End integration illustration of the PCM-based memory node using an oxide BJT .....	41
Figure 35.	(a) PCM-based 2x2 switchbox architecture (b) zoom on the cross point structure (c) example of programmed switchbox (memories in off-state are not shown for clarity) .....	42
Figure 36.	Example of switchbox programming sequence .....	43
Figure 37.	Input (a) and output (b) drivers for PCM-based crossbar .....	43
Figure 38.	Synopsis of the behavioural model with the five modules.....	44
Figure 39.	Comparison between the model and measurements of I-V characteristics of the amorphous and crystalline phases of a PCM cell.....	44
Figure 40.	Binary memory node transient simulation .....	45
Figure 41.	Elementary routing element transient validation simulation.....	46
Figure 42.	TSVs-based die stacking cross-sectional view.....	48
Figure 43.	Cross-sectional view of 3-D parallel integration (left) and 3-D monolithic sequential integration (right) .....	49
Figure 44.	Cross-sectional view of 3-D monolithic steps – a) optimized bottom FDSOI process b) high quality top film deposition c) low temperature top FDSOI process.....	50
Figure 45.	Contacts between metal layer and top/bottom layers [97] .....	50
Figure 46.	3-D memory for reconfigurable logic – a) schematic b) layout on top levels ..	51
Figure 47.	Two layer monolithic 3-D based 1bit Look-Up Table – layout and schematic .....	52
Figure 48.	Cross point schematic.....	52
Figure 49.	Two layer Configurable cross points (Configuration node based) – layout and schematic .....	53
Figure 50.	Two layer Configurable cross points (SRAM cell based) – layout and schematic .....	53
Figure 51.	Cross sectional schematic showing a BEOL FET and standard CMOS FET co-integration .....	57
Figure 52.	Boolean logic with vertical FETs.....	58
Figure 53.	NOT function implemented in a 3-D standard cell .....	59
Figure 54.	NAND function implemented in a 3-D standard cell – a) layout on two gate levels b) layout with only one gate per nanowire.....	59
Figure 55.	5T SRAM cell implemented in a 3-D standard cell – a) schematic b) cross view .....	60
Figure 56.	Cross point implemented using 3-D smart vias and 3-D SRAMS – a) schematic b) top view.....	60
Figure 57.	TCAD cylindrical representation of the vertical NWFET structure (scales are in micron) .....	61
Figure 58.	TCAD circuit characterization methodology flow chart.....	61
Figure 59.	I-V characteristics for vertical nanowire FET.....	62
Figure 60.	FPGA benchmarking flow diagram .....	69
Figure 61.	PCM-based FPGA organization.....	70
Figure 62.	Delay estimation for FPGAs synthesized with GST-PCM- and SRAM-based Switchboxes .....	71

Figure 63.	Delay estimation for FPGAs synthesized with PCM- and SRAM-based LUTs and MUXs.....	72
Figure 64.	Variation of the FPGA critical path delay with the PCM on-resistance (delay averaged over the whole benchmark set).....	72
Figure 65.	Monolithically 3-D integrated FPGA organization.....	73
Figure 66.	Area estimation for FPGAs synthesized with standard bulk circuits and monolithic 3-D integrated FDSOI circuits.....	74
Figure 67.	Delay estimation for FPGAs synthesized with standard bulk circuits and monolithic 3-D integrated FDSOI circuits.....	74
Figure 68.	Vertical-NWFET-based FPGA overall organization.....	76
Figure 69.	Delay estimation for 2-D and 3-D FPGAs.....	76
Figure 70.	Illustration of the 3-D and 2-D cell co-integration.....	77
Figure 71.	a) Schematics of a CNFET with a back-gate configuration [142] b) Measured characteristics of a typical CNTFET (CNT $\phi=1.4\text{nm}$ , Ti contacts and gate oxide 10-nm of SiO <sub>2</sub> ) [142] c) Band diagrams of the structure [142].....	85
Figure 72.	a) SEM image (top) and schematic cross-sectional diagram of a DG-CNTFET [142] b) Measured characteristics of a DG-CNTFET (Back-gates are polarized to shown p- and n-type unipolar behaviors) [142] c) Band diagrams of the structure (Left $V_{\text{gs-Si}} < 0$ , and right $V_{\text{gs-Si}} > 0$ ) [142].....	85
Figure 73.	DG-CNFET device schematic using the proposed process-flow and showing the source (S), drain (D), front- (FG) and back-gate (BG) contacts....	86
Figure 74.	DG-CNTFET device layout.....	86
Figure 75.	Specified I-V curve of a tuned DG-CNTFET.....	87
Figure 76.	Simulated I-V curve of a tuned DG-CNTFET [IMS].....	88
Figure 77.	a) Transistor level schematic and b) configuration table for CNT reconfigurable cell [122].....	89
Figure 78.	Reconfigurable logic gate waveforms in NOR configuration [122].....	89
Figure 79.	Delay analysis of a CMOS MUX in 65-nm and 22-nm.....	91
Figure 80.	Delay analysis of DG-CNTFET Reconfigurable Cell.....	92
Figure 81.	a) Generalized scheme for CMOS dynamic logic b) Generalized scheme for double-gate-based cells.....	93
Figure 82.	a) Buffer cell schematic b) associated waveform and c) XOR cell.....	94
Figure 83.	Silicon electronics evolution from bulk to nanowires.....	95
Figure 84.	Manufacturing process to build a nanowire crossbar. <i>a)</i> Deposition of silicon nanowire after growth, oxidation and Langmuir Blodgett alignment. <i>b)</i> Etching of the oxide shell. <i>c)</i> Salicidation of the nanowire regions which will not be transistor channels. <i>d)</i> Deposition, alignment and etching of the oxide shell of the second layer. <i>e)</i> Salicidation as step <i>c)</i> . <i>f)</i> metallization to contact the nanowires around the crossbar.....	96
Figure 85.	NOT function realized on <i>n</i> -type nanowires using dynamic logic. <i>a)</i> Pseudo-physical view <i>b)</i> Schematic view <i>c)</i> Associated waveform.....	97
Figure 86.	Representation of CB-NWFET dynamic reconfigurable logic cell.....	98
Figure 87.	Metal/via layout of the dynamic logic cell.....	98
Figure 88.	Delay analysis of Nanowire Reconfigurable Cell.....	99
Figure 89.	Manufacturing process to build a FDSOI crossbar: a) grating patterning and active regions doping. b) Passive regions definition. c) Gate deposit. d) Passive regions finalization and salicidation. e) Metallization to contact passive regions.....	102
Figure 90.	Crossbar inverter structure (a) equivalent circuit and (b) layout.....	102

Figure 91.	Stand-alone 4:1 MUX crossbar implementation (separate doping regions) and equivalent schematic .....	103
Figure 92.	FPGA-suited 4:1 MUX crossbar implementation (alternating doping regions) and equivalent schematic .....	104
Figure 93.	Modeling of parasitic devices .....	104
Figure 94.	Delay analysis of lithographic crossbar in a parasitic-free context.....	105
Figure 95.	Influence of P-type horizontal wire size ( $F_{HP}$ ) on MUX propagation delay (Output of the inverter) ( $F_{HN} = F_V = F_I = F = 60 \text{ nm}$ ) .....	106
Figure 96.	Influence of N-type horizontal wire size ( $F_{HN}$ ) on MUX propagation delay (Output of the inverter) ( $F_{HP} = F_V = F_I = F = 60 \text{ nm}$ ).....	106
Figure 97.	Influence of insulating wire size ( $F_I$ ) on MUX propagation delay ( $F_{HP} = F_{HN} = F_V = F = 60 \text{ nm}$ ).....	107
Figure 98.	Influence of gate oxide thickness ( $T_{OX}$ ) on MUX propagation delay ( $F_{HP} = F_{HN} = F_V = F_I = F = 60 \text{ nm}$ ).....	107
Figure 99.	Impact of scaling ( $F_{HP} = F_{HN} = F_V = F_I = F$ ) .....	108
Figure 100.	Delay analysis of lithographic crossbar with parasites .....	109
Figure 101.	Conventional hierarchical template (a), FPGA model (b) and modified levels (c) .....	115
Figure 102.	MCluster approach for reconfigurable architectures (MCluster_4_4).....	116
Figure 103.	Matrix of 16 reconfigurable gates with fixed interconnect topology (a)Banyan, b)Baseline, c)Flip, d)Modified Omega) .....	117
Figure 104.	MCluster-based CLB proposal.....	118
Figure 105.	Final FPGA layer organization .....	118
Figure 106.	Disruptive technology compatible benchmarking flow diagram .....	120
Figure 107.	Illustration of possible scenarios compatible with the proposed toolflow .....	120
Figure 108.	MPack model flow .....	121
Figure 109.	a) Banyan_4_4 topological arrangement (logic cells are in yellow, virtual input nodes are in green, virtual output nodes are in red) and b) associated cross-connectivity matrices .....	121
Figure 110.	Input level correction illustration .....	122
Figure 111.	Feedback correction illustration .....	122
Figure 112.	Jump correction illustration.....	123
Figure 113.	Multiple buffering/inverting path simplification illustration – a) incoming path – b) outgoing path.....	123
Figure 114.	Illustration of multiple output correction .....	124
Figure 115.	Output level correction illustration .....	124
Figure 116.	MPack mapping algorithm (pseudo-code) .....	125
Figure 117.	Function graph to map onto the Banyan_4_4 interconnect topology and its associated adjacency matrix .....	125
Figure 118.	Function graph after correction step.....	126
Figure 119.	MCluster after function packing .....	126
Figure 120.	MPack' clustering algorithm (pseudo-code).....	127
Figure 121.	Performance evaluation method for fixed interconnection topologies .....	128
Figure 122.	Programmability success rates for Banyan, Modified Omega, Flip and Baseline interconnect topologies within 4-deep 4-wide matrices.....	129
Figure 123.	Comparison of programmability success rates for Banyan and Modified Omega interconnect topologies within 4-deep 4-wide matrices in the case of faulty links and faulty cells .....	130

Figure 124.	Comparison of average interconnect length for Banyan, Modified Omega, Flip and Baseline interconnect topologies within 4-deep 4-wide matrices.....	130
Figure 125.	Area estimation for MCluster-based FPGAs with various granularities.....	132
Figure 126.	Area estimation for MCluster_3_3-based FPGAs with different latch depopulation scenarios .....	132
Figure 127.	Area estimation for 3 by 3 MCluster-based and CMOS-based FPGAs.....	133
Figure 128.	Critical routing delay repartition for 2 by 2, 3 by 3 MCluster-based and CMOS-based FPGAs .....	133
Figure 129:	Modified-Omega based MCluster with early internal outputs and cell depopulation .....	134
Figure 130.	Illustration of two MClusters_4_4 physical implementation on parallel carbon nanotubes layer.....	135
Figure 131.	MCluster organization for speculative CNLB .....	136
Figure 132.	Speculative pipelined organization of CNLB .....	136
Figure 133.	Double-stage MCluster organization with two-level hierarchic interconnection strategy .....	137

# *List of Tables*

---

Table I.	Linear scaling rules impact on device parameters [5].....	2
Table II.	State-of-Art architecture global metrics comparisons.....	28
Table III.	Specification estimation wrt. memory distribution through an FPGA (extracted from Xilinx Virtex6 architecture [20]).....	34
Table IV.	Detailed technology performance evaluation.....	46
Table V.	Technology performance evaluation (2x2 switchbox).....	47
Table VI.	Density survey of TSV technologies [96].....	49
Table VII.	Evaluation of SRAM Shifter performance.....	54
Table VIII.	Evaluation of Look-Up Table performance.....	55
Table IX.	Evaluation of Cross point performance.....	56
Table X.	Simulation results summary.....	62
Table XI.	Technology comparison (4-input LUT test case).....	64
Table XII.	Architectural evaluation summary.....	78
Table XIII.	Global evaluation of the DG-CNFET cell performances.....	91
Table XIV.	Power consumption vs. function (22-nm node – 4GHz).....	92
Table XV.	Dynamic DGCNFET-based cell transistor requirements.....	94
Table XVI.	Global evaluation of the CB-NWFET MUX performances.....	99
Table XVII.	Detailed Power Consumption figures for sublithographic cell (22-nm – 4GHz).....	100
Table XVIII.	Evaluation of a Parasitic-Free 4-to-1 MUX in lithographic crossbar technology.....	105
Table XIX.	Detailed Power Consumption figures for the sublithographic cell in a parasitic-free context (65-nm – 650MHz).....	105
Table XX.	Evaluation of scaling with parasitics.....	108
Table XXI.	Evaluation of lithographic crossbar-based structure with parasitics.....	108
Table XXII.	Detailed Power Consumption figures for sublithographic cell (65-nm – 650MHz).....	109
Table XXIII.	Global comparison of analyzed disruptive technology cells.....	110
Table XXIV.	Analysis summary.....	131
Table XXV.	Summary of global gain figures for the proposed architectures with respect to a baseline CMOS 65-nm FPGA architectures.....	139

# CHAPTER 1 *Introduction*

---

Our post-industrial way of life can be defined by an ever-growing need for mass consumer products. Every day, the industry generates novel innovations that continually boost the market. Planes are increasingly automated, cars assist the drivers in a various number of situations for security or leisure (e.g. speed control, parking assistance, etc.) and cell phones are able to communicate over the internet and to use several different networks. Video game consoles are controlled by handset motions or even player movements. Such evolutions mean that the systems have to be more intelligent, have simpler interface and be more communicative. The enabler behind all these novelties is the generalized use of embedded electronic systems, which have long become an integral and pervasive part of the human society. As the demand in terms of entertainment, data volume and connectivity is in constant growth, the system performance needs to improve more each day. The market is boosted by application requirements and remains in constant progress and evolution.

Electronic systems are complex systems, which aim to address computation in a specific target application. They are built around a physical hardware core, which can be programmed by loading and adding software modules. The hardware is the basic key for computation, in the sense that it is used to manage information at the physical level, i.e. in conventional systems and at the most elementary level, the electronic charge. The hardware is closely related to design and manufacturing techniques. The design defines the organization of the components in the system (i.e. the architecture). The manufacturing techniques define the physical components. This makes the related field highly multi disciplinary.

The basic unit of electronic systems is the transistor. Invented in 1947 by William Shockley, John Bardeen and Walter Brattain, the transistor is able to work as a switch at the smallest scale [1]. While the first demonstrated transistor was a bipolar transistor made of Germanium, the technology moved quickly to Silicon-based transistor and to the demonstration of the first *Metal-Oxide-Semiconductor Field Effect Transistor* (MOSFET) in 1959 [2]. The Silicon MOSFET transistor is still in use in the most advanced modern integrated circuits. By scaling the device and using a large number of transistors on a same substrate, it became possible to create integrated circuits of higher complexity and functionality. The era of the semiconductor industry has been defined by scaling, and consequent improvement for almost half a century: the industry has followed an exponential evolution from the first realization of integrated circuits to the current *Very Large Scale Integration* (VLSI<sup>1</sup>) with billions of transistors.

## 1.1 The Scaling Era

The way in which the semiconductor industry has evolved is unique with the scope of industrial history by its impact on society and rate of progress. Since the creation of the transistor, the semiconductor industry has grown exponentially, a trend which has been enshrined in Moore's Law [3] (Figure 1). Starting in 1965 and based on empirical

---

<sup>1</sup> This term was first coined at the 100ktransistor mark. Other terms were later proposed (e.g. ULSI) to reflect the continually increasing complexity of integrated circuits, but VLSI remains the term of choice within the community.

observation, this law predicts an exponential growth (a doubling) of the number of transistors per die for each new technology generation. Even the growth rate (or time between technology generations) has not been constant though during recent decades, it still has a factor of 2 every 2 years. As illustrated by figure 1, it is worth noticing that the transistor count of processors has followed a similar trend for the last four decades.

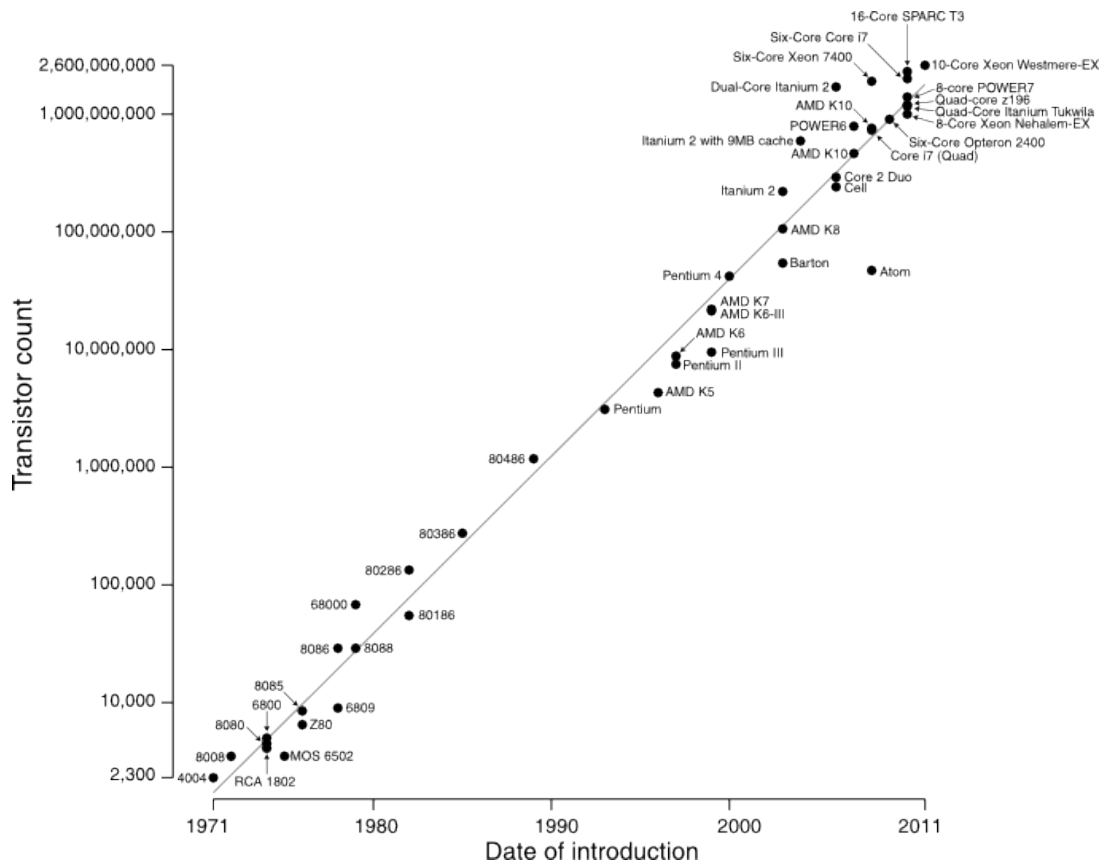


Figure 1. Processor transistor count and Moore’s Law since 1971 [4]

The increase in the number of transistors is an obvious consequence of the scaling. Scaling leads to the reduction of the dimensions of elementary drawn transistors, making it is thus possible to integrate more devices within the same area. Further scaling brings many more additional advantages than the sole matter of density. Since the majority of technological parameters are scaled according to a given scaling strategy. The impact of scaling on the main device parameters for constant electric field (the dominant scaling strategy for deep submicron technology generations) is presented in Table I. It is important to observe that scaling leads to a reduction of device power by a factor of  $\alpha^2$  and a speed up in the intrinsic delay by a factor of  $\alpha$ .

Table I. Linear scaling rules impact on device parameters [5]

Parameters	Scaling factor
Transistor length and width ( $L, W$ )	$1/\alpha$
Junction depth ( $x_j$ )	$1/\alpha$
Oxide thickness ( $t_{ox}$ )	$1/\alpha$
Doping concentration ( $N_d, N_a$ )	$\alpha$
Supply voltage ( $V_D$ )	$1/\alpha$
Drive current ( $I_D$ )	$1/\alpha$
Electric field ( $E$ )	1
Capacitance ( $\epsilon \cdot A/t_{ox}$ )	$1/\alpha$
Delay time ( $\tau=C \cdot V_D/I_D$ )	$1/\alpha$
Power dissipation ( $\sim V_D \cdot I_D$ )	$1/\alpha^2$
Device density ( $\sim 1/A$ )	$\alpha^2$

This fast evolution is unique. As an illustration, if other industries had followed a comparable rate of change, a one-way ticket between Paris and New York may would now cost only 0.01€, while the total flight time would have been reduced down to 0.25 s. Similarly, a car scaled by the same proportions would only cost 20€, reach a speed of  $3000 \text{ km.s}^{-1}$ , would only weight 10 mg and would only consume 1l of fuel per 100 000 km.

First only seen as a trend based on empirical observation, the Moore's Law has become much more than just a speculation. It has been regarded (and is still considered) as a real specification for growth, enabling the definition of objectives for research and development in the semiconductor industry. This has given rise in recent years to the *International Technology Roadmap for Semiconductors* (ITRS), compiled by a consortium of academic and industrial leaders in the fields of semiconductor, and whose goal is to survey the trends in semiconductor technology and predict its future evolution up to fifteen years ahead. As a consequence, it defines the objectives for future technology nodes and identifies roadblocks to overcome. Figure 2 shows the current tendency of scaling of the main technology node parameters (metal layer half pitch and printed gate length). It is essential to highlight that Moore's Law is still expected to be maintained for the next two decades.

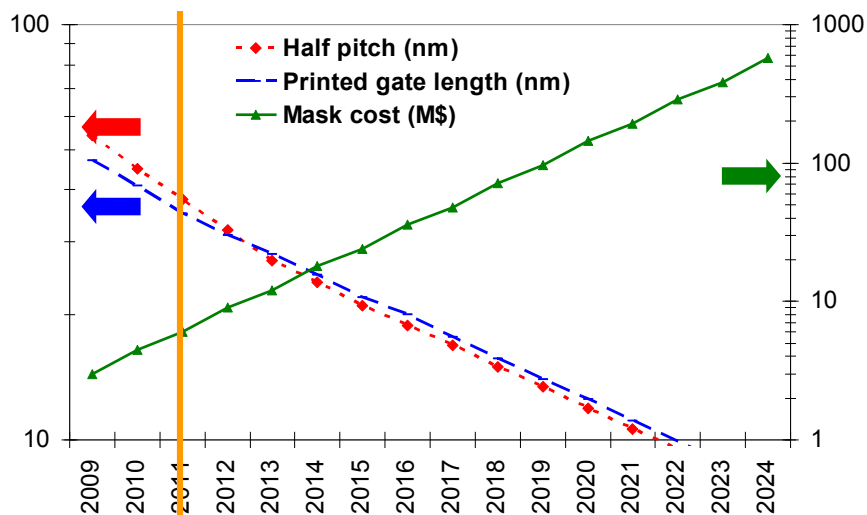


Figure 2. Projected evolution of dimensions and mask cost [6]

Nevertheless, it is also worth pointing out that ITRS recognizes the existence of physical and economical limits to this growth.

## 1.2 New Challenges

The ITRS has indicated significant future hurdles that may impede the original scaling law. Limitations are found at various levels, from fundamental device characteristics to advanced system design methodologies.

### 1.2.1 Process Integration and Devices

The process integration and the devices will suffer from leakage and quantum effects, as well as from intrinsic fabrication hurdles (for example lithography steps).

- **Electrostatic Channel Control**

With channel scaling, *Short Channel Effects* (SCE) are becoming increasingly dominant. The short channel effect reflects the lowering of the threshold voltage with a decreasing channel length. This effect is due to a two-dimensional distribution of the surface potential in the channel. Furthermore, in short channel devices, the *Drain-Induced Barrier Lowering* (DIBL) effect could be added. The DIBL consists of making the threshold voltage dependent on the drain bias. Both of these effects lower the threshold voltage and make the devices more vulnerable to variability.



- **Leakage**

It is worth noticing that when the lateral dimensions of a transistor are scaled, the oxide thickness is also scaled. This leads to an exponential growth of the tunnel current through the gate oxide, thereby an increase in the gate leakage.

Moreover, new physical phenomena appear at scaled dimensions and make a significant to leakage current. Indeed, high values of channel doping cause *Band-To-Band* (BTB) tunneling and *Gate-Induced Drain Leakage* (GIDL). These effects lead to higher power consumption [8, 9]

- **Lithography**

Microelectronics fabrication is based on lithography. Tools using *Deep UltraViolet* (DUV) light with wavelengths of 248 and 193nm allow minimum feature sizes down to 50nm. However, scaling below this will require several innovations in terms of design and equipments. The mask design must contain a high degree of regularity, in order to minimize the fabrication variability and impacts of lithography steps, while *Optical Proximity Correction* (OPC) is required to compensate inherent diffraction. Finally, new equipments are envisaged to reduce the wavelength down to *Extremely Ultra-Violet* (EUV).

### 1.2.2 Systems and Tools

On the design perspective, several obstacles must also be overcome.

- **System design**

Many device parameters do not scale exactly according to the scaling theory [10], due to device intrinsic hurdles. As a result, many device parameters which used to be fixed by the technology node are becoming design parameters. This is especially true for the supply and the threshold voltage. The designer can optimize these parameters in order to obtain the best trade-off in terms of area, power and delay.

Furthermore, other circuit parameters are not following the scaling rules at all. Indeed, the latency of on-chip wires does not follow the same trend as front-end devices. Global signals impact the performance metrics significantly, and therefore are very challenging. In order to both improve the energy-per-bit and to increase the bandwidth density, new on-chip connection paradigms are explored such as *Network-On-Chip* (NOC) [11, 12].

Another challenging issue is the distribution of clock signals. Conventional circuits are based on a single clock that is distributed through out the whole circuit. With the increase of the relative length of clock lines (compared to the technology node), the task of distributing the clock signal without skew (i.e. simultaneously) and without jitter (i.e. periodically) everywhere on the chip is extremely difficult.

- **Design tools and methodologies**

The possibilities enabled by current technologies are enormous. Nevertheless, the number of potential transistors that are available to chip design is much larger than the number of transistors (or complexity) that design teams can efficiently use. This difference between the technological possibilities and the design capabilities is called the Design Gap. In order to improve the productivity, design tools are developed to automate the design process, as much as possible. In this way, hardware description languages, high-level synthesis and hardware/software co-integration are a small illustration of the main design tasks that need to be addressed deeply to reach the full potential of a given technology.

## 1.3 Architectural and Technological Opportunities

The semiconductor industry has to face many issues to pursue the Moore's Law, which is, anyway, expected to end, when the MOSFET transistor reaches its physical limits. In this thesis, we will assess a twofold opportunity. This opportunity is based on reconfigurable

architectures that can be improved by disruptive technologies, in order to yield solutions for future architectures.

### 1.3.1 Reconfigurable Architectures

In order to keep following the Moore's law and to achieve the computing capacities necessary for future software applications, it is today widely recognized that *Systems-On-Chip* (SoC) will move initially towards *Multi-Processor Systems-On-Chip* (MPSoC), then towards *reconfigurable platforms*. These systems will be used in the majority of solutions and in particular for high-performance computing (analysis and modeling of complex phenomena, advanced human-machine interaction) and for low-power mobile systems (sensor networks ...).

The reconfigurable approach to computing systems offers several advantages. It allows volume manufacturing and thus constitutes a solution to the projected evolution of mask costs. The mask costs are expected to be above \$10M in 2011 and above \$100M in 2018, as depicted in figure 2. In all probability, “full-custom” dies will disappear or migrate towards the development of *Multi Project Wafer* (MPW) lines for very high performance applications; the appearance of “stacked” circuit based on 3-D integration will aim at heterogeneous applications; and the reconfigurable approach will allow mastering volume applications [6]. Such systems can cover a broad range of applications, and their performance levels exceed very clearly those of programmable systems in terms of computing speed, while requiring only one set of masks. Moreover, the natural association of these architectures with fault-tolerant design techniques enables to build robust architectures in the context of increasingly unreliable elementary nanometric MOS devices. Nevertheless, the various types of reconfigurable circuits (FPGA, coarse-grain reconfigurable systems) are at a disadvantage (compared to “full-custom” solutions) in terms of performance and device count necessary to fulfill a specific function (Figure 3).

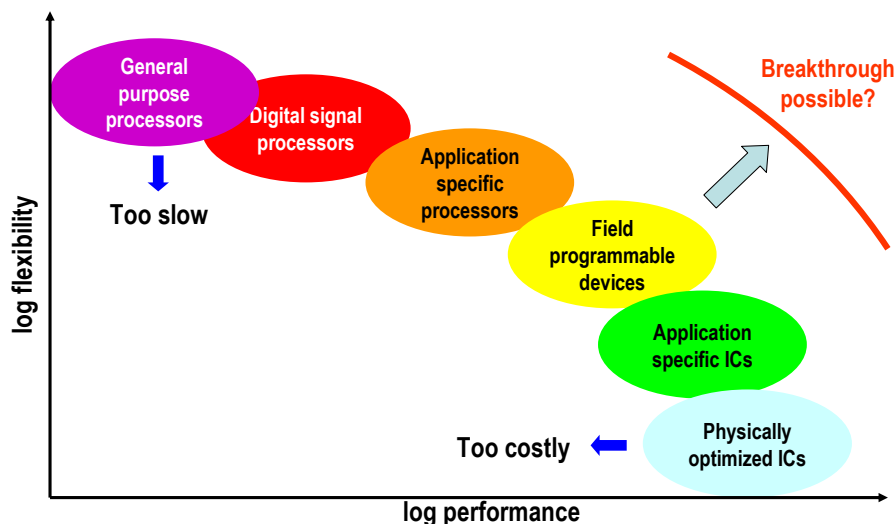


Figure 3. Spectrum of implementation solutions for data processing algorithms

### 1.3.2 Technological Evolutions

In this context, the emergence of new devices offers the opportunity to provide novel building blocks, to elaborate non-conventional techniques for reconfigurable design and consequently to reconsider the paradigms of architecture design. The ITRS *Emerging Research Device* (ERD) and *Emerging Research Materials* (ERM) chapters propose emerging technology fields for prospective research [7]. Two different directions are envisaged:

- i) the extension of MOSFET device to other shapes and materials, and
- ii) the use of other technologies and state variable for computing.

### 1.3.2.1 MOS Extensions

The straightforward evolution for MOSFET processes is to work on the channel materials and structure. The main idea consists of replacing the gate channel by new materials with high carrier mobility. Three different possibilities are open for consideration. The conventional approach consists of working on process engineering to improve the channel performance, with for example strained silicon or silicon on insulator wafers. The ultimate goal of such channel processing is to achieve the use of 1-D structures as enhanced channels. Thus, silicon NanoWires could be seen as promising candidates for this evolution. Finally, the use of unconventional materials could be seen as a replacement material for silicon. Carbon electronics is a good prospective approach, with promising electrical properties, coupled with the confinement which is achievable by the structure. Furthermore, it is important to highlight that carbon nanotubes and graphene also exhibit respectively 1-D and 2-D structures that have better electrical properties.

### 1.3.2.2 Emerging Devices and Memories

Rather than considering “only” an improvement of the conventional electric charge based MOSFET, the ITRS ERD and ERM chapters also suggest the possibility of using new state variables for computation as well as for the information storage. As show in figure 4, even if the standard way to carry out computation is based on Silicon MOSFET in Von Neumann multicore system architectures, many other ways could be explored. Indeed, unconventional state variables might be used, such as molecular state and phase state. It is also worth mentioning that unconventional solutions might be done at every hierarchy level, and that all these choice and/or combination of choice might give rise to new computation paradigms.

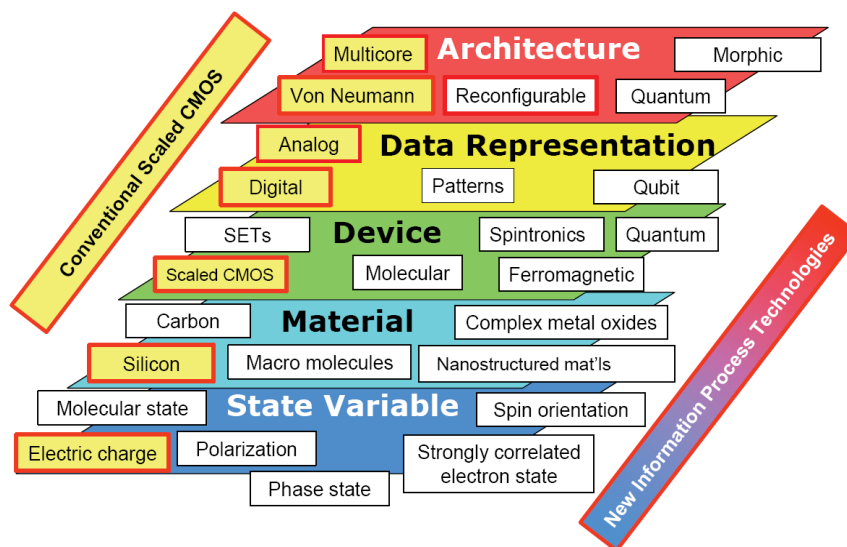


Figure 4. Electronic systems hierarchical organization and opportunities [7]

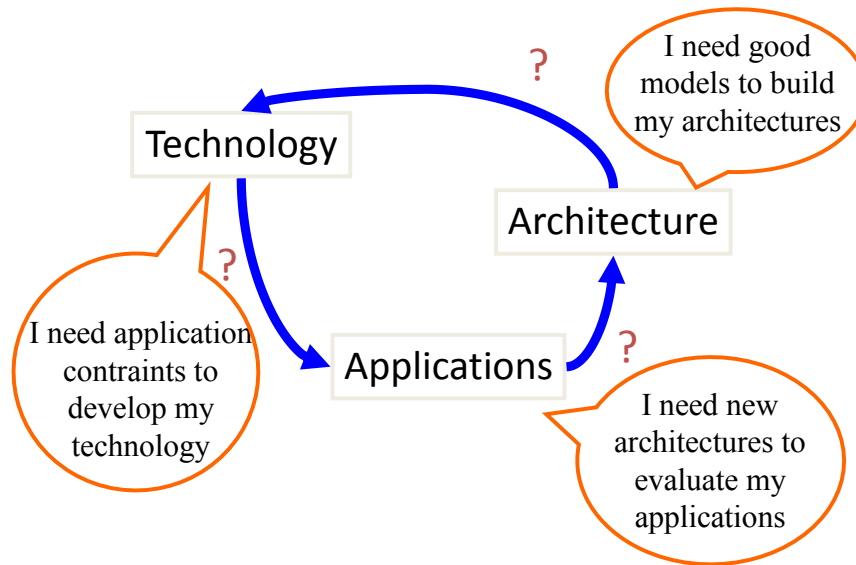
## 1.4 The Nanoscale Chicken and Egg Dilemma

While the use of emerging devices appears to be the correct approach to go beyond the limitations of CMOS technology, high cost and lack of maturity of fabrication processes render the task of finding a good candidate device difficult.

The evaluation of a new technology is at the interface of different domains, and generally leads to a Chicken and Egg Dilemma, as depicted in figure 5. A technology is generally optimized to a given class of application. For example, we could oppose low power memories and high performances applications. These applications have different constraints regarding the speed, the power consumption and many other parameters, which lead to different CMOS optimization strategies, such as HVT<sup>2</sup> and LVT<sup>3</sup> transistors. Nevertheless, it is difficult to

<sup>2</sup> High Threshold Voltage ( $V_t$ )

formulate application constraints without having an architectural model. This model is used to bench algorithms onto the architecture. The definition of an architecture requires information from the technology, in order to create an optimized system.



**Figure 5. The Chicken and Egg Dilemma of the disruptive technology**

We face several problematic, and it thus appears necessary to answer the following questions:

- i) How can we find a suitable technology for future applications?
- ii) Knowing that it is impossible to carry out technological developments for all candidate processes, as it would generate excessive costs. This indicates that we should employ a fast evaluation methodology for emerging technologies: But in this case, which architecture could we use? How will we evaluate the performance of a disruptive technology and which tools should we use?
- iii) How can the technology characteristics taken into account within the architecture? What should be the specific functions to address be? How can the computing blocks be improved? How can the peripheral circuitries be improved? What is the best arrangement of the logic blocks for efficient tuning of the architecture? In the following, we will present the global evaluation methodology that is used in this work and we will especially describe the contributions.

## 1.5 Work Methodology

The fast evaluation methodology is the core of this thesis and part of the ANR research project “Nanograin” to which this thesis work mainly contributes.

The “Nanograin” fundamental research project aims to devise and evaluate new design paradigms opened up by a family of novel reconfigurable cells based on *Double-Gate Carbon Nanotube Transistors* (DG-CNFET). Each reconfigurable cell is capable of realizing several logic functions (with an invariant hardware structure) according to the voltages applied to the back gates, and thus constitutes an original “ultra-fine grain” or “nano-grain” reconfigurable approach. The work outlined in this project will comprise: compact modeling of the double-gate CNFET device for accurate validation of developed cells and architectures; circuit-level research to

---

<sup>3</sup> Low Threshold Voltage ( $V_t$ )

- i) refine the structure and create a family of variants adapted to specific applications,
- ii) tackle non-volatile configuration memory aspects,
- iii) research inter-cell interconnect strategies; the development of a programmable architecture based on the various cells (“sea of gates” type), capable of realizing complex functions, with a special focus on the fault-tolerance aspect; analysis of a programming model for the architecture, and virtual demonstration of a system-level application running on the architecture; analysis of the characteristics of the architecture, particularly concerning the fault tolerance aspect, in the context of reduced device reliability.

The project covers all the parts of the formulated problem. In this work, we propose to pursue the idea and extend the methodology to many different technologies. The proposed methodology is presented in figure 6. A large set of available technologies will be evaluated under a generic architectural template. The chosen template will be a generalized reconfigurable architecture. Depending on the process and technology maturity, different models will be used and injected into the architectural template description. These models range from simple behavioral level to advanced physical TCAD extractions and real-circuit measurement. The envisaged applications will be standard testbenches sets in order to investigate a wide range of application domains and attempt to find the best one. A fast architectural optimization loop exists in the methodology. It allows identification of the best architectural trade-off for a given technology. Thus, it ensures that bad architectural sizing, which may lead to an erroneous estimation of technology performance, is avoided. Finally, the architectural evaluation is intended to feed back to the technology side, in order to find which parameters improve the targeted application.

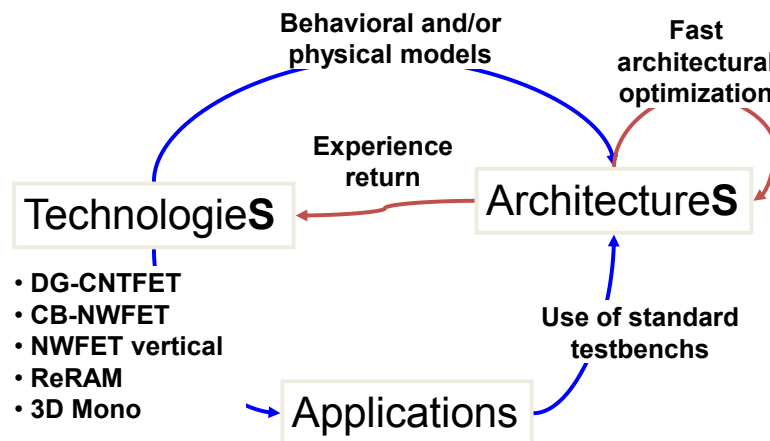


Figure 6. Template based Fast Evaluation methodology

## 1.6 Research Contributions

In accordance with the global goal of this research work, several questions have been assessed in different fields, in order to address globally the issues of fast evaluation and to propose realistic architectural solutions based on elementary technologies. The contributions have been made at several levels from the elementary circuits to the architecture, as well as benchmarking tools and methodologies.

### 1.6.1 Methodology and Tools

On the methodology and tools side, this thesis contributes to a state-of-the-art benchmarking tools by proposing a complete tool flow suited to the design exploration of emerging devices based reconfigurable logic (Chapter 6). Indeed, in order to evaluate and compare in a fair way the prospective circuits and architectures, it is mandatory to have a complete tool, able to

instantiate a set of standard circuits on the proposed technology and architecture. The proposed toolflow is based on extensions to the conventional tools from the reconfigurable community. The extensions allow a wide variety of structures and architectures for architectural exploration such as the management of logic cells instead of the conventional FPGA look-up table, the organization of logic cells into regular matrices and the use of fixed and incomplete interconnection topologies. To do so, a specific packing tool, called MPack, has been developed. It will be presented and its internal structure will be described in detail.

### **1.6.2 Memories and Routing Resources for FPGA**

As we will discover in the next chapter, the FPGA architecture suffers from its amount of resources required for programmability. In particular, the versatility of the architecture is handled by several pass transistors connected to programming memories. These circuits are distributed through out the whole circuit, and are known to be area and power hungry. In this thesis, we will assess three different technologies, in order to improve the memory and the routing aspects of reconfigurable systems (Chapter 3).

Resistive memories are promising devices which exhibit programmable non-volatile bi stable resistive states. This technology permits a passive programmable device to be embedded into the metal layers. Thus, we will propose a configuration memory node. This node will use a voltage divider to store a logic level intrinsically. Then, the impact of the structure will be studied in comparison to flash memories. Using the favorable on-resistance property of resistive memory, we also propose to integrate them directly into the data path. This is obviously of high interest to create non-volatile and reprogrammable switches with low resistance in the on-state and low area impact. In addition, a resistive memory based switchbox will be proposed and compared to an equivalent state-of-the-art non volatile implementation.

Instead of using only a passive device in the back-end layer, we will assess the opportunity to integrate active devices in a monolithic 3-D arrangement. Monolithic 3-D processes aims to stack several layers of active Fully-Depleted Silicon-On-Insulator. Since layers are processed sequentially, it is possible to obtain a high alignment quality, as well as a high via density between the two layers. We will propose an FPGA elementary block implementation using this 3-D technology. Effectively, we will split the memory part from the data paths of the FPGA blocks. This will allow distinct optimization of the sizing and the processes.

Finally, we will push the 3-D concept to the limits by proposing a “true” 3-D implementation of routing resources. Vertical NanoWire Field Effect Transistors have been demonstrated as a possible technology to create vertical transistors between metal lines. These transistors are large, which enable them to reach promising performance levels. We then propose in building a “smart” back-end methodology with complex configurable vias, as well as configuration memories for vias and signal buffers. We will propose a performance evaluation methodology based on TCAD simulations to open the way towards prediction of the most advanced technologies.

### **1.6.3 Logic Blocks for FPGA**

While improving memories and routing resources appears to be the most convenient way to improve the FPGA architecture, it is worth pointing out that the structure uses only a very small part for the computation real-estate. The most part of the circuit is used for “peripheral circuitry”. Thus, we could expect to find a more effective structure. In this way, it is interesting to work on logic blocks and see how emerging devices could lead to new paradigms for computing at fine-grain (Chapter 5). Two main device families will be studied.

Firstly, we will study the Double-Gate Carbon Nanotubes Field Effect Transistor technology. These devices might be configured between *n*- and *p*-type conduction by changing their back-gate voltage. This allows achieving an in-field reconfigurability at the device level. We called this kind of device an “enhanced-functionality” device. This technology has been used to provide a very compact logic cell, able to perform computation at an ultra-fine level of

granularity. We will refine the technological assumptions and generalize the proposal of the use of a configuration back-gate to a family of standard cells.

Secondly, we will focus on “density-increased” devices. Emerging technologies have opened the way towards the use of active devices in a high density crossbar fashion. We thus propose to use sublithographic crossbars of nanowires to realize a compact ultra-fine grain logic cell. We will then study the fabrication assumptions and study a potential structure for the device. Nevertheless, we must consider the current technological hurdles of sublithographic processes. Hence, we will also present a crossbar process which involves an FDSOI-based lithographic process flow. In the end, we will explain how circuits and layouts could be implemented based on this circuit.

Finally, it is important to note that all the expected blocks will lead to a significant reduction of granularity (i.e. towards finer levels of granularity) in the logic blocks.

#### **1.6.4 Nanoscale Architectures**

Ultra-fine grain logic provides a large reduction of the logic area real-estate. Nevertheless, the scaling of wires and peripheral resources is not so fast. Hence, it is required to envisage new organizations to assess an efficient architecture. We will propose regular arrangements of logic cells into matrices (Chapter 6). It should be noted that such arrangement is also motivated by regular patterning in lithography steps or in self-assembly. While the use of a standard FPGA interconnect scheme could lead to an unexpected overload in terms of resources, we propose a fixed and incomplete interconnect pattern. This pattern will ensure the maximum shuffling, reachable by scarce resource availability. This regular pattern is then used to realize logic blocks for reconfigurable architectures. A performance evaluation of the structure will be performed and it will be compared to standard FPGA. More disruptive architectures will then be discussed.

### **1.7 Organization of the Thesis**

The thesis is organized into five main chapters, not including the introduction and the concluding chapter.

In Chapter 2, we proposed to describe the background and motivation of the work with a state-of-the-art Field Programmable Gate Arrays architecture, tendencies and issues are described. Then, we will summarize existing computing architectures, based on emerging devices. We will particularly focus on post-CMOS heterogeneous architectures, as they appear to be the most relevant to our field. Finally, we will formulate and generalize the description of reconfigurable circuits in terms of hierarchical levels. We will also define the architectural template that will be the baseline of our architectural evaluation.

In Chapter 3, we will focus on the emerging technologies for routing and memory implementations in reconfigurable circuits. Resistive Memory technology, Monolithic 3-D integration processes and Vertical transistor technology will be studied to assess the impact of moving the routing resources to the back-end layer. In this chapter, the performance of each proposed circuit will be evaluated electrically and compared to their equivalent counterparts in conventional CMOS technology.

In Chapter 4, the architectural evaluation of a traditional FPGA architecture, enhanced by the technologies presented above, will be proposed. Then, benchmarking results will be discussed with regards to the equivalent CMOS circuit.

In Chapter 5, disruptive logic blocks will be assessed. The use of functionality-increased and density-increased devices is envisaged. The logic blocks are then compared in terms of electrical metrics and their overall potential is discussed.

In Chapter 6, we use the logic blocks designed in the preceding stage to develop an ultra-fine grain FPGA architecture. Thus, the reconfigurable template is enriched with specific architectural proposals. These proposals are expected to deal with the granularity of ultra-fine

computing and with interconnection overload. A specific benchmarking toolflow is described. In particular, a packer tool will be presented, in order to tackle the matrices of cells with fixed interconnections. Finally, benchmarking results will be presented and discussed.

In Chapter 7, the thesis is concluded and possible future works are outlined.





## **CHAPTER 2**    *Background and Motivation*

---

### **Abstract**

In this chapter, we aim to present the background and the motivation of this thesis work. We will first give an overview of the Field Programmable Gate Array architecture, which is today the most widely used reconfigurable circuit. After describing its conventional structure, we will detail current trends in architectural organization. Then, we will survey the literature to see how disruptive technologies are used to propose drastic evolutions in the field. We will in particular show how dense nanowires can be used to build logic fabrics in a crossbar organization, and also how the use of carbon electronics allows the construction of interesting logic functionalities. Finally, we will try to formalize the various approaches into a hierarchical representation and compare it to the conventional structure. This representation will help to define the objectives of this work. We mainly intend to propose a digital reconfigurable circuit based on real-life disruptive technologies. This is an important point, since even if a potential technology opens the way towards new phenomena, it is fundamental to work closely with technologists and to keep in mind its feasibility from an industrial perspective. In this context, we will continuously try, in this thesis work, to take into account the technology requirements when designing a circuit.

As introduced previously, reconfigurable logic architectures are generic and highly versatile. This makes them an excellent compromise between costs, development time and performances. Suited for a wide range of application, they offer an intrinsic regularity compatible with the most advanced technological processes.

In the preceding, a fast and global evaluation methodology was described. This method requires the definition of a generic architectural template. The template will be defined for existing reconfigurable architecture.

In this state of the art chapter, we will give an overview of the reconfigurable field. Hence, we will first introduce the conventional Field Programmable Gate Array architectures. From the basic FPGA scheme, we will move to the most recent evolutions and discuss the structural issues.

Subsequently, from the emerging technologies perspective in the context of computation, we will survey the nanodevices-based architectures relevant to our field and give a global comparison between the different approaches.

Finally, we will formalize our research methodology, by describing the chosen global architectural template, as well as the position of our work regarding the existing literature.

## **2.1 Conventional Reconfigurable Architecture Overview**

Reconfigurable architectures are today leads by the Field Programmable Gate Array. We will begin this overview with the survey of the homogeneous FPGA architectural scheme and its related sizing. Then, we will see how the homogeneous architecture has been improved to increase the structural performances and to solve some architectural. Finally, we will discuss the current limits of FPGAs.

### **2.1.1 The FPGA Architecture**

In this section, we provide an overview of the standard FPGA scheme. We will start with some generalities of the structure. These generalities will deal with a short history/overview of reconfigurable circuits and will present the basements of the FPGA architecture. Then, we will detail the logic blocks architecture. Logic blocks architectures can be based on logic gates or LUTs. Interconnection structures will then be detailed. Finally, sizing of the architecture will be presented.

#### *2.1.1.1 Generalities*

*Field Programmable Gate Arrays* belongs to the family of reconfigurable logic circuits. Its structure is currently the most advanced of the family.

Historically, the reconfigurability has been based on programmable diode logic. Second generation architectures used in *Programmable Array Logic (PAL)/Programmable Logic Array (PLA)* architectures [27]. The PAL approach focused on the use of a reconfigurable full interconnectivity pattern for the implementation of the signal routing between macro-logic blocks. Hence, the PAL approach is intimately defined by its large routing array. It is important to note that in such a circuit, the logic is fixed and only the routing part is programmed.

The FPGA breaks this model by using both programmable logic and programmable routing structure. The logic is distributed through the routing structure in an island-style manner. This distribution helps in handling the routing congestions with an optimum number of resources.

The PAL routing architecture is a very simple but highly inefficient crossbar structure. Every output is directly connectable to every input. Connection is made through a programmable switch. The FPGA routing architecture provides a more efficient routing where each connection typically goes through several switches. In a programmable logic device, the logic is implemented using two-level AND-OR logic with wide input for the AND gates. In an FPGA, the logic is implemented using multiple levels of lower fan-in gates, which is often

much more compact than two-level implementations. An FPGA logic block could be as simple as a transistor or as complex as a *Digital Signal Processor* (DSP) block. It is typically capable of implementing many different combinational and sequential logic functions.

An FPGA structure is defined by fine-grain logic. This name comes from the granularity which is achieved by the logic blocks. Here, each logic blocks is supposed to realize a part of combinational or sequential logic operations but is not able to handle complex logic operation. The granularity is then in opposition to Massively-Parallel Processor Arrays or MultiProcessor System-On-Chip. Another particularity of the FPGA architecture is the hierarchical logic stratification. Indeed, a logic block is build from smaller logic blocks and a local reconfigurable interconnect. This interconnect is generally complete. This means that, like in a PAL, each signals (inputs and outputs) can be routed everywhere.

The routing architecture of an FPGA could be as simple as a nearest neighbor mesh [30] or as complex as the perfect shuffle used in multiprocessors [36]. More typically, an FPGA routing architecture incorporates wire segments of varying lengths which can be interconnected via electrically programmable switches. The choice of the number of wire segments incorporated affects the density achieved by an FPGA. If an inadequate number of segments is used, only a small fraction of the logic blocks can be utilized, resulting in poor FPGA density. Conversely the use of an excess number of segments that go unused also wastes area. The distribution of the lengths of the wire segments also greatly affects the density and performance achieved by an FPGA. For example, if all segments are chosen to be long, implementing local interconnections becomes too costly in area and delay. On the other hand if all segments are short, long interconnections are implemented using too many switches in series, resulting in unacceptably large delays.

The storage of the configuration is in charge of memories. These memories drive logic gates or pass-transistors in order to configure the logic or the routing. Different type of technologies could be used like *Static Random Access Memories* (SRAM), antifuse or flash. SRAM is actually the most standard technology, thanks to its CMOS technological homogeneity. In all cases, a programmable switch occupies larger area and exhibits much higher parasitic resistance and capacitance than a simple via. Additional area is also required for programming circuitry. As a result the density and performance achievable by today's FPGAs are an order of magnitude lower than that for ASIC manufactured in the same technology.

The complexity of FPGA has surpassed the point where manual design and programming are either desirable or feasible. Consequently, the utility of FPGA architecture is highly dependent on effective automated logic and layout synthesis tools to support it. A complex logic block may be under-utilized without an effective logic synthesis tool, and the overall utilization of an FPGA may be low without an effective placement and routing tool.

Placement and routing tools mainly depend on the architecture and are specific to manufacturers [28, 29]. Some tools aims to be used as an exploration tool for research exploration of algorithms as well as architectures. As an example, the *Verilog-To-Routing* (VTR) toolflow is dedicated to this purpose [182].

### 2.1.1.2 General Architectural organization

FPGAs are built of three fundamental components: logic blocks, I/O blocks and programmable routing, as shown in figure 7 and in figure 15. In figure 7, CLB stands for *Configurable Logic Block*, which is the combinational and sequential logic block. CB and SB stand respectively for *Connection Box* and *Switch Box*. These circuits form the global routing resources. In figure 7, I/O blocks have not been shown. They are found at the periphery of the circuit and are fed by the routing lines (Figure 15).

In an FPGA, a circuit is implemented by programming each of the logic blocks. Each of the I/O blocks is configured to act as either an input pad or an output pad. Each block implements a small portion of the logic required by the circuit. The programmable routing is configured to make all the necessary connections between logic blocks and from logic blocks to I/O blocks.

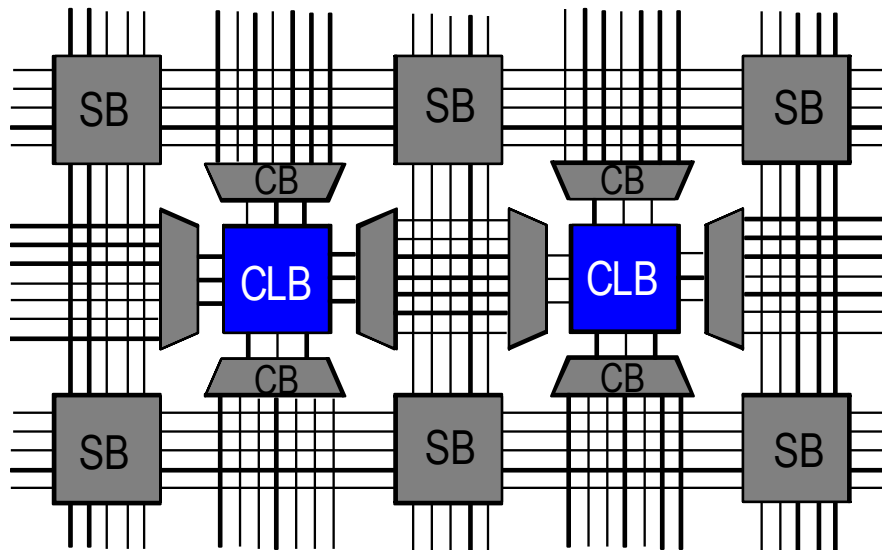


Figure 7. Island-style Field Programmable Gate Array organization

### 2.1.1.3 Logic Block Architecture

The logic block used in an FPGA strongly influences the circuit speed and area-efficiency. Many different logic blocks have been used in FPGAs, but it is possible to consider two main families: the gate based FPGAs and the *Look-Up Table* (LUT) based FPGAs. Most current commercial FPGAs use logic blocks based on LUTs.

Gate-based FPGAs closely resemble basic *Application Specific Integrated Circuit's* (ASIC) cells. The finest grain logic block would be identical to a basic cell of an ASIC and would consist of A few transistors that can be interconnected in a programmable way. The finest grain cell solution uses single transistor pairs. In [13], transistors pairs are connected together in rows. Within this pattern, the transistors are programmed to serve as isolation transistors or logic gates. Figure 8 illustrates how a function could be implemented. Figure 8-a shows the transistor pair tiles and figure 8-b shows a programmed function  $f=a.b+!c$ . The function is programmed by ensuring the correct connections between the different transistors. In figure 8-b, the dashed lines show the transistors that are turned off for isolation. The function is done by the two-input NAND gates formed on the left and right sides.

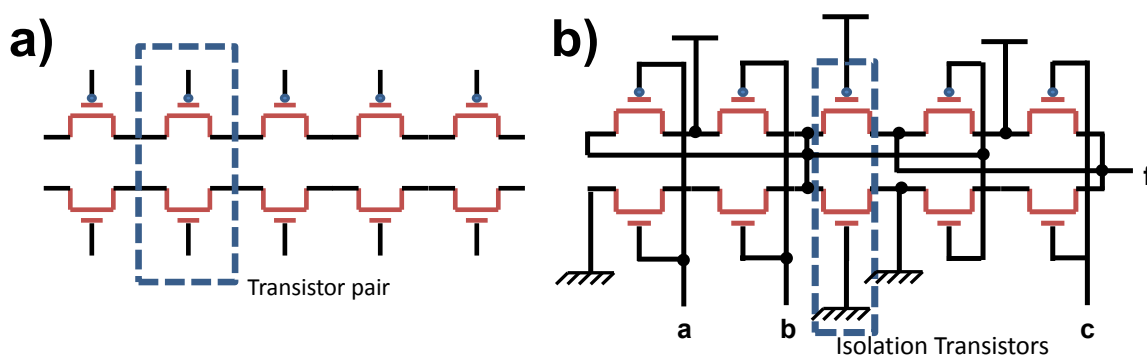


Figure 8. Transistor-level programmability for FPGA [13]

Instead of using programmability at the lowest transistor level, a two-input NAND gate has been used to realize the combinational block, as depicted in figure 9 [14]. The expected logic function is formed in the usual way by connecting the NAND gates together.

The main advantage of using fine grain logic blocks is that these blocks are fully utilized. This is because the logic synthesis techniques for such blocks are very similar to those for conventional mask-programmed gate arrays and standard cells. Then, it is easier to use small logic gates efficiently. The main disadvantage of these blocks is that they require a large number of wire segments and programmable switches. Such routing resources are costly in

terms of delay and area. As a result, FPGAs employing fine-grain blocks are in general slower and achieve lower densities than those employing coarse grain blocks.

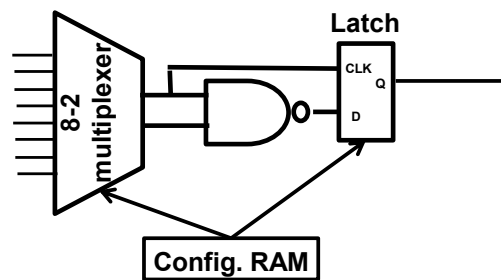


Figure 9. Simple gate FPGA logic block [14]

In [15], more complex logic blocks are proposed. Figure 10-a shows that these blocks are based on the ability of a multiplexer to implement different logic functions by connecting each of its inputs to a constant or to a signal. [16] presents a similar solution. Each input of the multiplexer and not just the select input is driven by an AND gate, as illustrated in figure 10-b.

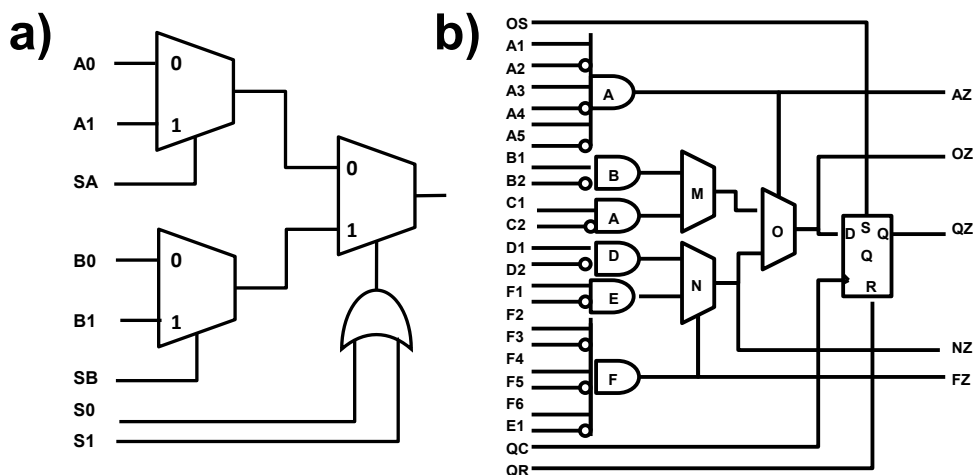


Figure 10. Complex gate FPGA logic blocks – a) Actel Act-1 LB [15] b) Quicklogic pASIC1 LB [16]

The alternating inputs to the AND gates are inverted. This allows input signals to be passed in true or complement form, thus eliminating the need to use extra logic blocks to perform simple inversions. Multiplexer-based logic blocks have the advantage of providing a large degree of functionality for a relatively small number of transistors. This is, however, achieved at the expense of a large number of inputs (8 in the case of Actel and 20 in the case of QuickLogic), which when utilized place high demands on the routing resources.

The other and most used approach is based on Look-Up Table. A LUT works thanks to memories driving the data inputs of the multiplexer. The truth table for a  $K$ -input logic function is stored in a  $2^K \times 1$  SRAM. The address lines of the SRAM function as inputs and the output of the SRAM provides the value of the logic function. For example, we consider the logic function  $f = a.b + !c$ . If this logic function is implemented using a three-input LUT, then the SRAM would have a 1 stored at address 000, a 0 at 001 and so on, as specified by the truth table. The advantage of look-up tables is that they exhibit high functionality. A  $K$ -input LUT can implement any function of  $K$ -inputs, i.e.  $2^{2^K}$  functions. The disadvantage is that they are unacceptably large for more than about five inputs, since the number of memory cells needed for a  $K$ -input lookup table is  $2^K$ . While the number of functions that can be implemented increases very fast, these additional functions are not commonly used in logic designs and are difficultly handled by the logic synthesis tools. Hence, it is often the case that a large LUT will be largely under-utilized.

Most modern FPGAs are composed not of a single LUT, but of groups of LUTs and registers with some local interconnect between them. A generic view for the LUT based logic block is shown in figure 11.

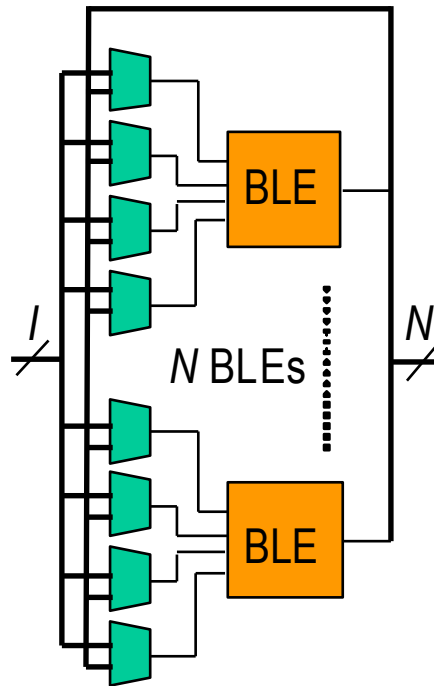


Figure 11. Configurable Logic Block architecture [17]

This Configurable Logic Block has a two-level hierarchy: the overall block is a collection of *Basic Logic Elements* (BLEs) [17]. As shown in figure 12, the BLE consists of a LUT and a register. Its output can be either the registered or unregistered version of the LUT output.

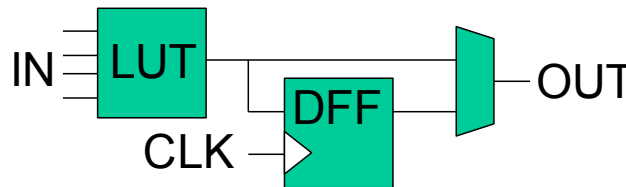


Figure 12. Basic Logic Element architecture [17]

This is how many commercial FPGAs combine a LUT and a register to create a structure capable of implementing either combinational or sequential logic. The complete logic block contains several BLEs and local routing to interconnect them. Such a generic logic cluster scheme is described by three parameters [17]: the number of inputs to a LUT ( $K$ ), the number of BLEs in a cluster ( $N$ ), the number of inputs to the cluster for use as inputs by the LUTs ( $I$ ). It is worth noticing that not all  $K.N$  LUT inputs are accessible from outside the logic cluster. Instead, only  $I$  external inputs are provided to the logic cluster. Multiplexers allow arbitrary connections of these cluster inputs to the BLE inputs. They also allow the connection of all the  $N$  outputs to each of the BLE inputs. All the  $N$  outputs of the logic cluster can be connected to the FPGA routing for use by other logic clusters. It is remarkable that each of the BLE inputs can be connected to any of the cluster inputs or any of the BLE outputs. Logic clusters are therefore internally fully-connected. This is a useful feature, as it simplifies *Computer Aided Design* (CAD) tools considerably.

The presented structure is a very generic representation of LUT-based FPGAs. In fact, logic blocks of FPGAs are more complex, as detailed in the next sub-chapter.

#### 2.1.1.4 Routing Architecture

The routing architecture of an FPGA is the manner in which the programmable switches and wiring segments are positioned to allow the programmable interconnection of the logic blocks.

Several routing architectures exist and come most of the time from results on the tradeoff between the flexibility of the routing architecture, circuit routability and density. Commercial routing approaches can be classified into three groups: row-based connections, island-style connections, and hierarchical scheme.

The row-based routing scheme is close to ASIC standard cells routing [15]. Effectively, logic blocks are organized in rows and a large number of horizontal wires are placed between the rows. Less vertical wires are used to connect rows to others.

Figure 13 depicts more precisely an island-style FPGA. Logic blocs are surrounded by routing channels of pre-fabricated wire segments on all four sides. A logic block input or output, which is called a pin, can be connected to some or all of the wiring segments in the channel adjacent to it via a connection block [18] of programmable switches. At every intersection of a horizontal channel and a vertical channel, there is a switch block [18]. This is simply a set of programmable switches that allows some of the wire segments incident to the switch block to be connected to others. It is worth pointing out that in figure 13, only a few of the programmable switches contained by switchboxes are shown. By turning on the appropriate switches, short wire segments can be connected together to form longer connections. In the figure, some wire segments continue unbroken through a switchbox. These longer wires span multiple logic blocks, and are a crucial feature in commercial FPGAs.

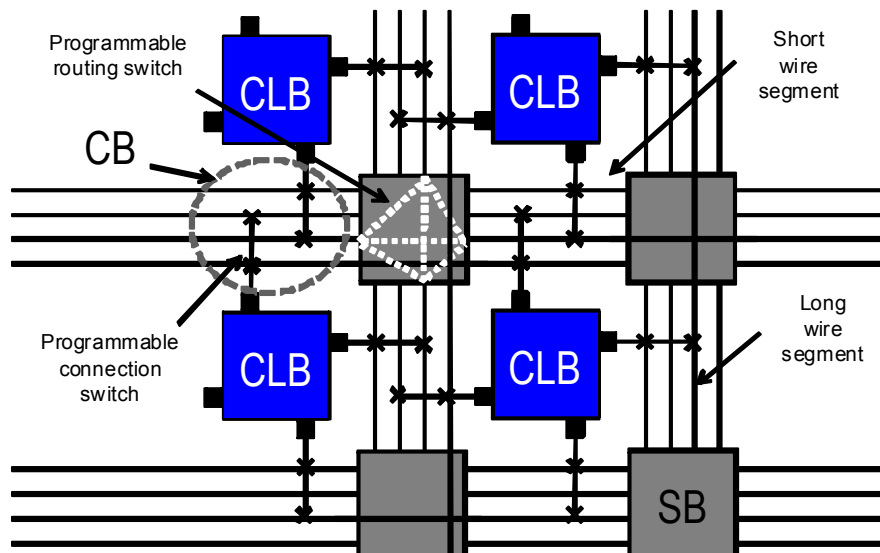


Figure 13. Details on the routing organization in an island-style FPGA [17]

The number of wires contained in a channel is denoted by  $W$ . The number of wires in each channel to which a logic block pin can connect is called the connection block flexibility, or  $F_c$ . The number of wires to which each incoming wire can connect in a switch block is called the switch block flexibility, or  $F_s$ .

Inspired from generic programmable logic devices routing schemes, the hierarchical routing scheme try to use some form of locality to obtain better density and performance. This hierarchy could be realized with several different interconnect schemes, like PAL or island-style. This is the most currently used interconnection in FPGAs. For example, considering the structure presented in figure 7 and figure 11, we could remark that a full-interconnectivity scheme is used for CLB, while island-style organization is used for the global FPGA.

#### 2.1.1.5 Architectural Parameterization and Optimum

The programmable switches introduced for routing purpose impact the performances. Thus, the FPGA architecture is the results of a trade-off between high versatility and cost/performance.

The adverse effects of the large size and relatively high parasitic of programmable switches can be reduced by careful architectural choices. By choosing the appropriate granularity and



functionality of the logic block, and by designing the routing architecture to achieve a high degree of routability while minimizing the number of switches, both density and performance can be optimized. The best architectural choices, however, are highly dependent on the programming technology used as well as on the type of designs implemented, so that no unique architecture is likely to be best suited for all programming technologies and for all designs.

A complete study of architectural parameterization has been conducted in [17, 19, 175]. These studies assume SRAM based homogeneous FPGAs. The principal results are summarized in the following.

- **Combinational Granularity Impact**

As the granularity of a logic block increases, the number of blocks needed to implement a design should decrease. On the other hand a more functional (larger granularity) logic block requires more circuitry to implement it, and therefore occupies more area. This tradeoff suggests the existence of an “optimal” logic block granularity for which the FPGA area devoted to logic implementation is minimized. It has been shown in [19] that the most suited LUT size is reached for  $K$  equal to 4.

- **Logic Block Sizing**

Several interesting sizing results can be taken from [17]. Firstly, the number of distinct inputs required by a logic cluster grows fairly slowly with cluster size,  $N$ . A cluster of size  $N$  requires approximately  $2N+2$  distinct inputs (for  $N \leq 20$ ). Secondly, because all the input and output pins of a cluster are logically equivalent, one can significantly reduce the number of routing tracks to which each logic cluster pin can connect,  $F_c$ , as one increases the cluster size. A good value for  $F_{c,output}$ , is found with  $W/N$ , while  $F_{c,input}$  is somewhat higher. Thirdly, logic clusters containing between 4 and 10 BLEs all achieve good performance, so any clusters in this range is a reasonable choice.

- **Routing Organization**

In [17], simulations have been carried out to study the impact of the routing structure. It has been shown that the most area-efficient routing structure is one with completely uniform channel capacities across the entire chip and in both horizontal and vertical directions. The basic reason is that most circuits “naturally” tend to have routing demands which are evenly spread across an FPGA. Furthermore, it has been shown that it is most important for FPGAs to contain wires of moderate length (4 to 8 logic blocks) even if commercial FPGAs are using some very short and some very long wires.

### 2.1.2 Market Trends

In the previous sub-chapter, we focused on the FPGA architecture and especially the generic homogeneous island-style architecture. While the island-scheme is used in most FPGAs, their structures are relatively complex. Indeed, the generic architecture is useful for understanding and research purpose, but it suffers from several performance issues in some application classes. Thus, several improvements have been proposed for modern commercial FPGA. The circuits are enhanced in many ways.

#### 2.1.2.1 *Increase of the CLBs Complexity*

In order to increase the logic block functionality, the structure has been customized and new features have been added. The logic block of a modern Xilinx Virtex-6 FPGA is shown in figure 14.

The block is still organized around LUTs and FFs. These elements give the combinational and sequential ability. The first novelty is the fractional behavior of the LUT. A large LUT is useful to realize a large and complex combinational function. However, this situation is not frequent. Instead of wasting a large amount of logic resources when the LUT implements a

small function, it is possible to split it into two independent smaller LUTs. It is then possible to optimize the logic fabric to the application, during the synthesis and packing operations.

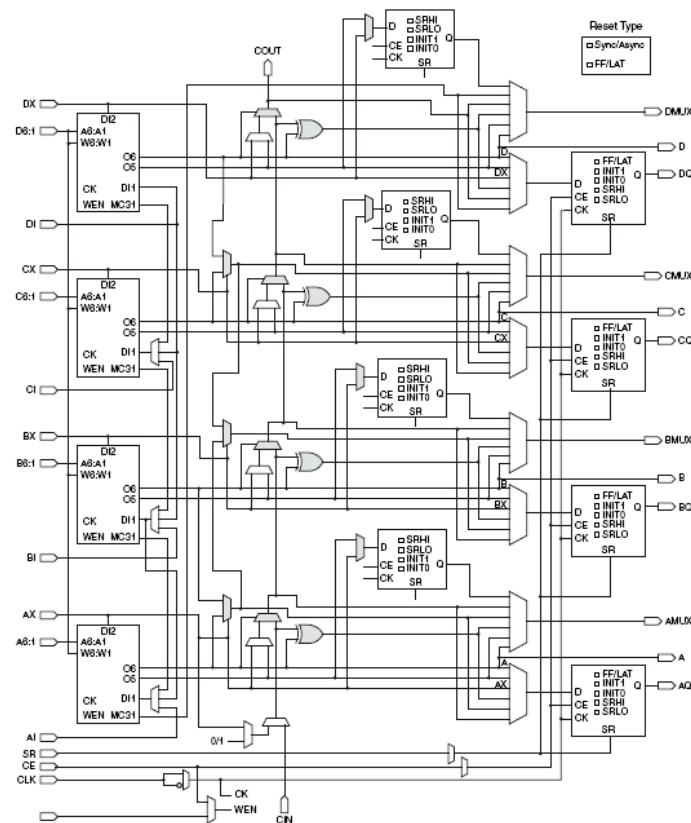


Figure 14. Virtex-6 logic block [20]

Within the CLB, extra circuitries have been added to address specific functionalities. We can point out in particular notice the carry lines, the SRAM decoders and the shift-register.

The carry lines extend through several LBs, in order to implement logic multipliers efficiently. Effectively, the specific track avoid to implement the carry through the global interconnect. Since global interconnect must be routed, delays for carry will be unknown and in any case larger than with the specific track. This allows significant increase in the performance of such a block without costing too many resources.

The use of FPGAs keeps increasing to implement complex SoCs. Thus, the requirement in terms of embedded memories is growing. LBs embed a large amount of SRAM memories for the LUT configuration. It is then of high interest to give the possibility to use them directly as standalone memories. In this configuration, the LUTs multiplexers are used as the address decoders, and a specific extra circuitry deals with the read management of the read/write control signals.

Finally, shift register behavior is also achievable, thanks to a dedicated module. Indeed, SRAMs can be cascaded and the writing is driven by a shifting register clock signal. This is obviously precious for specific applications.

### 2.1.2.2 Transition from Homogeneity to Heterogeneity

Associated with the complexity and diversity trends of the CLB architectures, the FPGA scheme tends to add more dedicated blocks and to heterogenize its structure. Indeed, the homogeneous FPGA is known to be slow and area costly for mathematical computation or floating point arithmetic. Hence, vendors have added dedicated co-processing logic units. These units are distributed through the logic grid. Figure 15 depicts the internal organization of a commercial Xilinx Spartan-3A FPGA. It is then possible to find *Digital Signal Processing* (DSP) blocks or *Multiplier-Accumulator* (MAC) units for computation, specific

blocks for encryption with specific models or even USB controller, Ethernet controller for communications.

Furthermore, requirements on memories are wide in current applications. FPGAs then integrate dedicated standalone RAM blocks directly reachable by the routing part. This allows the creation of complex SoC and the efficient instantiation of soft microprocessors on the structure.

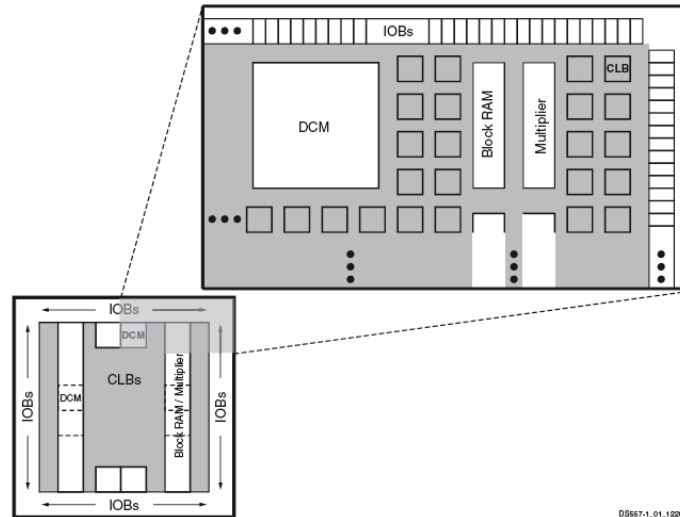


Figure 15. Xilinx Spartan-3A architecture organization [31]

### 2.1.2.3 Non-Volatility Features

In high production FPGAs, the configuration is stored by SRAM memories. These memories are found distributed through the entire circuits to store the information as close as possible to the data path logic. SRAMs circuits use the same integration process than the other logic parts. Thus, they are the simplest solution to implement distributed configuration logic. Nevertheless, they suffer from their large size and, in particular from their volatile behavior.

Due to the volatility, the configuration must be reloaded into the FPGA circuits at each power-up. This obviously leads to a large loss of efficiency in terms of delay and power consumption. Furthermore, a non-volatile storage is still required outside of the chip. In current circuits, standalone flash memories are used to store the configuration bitstream, and specific programming circuits have in charge the programming sequence. This is highly area and power consuming.

Then, it appears suited to use on-chip non-volatile memories directly. The current dominant non-volatile memory technology is the flash. Flash requires several technological steps in addition to CMOS. To take into account the complexity of flash/CMOS co-integration, the simpler solution is to co-integrate in the same package or directly on the same die the SRAM FPGA and its configuration flash memory. This solution is used in [31]. Such integration decreases the extra-chip requirements. Nevertheless, the most important hurdles still exist. Information is still duplicated in the flash and in the FPGA's SRAMs. This means that power consumption remains the same with a large wasted power in the distributed RAMs and a long power up time.

Thus, solutions have been envisaged to distribute the non-volatile memories through the logic grid. Flash based solution is proposed in [32, 33, 34, 35], while emerging solutions are shown in [37].

In [32], co-integration of flash and MOS transistors yield in a compact configuration memory nodes. The structure uses a non-volatile pull-up (pull-down respectively) network and resistive pull-down (pull-up respectively) network. This voltage divider arrangement allows the storage of the configuration and the drive of logic gates. This enables the fabrication of a non-volatile base LUT, as introduced in [33]. Complex co-integration permits to create

compact non-volatile switches for FPGAs. In fact, flash transistor could be seen as a programmable switch. The programming of a flash transistor needs specific voltages and circuitries. Traditionally, these circuits would be introduced into the data path. In [34, 35], a simple circuit composed of two flash transistors is presented (Figure 16). The particularity is that the floating gate is shared in order to connect a dedicated programming transistor with the data path one. This is of high interest for logic compactness, and distinct separation between programming and logic path.

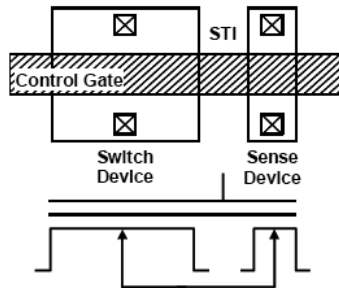


Figure 16. Schematic of a merged programming-data path floating gate transistor [35]

In [37], magnetic RAMs are used instead of flash. Indeed, emerging resistive memories are expected to reduce the production costs of the circuits. Resistive memories could be integrated after the back-end process. This means that the resistive devices are realized only after the costly CMOS process. This is obviously interesting for cost reduction and process simplification. Figure 17-a depicts the circuit that is used to store the information. The circuit is unbalanced flip-flop. Two magnetic tunnel junction memories store complementary informations that are used to start the flip-flop in a good configuration at power up. The programming of the magnetic memories is in charge of special writing lines. The magnetic memories are placed above the IC as depicted in figure 17-b. Nevertheless, we should remark that this solution is area consuming due to the presence of front end transistors in addition to the size of the back-end memories.

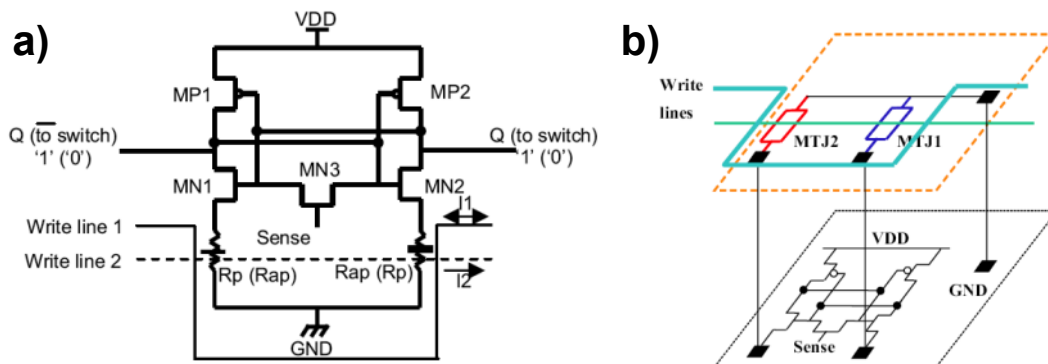


Figure 17. Magnetic Tunnel Junction unbalanced flip flop (a) and associated above IC structure (b) [37]

### 2.1.3 Limitations of FPGAs

The impact of area of each FPGA part has been objectively studied in [38]. This evaluation has been realized with a precise methodology using a stick diagram of the implemented circuits, instead of using a simpler minimum-width equivalent size transistor count. Figure 18 presents the area breakdown of the various components of a baseline Xilinx Virtex island style FPGA. It is worth noticing that the configuration memories occupy roughly half of the area in both the logic blocks and the routing resources. The logic blocks occupy only 22% of the whole area including its own configuration memory. Only 14% is then used for actual computation.

In addition to consuming most of the die area, programmable routing also contributes significantly to the total path delay in FPGAs. In [21, 22], interconnect delays are estimated and found to account for roughly 80% of the total path delay. Programmable routing also

contributes to the high power consumption of FPGAs. This problem has recently become a significant impediment to the FPGA adoption in many applications. The power consumption measurements of some commercial FPGAs have shown that programmable routing contributes more than 60% of the total dynamic power consumption [23, 24, 25]. As a result of these performance degradations, FPGA performance is significantly worse in terms of logic density, delay, and power than cell-based implementations. Finally, it is commonly admitted that FPGAs are more than 10 times less efficient in logic density, 3 times larger in delay, and 3 times higher in total power consumption than cell-based implementations [26].

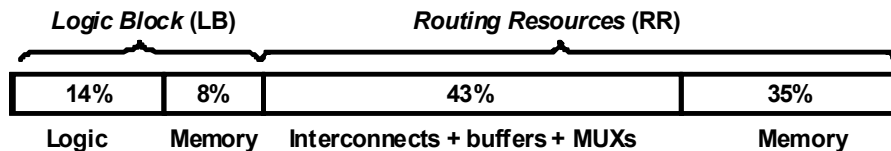


Figure 18. Field Programmable Gate Arrays area repartition per block [38]

## 2.2 Emerging Reconfigurable Architectures Overview

While we expect that reconfigurable architecture represent the future of computation architectures, we will assess the ways in which emerging devices could improve or even break the current FPGA model. This represents a difficult challenge because in many cases, circuit-level models and/or architecture-level models for these devices and their interconnect systems either do not exist or they are very primitive. Moreover, the applications under considerations for these new devices can take many forms; i.e.

- i) as a drop-in replacement for CMOS,
- ii) as additional devices that complement and coexist with CMOS devices, or
- iii) as devices whose unusual properties can provide unique functionality for selected information processing applications.

The ITRS – Emerging Research Devices [7] and especially the *Emerging Research Architectures* section aims to identify possible applications for emerging logic and memory devices. Two main families of architectures are defined: the morphic architectures and the heterogeneous architectures. We will survey the principal architectures that have been devised for computing at the nanoscale.

### 2.2.2 Morphic Approach

As indicated by the etymology, the morphic approach aims to develop architectures that are inspired from other systems. They are often inspired from biology as far as biological systems are highly efficient due to very limited power consumption and high system reliability. In this class of application, neural networks immediately spring to mind. An *Artificial Neural Network* (ANN) is a model inspired by the human brain [40]. Its goal is to reproduce some properties found in the biological organ. Globally, ANNs are used in the following applications which benefit from its properties related to memories and/or learning: sorting, associative memories, compression... The conventional implementation of neural networks is based on analog neuron and a digital control of the interconnectivity and edge weighting [39]. Such an implementation is far from the biologically inspired models. On the contrary, some implementations copy to emulate the behaviour and the structure of complex biological neural systems [41, 42]. They are called neuromorphics. Emerging devices are in this topic highly adapted for the hardware realization. In particular, [43] proposes to create reconfigurable hybrid CMOS/nanodevice circuits, called CMOL. In such a circuit, a CMOS subsystem with relatively large silicon transistors is used for signal restoration, long-range communications, input/output functions, and testing/bootstrapping. An add-on nanowire crossbar with simple two terminal nanodevices at each cross-point provides most of information storage and short-range communications.

The emergence of new devices and integration paradigms will certainly lead to several improvements of neuromorphic circuits. In this work, we are focusing on standard computation paradigms. Thus, we will consider that architectures suited for morphic applications are outside of our scope.

### 2.2.3 Heterogeneous Approach

Emerging technologies are expected to supersede CMOS in terms of functionality and performance metrics. Nevertheless, the transition to a new disruptive technology will not occur abruptly. In the near future, it seems reasonable to consider that standard CMOS circuits will be improved by new technologies. Two different means of improvement can be pointed out: the improvements coming from the increase of the device functionality (i.e. the use of devices with more functionality in the same area) and an improvement coming from the increase of integration density (i.e. more devices in the same area).

#### 2.2.3.1 Regular Architectures

Mono-dimensional devices improve the performance of transistors channel. Nevertheless, the ultra-scaled dimensions represent a real challenge for the integration of complex circuits. In particular, photolithography unreliability requires the use of high regularity. Regularity can be envisaged at the transistor level. The use of micro-regularity is of high interest to reduce the size of the circuits drastically. Furthermore, regularity is compatible with bottom-up fabrication techniques. These techniques open the way towards complex arrangements at the nano-scale and lead to the emergence of crossbar circuits. A crossbar is defined by the regular arrangement of devices in the array. This leads to the most integrated structures achievable by the technology.

The first crossbars realized from emerging devices were proposed in [44, 45, 46]. The basic function that those circuits are implementing is information storage. Subsequently, the use of dense crossbars computational units has been conceptually proposed in [47, 48, 156]. More precisely, such computation fabrics are based on semiconducting *Silicon NanoWires* (SiNWs) organized in a crossbar fashion. Active devices are formed at the cross-points. In [154], a Programmable Logic Array is proposed. The cross-points are realized by molecular switches. The switches can be programmed in order to perform either signal routing or wired-OR logic function. The structure is arranged in several sub-crossbars and presented in figure 19.

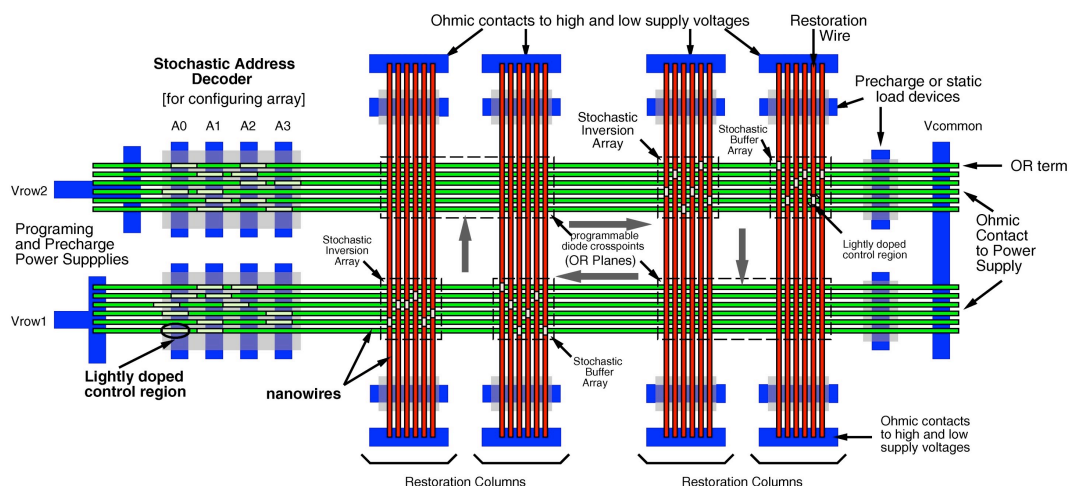


Figure 19. NanoPLA architecture [154]

The input of the structure is a decoder interface between the micro/nano worlds. The decoder is used to address every nanowire independently of the others. The decoder design assumes that the nanowires are differentiated by a given doping profile [49]. The output of the crossbar is routed to a second crossbar. The signals can be inverted by gating the nanowires carrying the signals. A cascade of these two planes is equivalent to a NOR plane. It is worth noting that, due to diode logic, a logic restoration stage is required. Indeed, a more than unity gain is

required to ensure a good cascade. At the end, the structure is duplicated many times and several stages are connected to each other in order to perform complex logic functions.

Instead of using diode logic, [64] uses FETs realized at the cross-points. While the previous approach implements a reconfigurable PLA circuit, this technology based on FETs address specific application-driven designs. Within this organization, a *Nanoscale Application Specific Integrated Circuit* (NASIC) is introduced in [65]. A NASIC tile consists of basic circuits such as adders, multiplexers and flip-flops. Circuits are realized using a dynamic logic style. Two clock transistors are placed between the power lines and a stack of transistors, which realize the logic function. It is possible to implement a standard AND/OR functions and their inverted counterparts. The implementation of the different logic functions is depicted in figure 20.

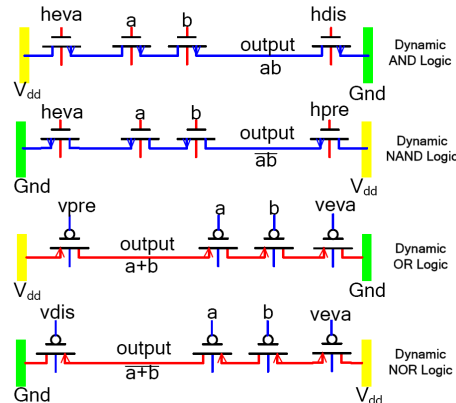


Figure 20. Dynamic logic implementation of AND, NAND, OR and NOR functions [65]

From this logic organization, two-planes are cascaded (AND/OR, NAND/NAND...). Figure 21 depicts the complete NASIC tile. A first set of AND logic functions is realized in the horizontal direction. Nano-scale wires are connected to micro-scale power lines and to the other blocks that are surrounding the crossbar core. The second horizontal AND functions are driving transistors into the vertical orientation. The vertical line set implements OR functions. As a complete illustration, the circuit, presented in figure 21, implements a 1-bit full adder.

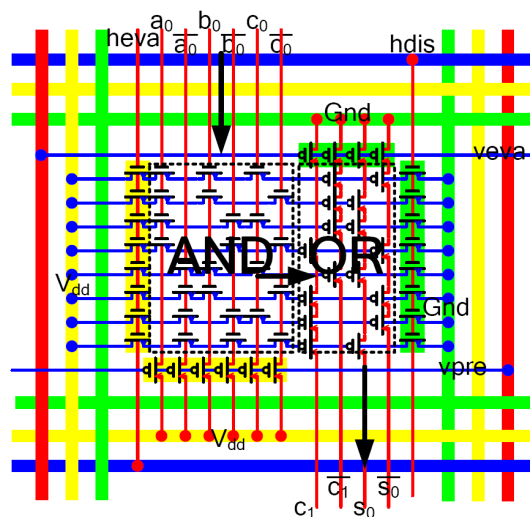


Figure 21. NASIC structure (implementing a 1-bit adder) [65]

The doping of nano-grid strips, the size of the NASIC tiles, the use of certain nano-scale (i.e. sub-lithographic) wires as interconnect between tiles and the micro-level interconnects are chosen in an application / architecture-domain specific manner. These aspects determine a *NASIC* fabric and are the key differentiators between PLA type of nanoscale designs [154] and *NASICs*.

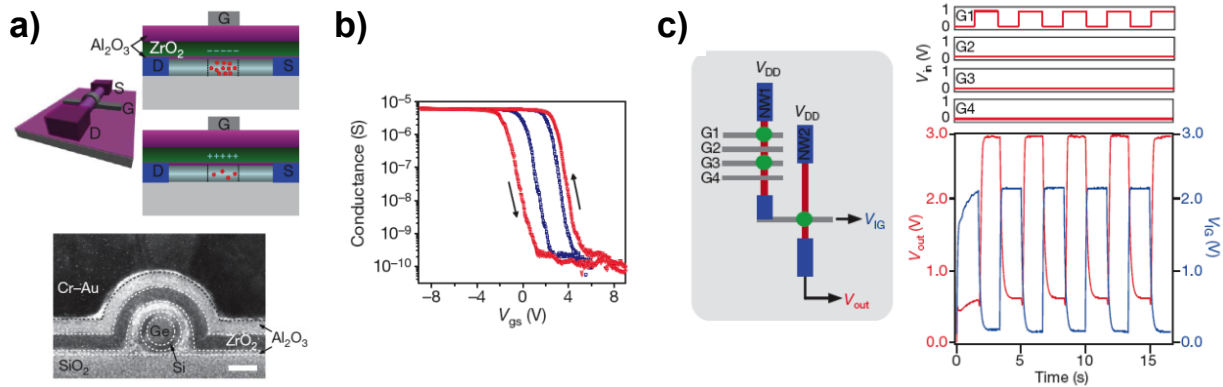
The most sensitive issue in these crossbar proposals remains the fabrication assessments. In fact, realization of 1-D structures and their alignment over long distance with a good aspect



ratio is extremely difficult to perform.

Several works have been published in order to assess the technological credibility of the structure [66]. Recently, a simple programmable crossbar-based processor has been fabricated. In [50], the authors have demonstrated a programmable and scalable architecture based on a unit logic tile consisting of two interconnected, programmable, non-volatile nanowire transistor arrays. The transistor structure is depicted in figure 22-a. S, D and G correspond to source, drain and gate, respectively. On top right, the hole concentration in a p-type Ge/Si NWFET for two charge-trapping states is presented. This charge accumulation is in charge of the hysteresis behavior of the conductance, depicted in figure 22-b.

Each NWFET node in an array can be programmed to act as an active or an inactive transistor state. This is done by charge trapping into the floating layer. By mapping different active-node patterns into the array, combinational and sequential logic functions including full adder, full subtractor, multiplexer, demultiplexer and D-latch can be realized with the same programmable tile. Figure 22-c presents a programmed logic operation on two multigate nanowires. Cascading this unit logic tile into linear or tree-like interconnected arrays is possible given the demonstrated gain and matched input–output voltage levels of NWFET devices. This provides a promising bottom-up strategy for developing increasingly complex nanoprocessors with heterogeneous building blocks.



**Figure 22. Structure of the programmable NWFET (a), associated characterization (b) and nanowire-nanowire coupled multigate device (c) [50]**

### 2.2.3.2 Enhanced Reconfigurable Architectures

In the previous part, emerging devices have been used to realize crossbar circuits. In fact, 1-D structures have been mainly envisaged to create dense interconnection networks or dense substrates to build active devices at the nanoscale. Nevertheless, some 1-D materials could be used efficiently to obtain new functionalities at the device level. For example, *Carbon NanoTubes Field Effect Transistor* (CNFET) exhibits an ambipolarity property [142]. This means that the same device could be controlled between *n*- or *p*- type, only thanks to the voltage applied to back-gate electrode. This property of CNTs is an opportunity that does not exist in CMOS technology. In [51, 52, 122], the benefit on logic circuit design is assessed.

In [51, 52], the main novelty is to leverage the ability of performing logic operations between the signals feeding both gates of ambipolar CNFETs. The design with such operations is demonstrated in a set of static logic families including combinations of transmission-gate/pass-transistor on the one hand and complementary/pseudo logic on the other hand. This yields to a natural, simple and efficient implementation of the XOR function in ambipolar CNT technology with almost no cost, as shown in figure 23. The basics of the proposition are to configure the polarity of the input signal. The polarity choice is obviously used by the configuration voltage and the transistor type.

In [122], a dynamic reconfigurable logic cell is introduced. As presented in figure 24-a, the cell is composed of seven double gate CNFETs organized in two logic stages: logic function and follower/inverter [122]. The polarities (*n*-type/*p*-type) of double gate devices  $T_1$ ,  $T_2$  and



$T_3$  are controlled by the corresponding back-gate bias voltages  $V_{bA}$ ,  $V_{bB}$  and  $V_{bC}$ . The cell may thus be configured to one of fourteen basic binary operation modes, as shown in figure 24-b. This cell will be presented in more details in chapter 5.

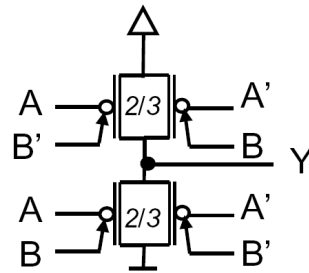


Figure 23. XOR function example built with four ambipolar transistors [52]

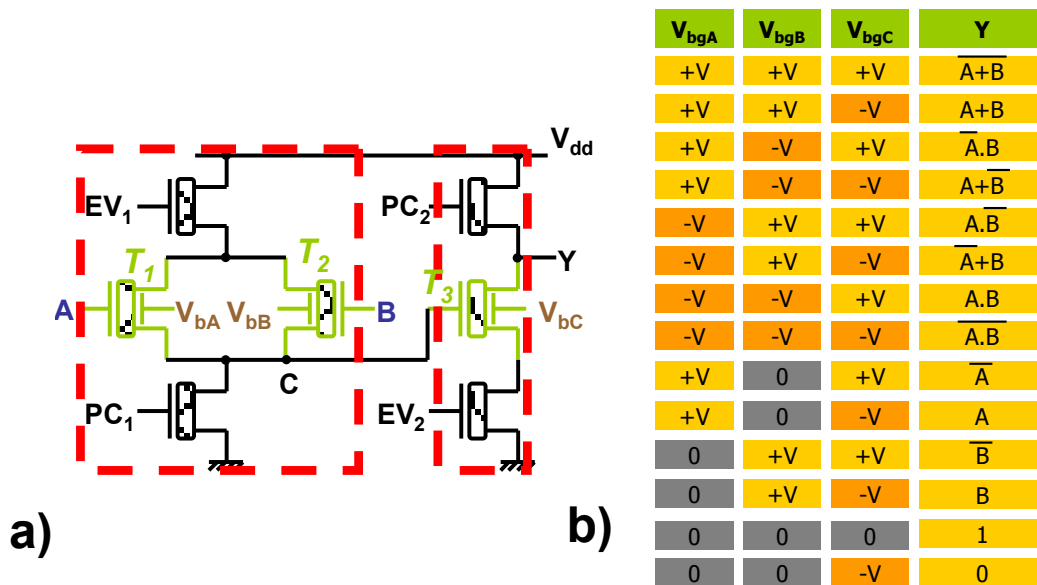


Figure 24. a) Transistor level schematic and b) configuration table for CNT reconfigurable cell [122]

### 2.2.4 Global Comparisons and Discussions

In order to compare the different architectures presented above, we will use five global comparison metrics: the area, the performance, the power consumption, the technological maturity and the fault tolerance ability. Table II presents the global results.

Table II. State-of-Art architecture global metrics comparisons

	Reference publications	Area	Performances	Power consumption	Maturity	Fault tolerance
CMOL	[43]	+	+	+	- - -	+++
NanoFabric	[47, 48]	++	-	-	-	-
NanoPLA	[154]	++	-	-	-	- -
NASIC	[64]	++	+	-	-	-
DG-CNFET	[52, 122]	+	+	++	-	+

It is worth noting that a solution based on sublithographic nanowires seems to increase drastically the area of the final circuits. This seems obvious considering that the integration density of the elementary devices is enhanced. Nevertheless, these solutions are quite worse in terms of other metrics. In fact, NW-based solutions require the use of complex clocking signals to restore the logic levels. This means that a large part of performance and power are wasted only for the periphery. In this context of intense power consumption, it is remarkable

that carbon electronics is promising due to the good intrinsic properties of the devices. Thanks to carbon technology, it is possible to build architectural solutions that are globally efficient in all the proposed metrics, even in terms of fault tolerance and reliability. Carbon electronics have been widely studied, to enhance the reliability of the fabrication processes, and to provide complete methodologies for robust designs [55, 56]. Finally, we should consider that morphic approach provides naturally a high robustness. Effectively, these approaches are implementing neural networks. This computation scheme is well-known for its versatility and performance regarding faulty tolerance.

## **2.3 Architectural Formalization and Template**

### **2.3.1 Global Statement**

From the existing state-of-the-art, it is possible to extract some global comments, regarding the design strategy and methodology.

Several works have been conducted on the fabrication technology process, reliability, circuit design, architectural definition and associated tools. Design with emerging technologies suffers from several unknown parameters. It is challenging to evaluate the expected performance of a whole circuit. Only a few works actually merge more than two aspects. Nevertheless, it is desirable to evaluate the potentiality of the technology within a complete architecture. The term complete covers the questions regarding the technological credibility of assumptions, the design of the computation part itself, the design of all the peripheral circuitries and the associated methodologies. This is required addressing all the previously defined aspects. As an illustration, the design of a compact logic gate does not make sense if the gate is not robust enough or requires significant additional circuitry.

Along the same lines, it is worth noticing that the work on global architectures is often completely uncorrelated from the technological assumptions. This is standard architecture design with mature CMOS process. In the case of an advanced process, it is hard to handle the design of architecture without credible technological assumptions. As an example, while the use of bottom-up aligned SiNW is credible, it is difficult to imagine the wires aligned over several micrometers at only a pitch of some nanometers. This will lead to an angle deviation less than  $10^{-5}$  degrees and is hard to handle today, as well in the near future. This makes that the architecture study must not be done under unrealistic assumptions, to ensure its credibility.

Finally, while it is recognized that variability and technological issues will increase drastically in the future, only a few works actually consider the questions of reliability. Conventional robustness techniques could be used, such as defect avoidance [53], error correction coding [54], or redundancy. Generally, unreliability in systems is overcome by duplicating the unreliable resource. Such spare circuits will then be used as a replacement part if the primary circuits are defective. However, a robust design approach is preferred to increase the reliability of a circuit early in the design [55, 56]. This allows to directly correct the misperformance at the level where it occurs, instead of a much higher level. Indeed, the correction will have a much larger and detrimental impact on the whole circuit if it is performed far from the origin of the issue. Even if the question of reliability must be tackled in future systems, we consider this specific point as out of the scope of this thesis.

To avoid this lack in methodology, it appears interesting to address the question of the architecture globally. Nevertheless, it is of course, not possible to handle it directly, due to the complexity of the design. Indeed, the problem concerns the handling of a complex mix between gates that leads to computational units, memories and routing with so many unknown parameters. In this thesis, we will use a generic template for reconfigurable architecture in order to organize the contributions into a hierarchy. The template will be presented in the following. Each layer will then be studied and optimized independently. A correct hierarchy organization ensures that the levels will not impact the others as long as interface

requirements remain the same. The optimization of specific parts will be examined in the following chapters. Finally, even if a layer is improved in terms of area or functionality, it is important to assess the impact on the architectural level. Thus, it will be necessary to define a benchmarking tool which is compatible with our architectural template, and which is able to implement much different architecture for each layer. This will allow exploring the architectural design space and permitting fast track optimization for the designs. This will be presented in the following.

### 2.3.2 Generic Architectural Template

In this thesis, we are focusing on reconfigurable architectures and the enhancement that could come from emerging technologies.

Reconfigurable architectures such as FPGAs are highly versatile and adaptive in terms of target application. They use regular architecture, with several blocks that are replicated through the circuit. We have previously seen that the FPGA structure is organized hierarchically, which seems a sound approach to manage the large amount of replicated blocks. In this context, it is difficult to see where emerging technologies will be useful to improve the performance of the structure. Furthermore, the FPGA structure has been designed for CMOS-based LUT logic. The structure is optimal for this application, but we could expect that the structure will not be the same in the case of emerging technology. Thus, it is of great interest to propose an architectural template based on a hierarchical organization of configurable and routing blocks, but also to generalize each hierarchical level to its global behaviour.

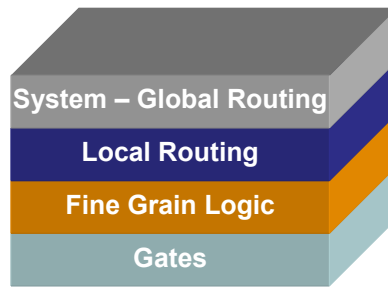
The generic template is shown in figure 25. The template is organized into four levels of hierarchy: the gate level, the Fine Grain Logic level, the Local Routing and the System-Global Routing level.

The Gates level corresponds to the elementary gates that are used to perform the computation. By computation is meant the combinational and the sequential operations. Computation in FPGAs is generally realized by LUTs, but it also could be, as we have seen in the literature review, custom logic from the simple transistor to the bigger logic function.

The Fine Grain Logic corresponds to the smaller autonomous block that could perform both combinational and sequential logic. We consider also at this level all the configuration memories that are used to program the structure. These memories give the fully programmable functionality of the block. In the FPGA scheme, this level corresponds to the BLE, i.e. a LUT which is used to perform a fine grain operation. The LUT output is routed to a latch but also directly to an output multiplexer, which is in charge of the path selection between the latches and the unlatched version of the signal.

The Local routing level aims to provide an arrangement of Fine Grain Logic and to provide it with a local connectivity pattern. Generally, this connectivity pattern is full and aims to provide a large kind of programmable paths at the lower level. This allows a huge simplification of the packing tools. In the FPGA scheme, this level corresponds to the CLB. Indeed, an assembly of BLEs are fully interconnected by a multiplexer based interconnect. Thus, it is possible to realize with such a block a coarser function with large synthesis simplicity.

At the final system level, the circuit is organized with a macro-regularity. Logic blocks defined at the previous level are regularly arranged and interconnected. The organization is generally an island-style scheme and the interconnect resources handle the resource limited interconnection pattern. This means that each logic blocks could be interconnected to others, but it is not possible to realize all the connections at the same time. In the FPGA scheme, we can also find the CLBs surrounded by the programmable interconnect. Several other interconnect patterns can be found with for example a local connection set between the logic blocks, such as those defined in [47].



**Figure 25. Layered organization of the generic template architecture**

This generic template represents a real opportunity for design exploration. Indeed, the denomination of a very generic template helps in enlarging each level. In fact, it allows extension of both technologies and architectures of each level, while keeping the others layer unchanged. Furthermore, it allows the study of each layer independently, while the assumptions on all the other part remain unchanged. This means that working on fine grain logic implementation may not have a negative impact on the system if all the assumptions required by this level are followed. In this way, the hypothesis driven by the different levels are useful to define the specifications of all layers and thus to constrain the development of a level.

## **2.4 Conclusion and Work Position**

In this chapter, we surveyed the state-of-the-art regarding reconfigurable architectures. The FPGA structure is the traditional reconfigurable architecture. The basement of the original FPGA is the hierarchic and homogeneous arrangement of logic blocks. The architecture suffers from the programming circuit in area. Indeed, the computation part of the FPGA structure occupies only a small amount of area, conversely to routing structures. Today, the structure evolutes to heterogeneity. Thus, it will reach new application classes and improve the computation/routing ratio. The heterogeneity is found at design level. Adjunction of several specific logic blocks such as memory blocks and DSPs allows increasing the versatility and the computation performance. Also found at the technology level, heterogeneity leads to the use of flash/MOS co-integration to non-volatile store the configuration.

The standard architecture is optimized for CMOS. Emerging technologies could leads to new architectural paradigms for reconfigurable computation. Emerging reconfigurable architectures have been previously envisaged through two approaches: the use of density-increased devices (i.e. ultra-scale devices) and the use of functionality-enhanced devices (i.e. devices with new highly functionality within the same area). Globally, the use of emerging technology leads to the improvement of performance compared to CMOS structures. Nevertheless, their immaturity and the technological hurdles make these solutions prospective.

In this work, we will use the emerging technologies to improve the standard FPGA approach. Two approaches will be used. We will first focus on the improvement of peripheral circuitries, i.e. memories and routing resources (chapter 3). These peripheral circuits are expected to improve the standard FPGA scheme by their direct use within the routing resources (chapter 4). Then, we will break the logic block paradigms (chapter 5) and propose a new seed for architectural organization (chapter 6). In addition, we will have in mind two requirements: the proposition of a credible fabrication process and the architectural compatibility with existing FPGAs.



## **CHAPTER 3** *Innovative Structures for Routing and Configuration*

---

### **Abstract**

The goal of this chapter is to demonstrate how emerging technologies can help to improve performance metrics of conventional Field-Programmable Gate Arrays structures. It is widely recognized that in traditional FPGAs, both the memory and the routing circuitry (with 43% of area for each contribution) represent the principal bottleneck to scaling and performance increase. In this context, we investigated three different technologies: a resistive memory technology, a monolithic 3-D integration and a vertical 1-D transistor technology.

The resistive memory technology, with its *Phase Change Memory* (PCM) incarnation, show significant promise with characteristics such as non-volatility, Back-End compatibility and low-on resistance (up to  $50\Omega$ ). We investigated the use of PCM to build an elementary configuration memory node for reconfigurable logic. An elementary circuit made of 2 resistive memories and 1 programming transistor able to store a configuration voltage was studied. We examined the proposed node in terms of area and write time and we assessed its impact on complex circuits. We show that the elementary memory node yields an improvement in area and write time of 1.5x and 16x respectively vs. a regular Flash implementation. Finally, we investigate the design of a switchbox, where programmable connections are formed of a unique resistive memory. This resulted in an area reduction up to 3.4x, while the low-on resistance of the memory is directly introduced into the data path.

The monolithic 3-D process opens the way towards the dense integration of configuration memories close to the associated logic function. Hence, an implementation of FPGA blocks, in a 3-D implementation scheme was proposed. The envisaged blocks are either routing structures or LUTs. It was possible to reduce the area up to 2x compared to equivalent 2-D bulk counterparts, while the performances are increased by 1.6x.

Finally, we examined the opportunity of a clearly disruptive vertical silicon field effect transistor. Their Back-End integration induced a low area impact, while their large channel size allowed the implementation of many high-performance configurable vias. In this regard, a novel implementation of logic gates fully benefiting from nanowire-based vertical transistors embedded within the metal lines was described. The logic design in this technology was explored and its performance was evaluated. A comparison made on an equivalent technology node showed that our cells reduce area and delay by a factor of 31x and 2x respectively.

In the previous chapter, we have seen that most area and performance metrics of modern Field Programmable Gate Array are limited by configurable interconnect and configuration memories. In this chapter, we will see how emerging technologies might be used to improve these two aspects. 3-D integration techniques will be surveyed.

### 3.1 Context and Objectives

#### 3.1.1 Context Position

In 2.1.3, the distribution of the area occupation between logic, memory and routing resources within an FPGA has been described. It shows that almost 45% of the silicon area is used only for the configuration memories, while the routing resources occupy 78% of the total area. In this chapter, we will address these two specific parts of the FPGA: memory and the routing resources.

The requirement for storage in an FPGA could be considered at several levels. This distribution means that designers face various constraints. Table III shows a study on memory requirements in an FPGA and gives an overview of associated constraints. We consider three different types of storage: high speed-data storage (fast flip-flops in the BLEs); configuration memories for logic (LUTs) and routing (Connection boxes and Switchboxes); and standalone Random Access Memories (Dedicated RAM Blocks). We focus on configuration memories (LUTs, CBs and SBs). These memories represent the biggest part of the FPGA area and share the same properties. They are distributed throughout the logic circuit and are programmed only a small amount of time. Traditionally, they are realized by SRAM circuits. SRAMs ensure a technical homogeneity between the logic and the configuration part.

Table III. Specification estimation wrt. memory distribution through an FPGA (extracted from Xilinx Virtex6 architecture [20])

	<b>Registers</b>	<b>LUTs</b>	<b>Connection boxes</b>	<b>Switchboxes</b>	<b>Memory blocks</b>
<i>Read operations</i>	$10^8 \text{ s}^{-1}$	$10^8 \text{ s}^{-1}$	10	10	$10^7 \text{ s}^{-1}$
<i>Write operations</i>	$10^8 \text{ s}^{-1}$	10	10	10	$10^7 \text{ s}^{-1}$
<i>Number</i>	$1.10^6$	$30.3 \cdot 10^6$	$20 \cdot 10^6$	$20 \cdot 10^6$	$38 \cdot 10^6$
<i>Distribution</i>	Even distribution throughout the architecture (fine-grain blocks)				Standalone blocks

Nevertheless, SRAM memories are power- and area-consuming, as well as volatile circuits. This means that the configuration must be loaded at each power-up, wasting time and power. Several circuits proposed the use of flash memories to create a non-volatile configuration [32]. However, the flash technology has a long programming time, and requires process co-integration. Indeed, floating-gate transistor processes require more steps than high-performance CMOS processes. This obviously incurs extra fabrication costs and technical difficulties. Such a non-volatile technology is thus adapted only to niche applications. Hence, the technology of configuration memories should be improved by using low-cost and non-volatile technologies.

Furthermore, the largest part of the configuration memories are used to configure the routing circuits (82% of the memory area). The reconfigurable interconnect alone occupies 45% of the area and introduces many active devices within the data paths. These devices impact directly the performance metrics of the structure, by increasing the critical path delay of the reconfigurable architecture. Thus, it is of great interest to consider the problem of “routing” in a global way by addressing the question of memory and the question of active configured devices at the same time. Hence, while memory and routing could be addressed separately, it is reasonable to work on both sides try to compact this entire periphery.

### 3.1.2 Objectives

In this chapter, we will propose the use of 3-D techniques, to place devices in the back-end layers. Three different technologies will be surveyed. The difference resides in the type of device that is placed in the back-end.

Firstly, *Phase-Change Memories* (PCM) will be used to embed a passive resistive memory above the IC. Such a device is non-volatile and technologically compatible (i.e. indicating homogeneous integration) with the CMOS technological process. Hence, we will propose an elementary memory node, able to store a configuration in the resistive state of the memory and to provide it intrinsically to a logic gate. The storage of a configuration for a resistive memory is quite obvious. Nevertheless, such a technology is of high interest due to the low on-resistance characteristic of the memory. This makes it possible to directly use a memory as a high-performance switch and to embed it directly within the logic data paths. We will thus propose a switchbox circuit that uses simple resistive elements to replace SRAMs and routing pass gates. All these circuits will be compared to their elementary CMOS FPGA counterparts.

Secondly, we will use monolithic 3-D integration technology to stack active devices with a high via density. Such a process allows the stacking of 2-D active devices. We propose to split the configuration memories and the data path transistors. This allows technological improvement of both classes of circuit. We suggest a complete integration of simple FPGA blocks, such as configuration memory, LUT and pass-gates, down to the layout level. In this way, we can provide a performance evaluation of elementary nodes with regards to standard CMOS FPGAs.

The previous technology proposes an integration of devices in a 3-D manner. Nevertheless, only 2-D devices are stacked. We therefore propose an integration process, which aims to realize a transistor (channel) in the vertical direction. A vertical NanoWire Field Effect Transistor process allows a vertical orientation of the active part of the transistor. It is then envisaged that several routing circuitries (such as programmable vias and signal buffers) can be embedded in the back-end layer. In order to evaluate the performance metrics of the technology and compare it to that of CMOS, we present a methodology based on TCAD simulations. TCAD will be used to model the elementary device, and electrical simulations of simple circuits will be performed.

Finally, we will draw a global comparison between the technologies and extract some overall conclusions.

## 3.2 Proposal 1: Resistive memory technologies

While Static Random Access Memories, *Dynamic Random Access Memories* (DRAMs) and Flash memories are predominant in microelectronics systems, thanks to their CMOS process compatibility, a large number of new memory devices have been highlighted by the International Technology Roadmap for Semiconductors [7]. These memories are generally based on new physical phenomena to retain the information and lift roadblocks to high density integration. In this sub-chapter, we will focus on non-volatile resistive memories.

### 3.2.1 Introduction

Next-generation *Non-Volatile Memory* (NVM) has attracted extensive attention due to conventional memories approaching their scaling limits. Several types of NVMs, such as ferroelectric random access memory, magnetic random access memory, and *Resistive Random Access Memory* (ReRAM), are being investigated. Among various NVMs, ReRAMs are typically composed of a simple metal-switching element-metal structure, which has the merits of low power consumption, high-speed operation, high-density integration and CMOS process compatibility.

Resistive memories, which can see their resistance vary depending on the applied voltage, were intensively studied from the 1960s to early 1980s for device applications [57]. Several materials can be envisaged to execute this desired functionality. The type of material



determines the physical phenomena that are used in the resistive change. For example, chalcogenide materials, semiconductors, various kinds of oxides and nitrides, and even organic materials were found to have resistive memory properties. Hence, the architecture will depend to a large extent on the technology. We can classify the technologies into two main families: *Oxide Memories* (OxM) and Phase-Change Memories.

Oxide memory state change is accomplished by the creation or the destruction of a conductive bridge through an oxide layer. This property is due to different physical phenomena which depend on the material. The conduction forming mechanism is still not fully understood, and is currently under investigation. The structure is composed of a changeable resistance material sandwiched between two terminal electrodes. Resistance change can be achieved by controlling the current or voltage pulse applied to the electrodes, and the resistance state remains stable without being refreshed. To date, a number of different switching characteristics have been observed in a variety of material systems; including NiO<sub>2</sub> [58, 59], TiO<sub>2</sub> [60], HfO<sub>2</sub> [61], WO<sub>x</sub> [62], CuO<sub>x</sub> [63], TaO<sub>x</sub> [67]. In fact, it has become well understood that a number of combinations of an oxide with metal electrodes can exhibit some kind of resistance switching behavior.

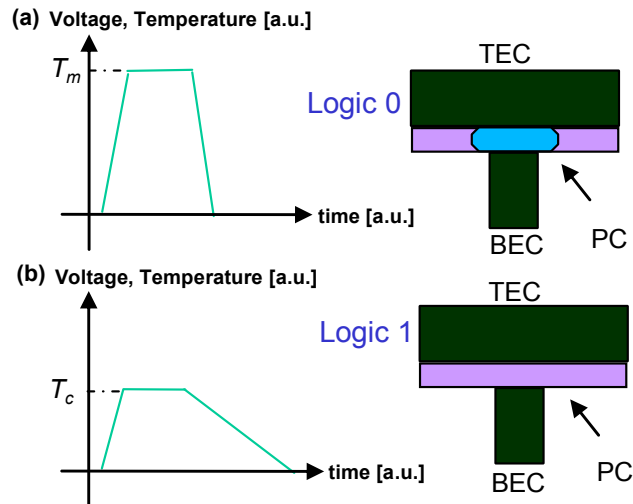
Phase-Change Memories are, as the name indicates, based on a material having two different stable physical phases leading to two different resistances. As with OxM, several materials might be used such as GeSbTe [82], GeTe [83], GeTeC [116] ... The PCMs are considered today to be one of the most promising candidates for the next generation of non-volatile memory applications [70]. The interest in PCMs is due to various advantages, including: better scalability (down to a few nanometers) [71], faster programming time (of the order of few nanoseconds) [72] and improved endurance (up to 10<sup>9</sup> programming cycles) [73]. Some prototypes (such as a 60-nm 512-Mb [74] and a 45-nm 1-Gb [73] PCM technology) have been presented recently to showcase the viability of high-density standalone memories based on PCM technology from an industrial point of view. The PCM technology achieves the maturity required for large applications. We will focus on this technology. Nevertheless, it is worth pointing out that this work can be generalized to any other resistive memory technology.

### **3.2.2 Phase Change Memory Properties and Technological Assumptions**

#### **3.2.2.1 Physical Phenomena**

A PCM device is based on the electrothermal-induced reversible phase transition of a chalcogenide alloy between an amorphous insulating state (RESET) and a polycrystalline conductive state (SET). The polycrystalline phase is inherently stable, as it is the lowest possible energy state of the system. On the other hand, retention instability affects the amorphous phase through two physical phenomena: spontaneous crystallization and low-field conductivity drift [77]. Recently, many efforts have been made in order to achieve a better understanding of the physical mechanisms which govern the behaviour of amorphous chalcogenides integrated in PCM cells [77, 78, 79, 80]. Chalcogenide alloys are semiconducting glasses made by elements of the VI group of the periodic table, such as sulphur, selenium and tellurium. The best known and most widely used chalcogenide alloy is Ge<sub>2</sub>Sb<sub>2</sub>Te<sub>5</sub> (GST). GST can guarantee stability of programmed amorphous bits for more than 10 years at 85°C [81]. While this can be considered to be sufficient for consumer applications, many efforts are today devoted to the development of new chalcogenide materials to improve the high-temperature reliability of PCM technologies in order to address the embedded memory market as well. Recent findings show that GeTe thin films demonstrate a higher crystallization temperature than GST [82], as well as superior data retention performances when integrated in memory cells [83]. Furthermore, it is known that the crystallization process is affected by the presence of foreign atoms in the material. For example it has been demonstrated that the Nitrogen doping of GST dramatically increases the stability of amorphous phase [84, 85].

Indeed, by means of a careful control of Joule heating through the cell, it is possible to electrically switch the chalcogenide layer between its two stable configurations, i.e. the high-conductive polycrystalline state and a low-conductive amorphous one, as shown in figure 26. A sufficiently high voltage pulses heat into the *Phase-Change* (PC) layer above the melting temperature of the material ( $T_m$ ). A rapid quench follows and part of the chalcogenide alloy (depicted as an oval in the PC layer) is stuck in the amorphous phase. The resulting memory cell is in a high resistance state (Figure 26-a). A lower but longer pulse is used to crystallize the amorphous region of the PC layer in order to achieve a low resistance memory cell (Figure 26-b).



**Figure 26.** Schematic of PCM device in Logic 0 (named RESET) and Logic 1 (named SET) configurations and of the programming pulses suitable to obtain the states.

### 3.2.2.2 Technological Assumptions

PCM technology is CMOS-compatible. As in Flash-NOR arrays, each memory cell includes a storage phase-change node and a selector transistor in series (i.e. *1-resistor-1-transistor* configuration). The memory element may be fabricated either just after the Si contact forming step at the *Front-End-Of-Line* (FEOL) level or after the first steps of interconnections at the *Back-End-Of-Line* (BEOL) level, (e.g. on top of the Metal 0 or Metal 1 interconnect level) [75]. A schematic cross-section of the storage element architecture is shown in figure 27. The PCM device, formed of a PC layer with Bottom (BEC) and Top (TEC) Electrode Contacts, is integrated between M0 and M1 interconnection level in the back-end-of-line. The MOSFET selector (bottom) is fabricated in the front-end-of-line. This figure depicts a pillar structure. The pillar approach is the simplest way to create a PCM device. First, a metallic heater is built. The heater is made by etching a via into the inter-layer dielectric and by filling it with a metal. The role of the heater is to help to channel the current in order to increase its density and thus maximize the heat control in the memory node. To improve the heater fabrication, several sublithographic techniques have been proposed [68, 69]. After the heater metal deposition, the via is filled by chalcogenide alloys with a room temperature deposition. The top electrode is obtained by a final metal deposition.

### 3.2.2.3 Opportunities

Resistive memories, and especially the envisaged Phase-Change technology, represent truly promising opportunities for several aspects of design. Indeed, PCMs demonstrate non-volatile behavior at low cost. Such a property is obviously of high interest for all types of reconfigurable circuits, where a permanent configuration circuit is strongly desirable. Furthermore, the technology is fully compatible with Back-End Of Line and able to integrate the memories into the 3<sup>rd</sup> dimension. This makes the resistive memories highly relevant as configuration points, since we could expect promising size reduction through the integration, above active silicon, of all area-hungry memories. Furthermore, the resistance of the on-state is typically below 1k $\Omega$  (for example that of GeTe is around 50 $\Omega$ ). This is far lower than any

MOS switch. Thus, it makes sense to use them as a high performance switch replacement element for FPGAs, by directly introducing them into the logic data path. In such a way, we expect not only to improve the size of the routing elements, but also to drastically reduce the delay of implemented circuits.

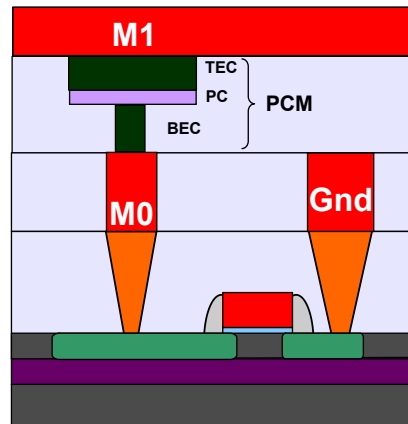


Figure 27. Cross sectional schematic showing a PCM device integration.

### 3.2.3 Elementary Memory Node

In this section, we present an elementary circuit, based on a PCM non-volatile resistive memory, used to move most of the configuration part of reprogrammable circuits to the back-end, reducing their impact on front-end occupancy. Such a memory node is dedicated to drive *multiplexer (MUX)* inputs or pass-gates. The memory node is programmed by injecting a certain current through it; while the information has to be read as a voltage level. Furthermore, it shall allow a layout-efficient line sharing.

#### 3.2.3.1 Concept

The elementary memory node is presented in figure 28. The circuit consists of 2 resistive memory nodes connected in a voltage divider configuration between 2 fixed voltage lines. A transistor is also connected between the ground and the output node of the cell. It is used to select the node during the programming phase. The output is designed to place a fixed voltage on a classical standard cell input. Read operations are intrinsic to the structure, while programming is an external operation to be performed on the cell.

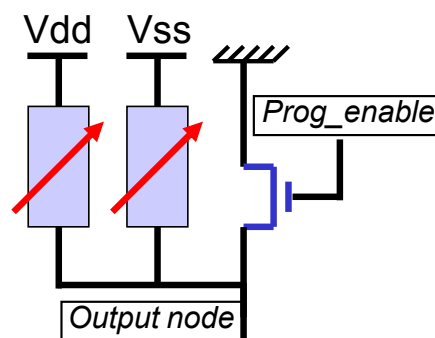


Figure 28. Logic-in-PCM elementary memory node

The voltage divider is used in this topology to execute intrinsically the conversion from a data stored in a variable resistance to a voltage signal. Figure 29 shows a configuration example where the node stores a '1'. The programming transistor is placed in the off-state by the non-active *Prog\_Enable* signal, so that the ground is disconnected from the output. The resistive memory (1) that is connected to the  $V_{dd}$  line, is configured into the crystalline state, so its associated resistivity is low (a few  $k\Omega$ ). The other memory (2), connected to  $V_{ss}$ , is in the amorphous state with high resistivity (close to  $1M\Omega$ ). As a consequence, a voltage divider is configured and the output node is charged close to the voltage of the branch with a high conductivity. The logic levels depend on  $R_{ON}$  and  $R_{OFF}$  as in the following relations:

$$'1' = V_{dd} - \frac{R_{ON}}{R_{ON} + R_{OFF}} (V_{dd} - V_{SS}) \quad '0' = \frac{R_{ON}}{R_{ON} + R_{OFF}} (V_{dd} - V_{SS})$$

It is also worth noticing that in continuous read operation, a current will be established through the resistors. This leads to a passive current consumption through the structure depending to the following relation:

$$I = \frac{V_{dd} - V_{SS}}{R_{ON} + R_{OFF}} \approx \frac{V_{dd} - V_{SS}}{R_{OFF}}$$

As an illustration, the PCM technology off-resistance is around the M $\Omega$  value. At  $V_{dd}=1V$ , the technology yields to a leakage current of 1 $\mu$ A. This is obviously too high for a viable industrial solution. Nevertheless, this static current could be reduced by the choice of a memory technology maximizing the  $R_{OFF}$  value (e.g. OxRAM technology exhibits off-resistance bigger than the G $\Omega$ )

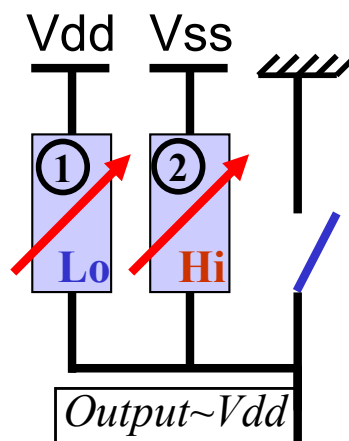


Figure 29. Node in read configuration

### 3.2.3.2 Programming Circuitry

Figure 30 presents the programming phase of the node. While the programming transistor is placed in the on-state by setting the *Prog\_enable* signal, the fixed read voltage sources are disconnected from the top lines and replaced by the programming unit. Then, a programming current is applied sequentially into the resistive memories to change their states. Programming currents are drained to the ground.

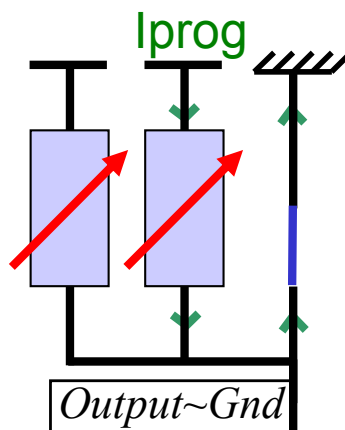


Figure 30. Node in write configuration

As each cell has its own selection transistor, the programming lines can be shared in a standalone-memory-like architecture, as shown in figure 31. The programming unit is composed of 3 different elements. A programming pulse generator handles the creation of the programming pulse with the correct waveforms. The programming signals are then routed by two stages of multiplexers. They are routed through  $BL_{XA}$  or  $BL_{XB}$  lines in order to program memory A or B respectively. We also note the Program/Operation selectors. Their aim is to

route static voltages when the nodes are not under programming. During the programming, the selection of a node is ensured by the  $WL_X$  signals. Thus, the choice of the memory node to program is made at the programming unit level, through the selection realized by the Program/Operation selectors and the Memory Program selector.

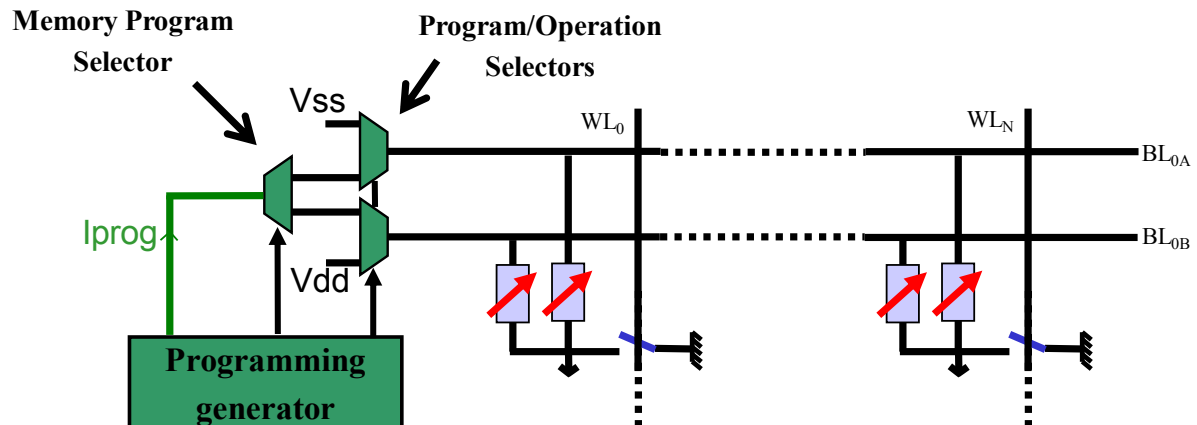


Figure 31. Line sharing illustration in standalone-memory-like architecture

### 3.2.3.3 *N*-ary Logic Generalization

In the preceding section, the presented memory node aims to store a Boolean value, represented by  $V_{dd}$  or  $V_{ss}$ . It is worth noticing that a generalization of the structure could be done, in order to address the *N*-ary logic. Figure 32 presents a ternary logic node. The structure is again based on a voltage divider scheme. For each logic voltage, a branch is added. Each branch contains a resistive memory and is connected to the logic power supply corresponding to the expected voltage to place at the output. Thus, the application of the expected logic value at the output node is done by programming only the corresponding memory into its low resistance state, while the others remain in the high resistance state. With this configuration, the output node is pulled close to the expected value. The programming of the structure is performed in the same way as for the binary node. In addition to the previous programming circuitry, a Program/Operation selector is added to each power supply branch, in order to route the programming pulses.

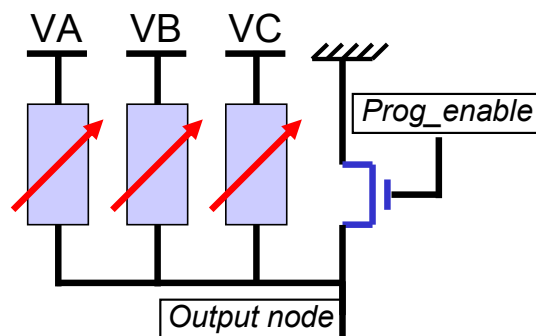


Figure 32. Ternary logic node based on PCMs

### 3.2.3.4 *Technological Integration Opportunities*

The process flow described in 3.2.2.2 is able to integrate the proposed elementary memory node. Memories will be realized in Back-End Of Line, while the access transistor is built in the active Front-End silicon. The cross-view of the fabricated circuit is depicted in figure 33. Nevertheless, we have to keep in mind that this node is intended to place the configuration memory of an FPGA above the programmed logic. Thus, it is of high interest to look at full back-end integration of the node, by pushing the transistor up to the metal lines.

Several disruptive techniques can be envisaged to further this approach. For example, vertical FETs could be envisaged to build a transistor between two metal lines (this technology will be studied further in this chapter). Although the chalcogenide alloys are deposited as oxides at room temperature, we will try to integrate the transistor with oxide technology. The

realization of oxide heterojunctions has been demonstrated in [88]. The junction is a  $\text{CuO}_x$ - $\text{InZnO}_x$  junction. In this demonstration,  $\text{CuO}$  acts as a *p*-doped material and  $\text{InZnO}_x$  acts as an *n*-doped material. Thus, the realization of a *Bipolar Junction Transistor* (BJT) could be envisaged by the deposit of another layer of oxide. To achieve the base behavior, the inner layer must be thin enough. In this context, the proposed process integration envisages a vertical BJT.

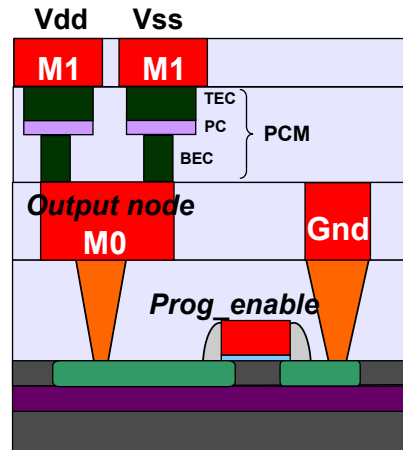


Figure 33. Cross-sectional view of the memory node using standard process integration

Figure 34 presents the cross-sectional view of the memory integrated with the oxide BJT. First of all, the used selector depicted is a *pnp*-heterojunction formed of the (*p*-doped) GST layer, an  $\text{InZnO}_x$  layer and a  $\text{CuO}_x$  layer. In this structure, the oxide layers are deposited successively. Contacts are patterned lithographically to obtain the base and the collector contacts. It is also worth noticing on this figure that the phase-change material is shared between all the memory nodes. This allows avoiding the bottom metallic connection between the nodes and the selector. To do so, it was required to invert the traditional scheme for the heater. Heaters are now built above the phase-change material. Finally, such a structure is connected to logic underneath with a standard via process. Thus, it is possible to completely embed the memory node into the back-end layer.

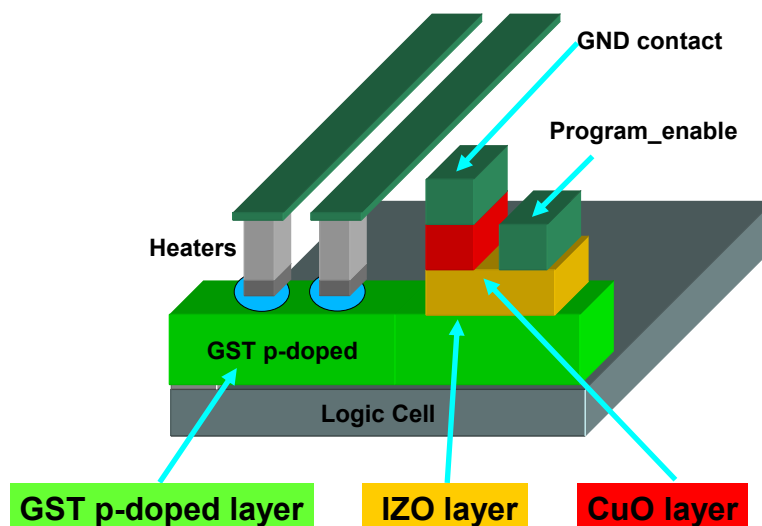


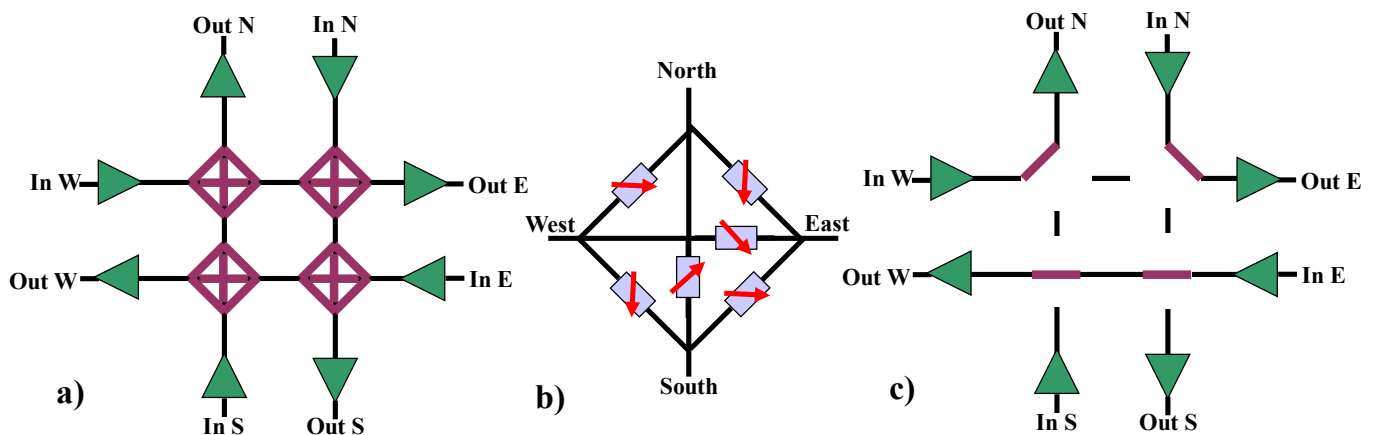
Figure 34. Full Back-End integration illustration of the PCM-based memory node using an oxide BJT

### 3.2.4 Routing Elements

In the previous section, we presented a memory node, which aims to replace the configuration memories in FPGAs. The node is intended to drive any logic gate and is thus a straightforward replacement part for SRAMs. In 3.2.2.3, we observed that the use of Resistive RAMs is also of great interest for performance improvement of data paths, thanks to their low on-state resistance.

### 3.2.4.1 Concept

The most important structures for FPGA routing are switchboxes. We thus propose to replace the traditional pass gates with the structure shown in figure 35. Figure 35-a shows a schematic diagram of a 2x2 crossbar structure using PCM resistances. At each cross point in the crossbar, routing elements are placed. They are able to create any combination of connections between the North/South/East/West terminals of the cross point. These routing elements are built with a similar structure as those of CMOS, whereby two-terminal PCMs replace pass-transistors (Figure 35-b). An on-connection is realized by programming the PCM, situated between the two wires that it should connect, to a low resistance state. With this structure, we build an “intelligent cross point”, which is able to merge the programming node and the pass switch in a single device. The device is embedded in the Back-End levels as a via, and replaces a 5-transistor SRAM and a 4-transistor pass gate. Furthermore, since the PCM itself performs the switching in the data path, its own on- and off-resistance values also have a direct impact on the circuit performance.



**Figure 35.** (a) PCM-based 2x2 switchbox architecture (b) zoom on the cross point structure (c) example of programmed switchbox (memories in off-state are not shown for clarity)

Once programmed, the cell behavior is purely static. Conductive paths are created by the resistive networks as illustrated in figure 35-c. In this figure, we observe 3 conductive paths through the box:  $InW \rightarrow OutN$ ,  $InN \rightarrow OutE$  and  $InE \rightarrow OutW$ . In this structure, the number of reachable inputs/outputs depends on the  $R_{on}/R_{off}$  ratio. In such crossbar architecture, the lowest resistive memory paths define the connections. Nevertheless, currents will also be established in the high conductive bridges. Thus, the discrimination between a conductive and a non-conductive path rely only on the path resistance difference. We consider the situation of figure 35-c. The path  $InE \rightarrow OutW$  goes through two on-resistances. This resistance is thus  $2R_{on}$ . The path  $InE \rightarrow OutS$  goes through a unique off-resistance. The resistance is thus  $R_{off}$ . In this example, it is not possible to discriminate the path if  $R_{on}/R_{off}$  is too low (i.e. 2 in the extreme bound) More precisely, the longest on-resistance path of the structure should be compared to the shortest off-resistive path. The tolerable  $R_{on}/R_{off}$  ratio is thus given by:

$$\frac{\text{Longest Path } R_{on}}{\text{Shortest } R_{off}} \geq 1000$$

(with 1000 an arbitrary choice to ensure a significant discrimination). More considerations on the structure can be found in [89].

### 3.2.4.2 Programming Circuitry

A PCM is programmed by applying a pulsed signal between its two terminals. This conducts the PCMs of the switchbox to be addressed sequentially. In the structure, the PCM is selected by connecting one terminal to the programming unit, while the other is grounded. When a unique memory is select, the other PCM terminals are left floating to avoid parasitic programming. After the selection, the programming unit drives the desired set or reset pulse to program the resistivity state. An example of sequential programming is shown in figure 37.



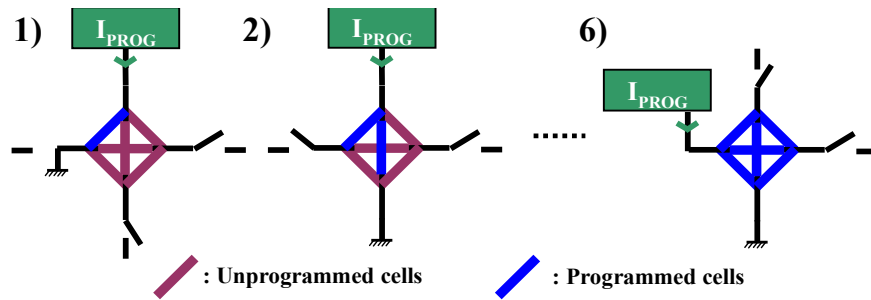


Figure 36. Example of switchbox programming sequence

The programming voltages and timing pulses required to program the PCMs may be applied through the drivers at the inputs and the outputs of the switchboxes (Figure 35-a). Figure 37 shows a possible structure for these drivers. They are used for the electrical interface between signal channels and the programming unit, which generates the configuration waveforms. Figure 37-a shows a possible input driver, while figure 37-b shows the implementation of the output driver. As explained, the buffers must allow the connection of the nodes to the programming unit, to the circuit or to a high-impedance node. This last possibility is handled by the 3-state multiplexers and buffers. The programming unit is routed by the multiplexer for the input, and through a single pass transistor for the output.

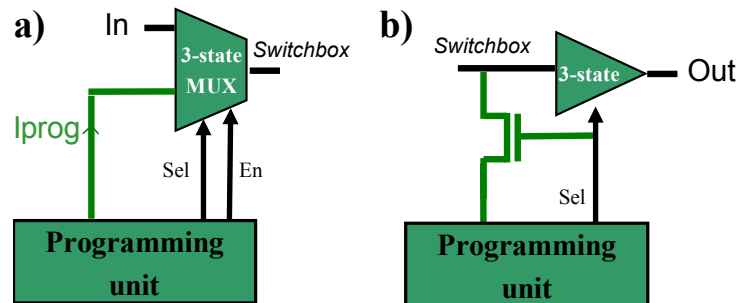


Figure 37. Input (a) and output (b) drivers for PCM-based crossbar

### 3.2.5 Functional validation and performance characterization

#### 3.2.5.1 Methodology

To validate and characterize the elementary memory node and the switchbox implementations, electrical simulations have been performed using a compact model. We first validated the node behaviour in transient configuration. Then, we characterized the performance of the node in terms of area and write time. Finally, comparison with memory elements traditionally used in FPGA, such as MOS SRAM 5T [17] and flash memories storage LUT elements [76], are used to benchmark the structure.

Since the technology enables most of the configuration circuit to be placed in the back-end levels, the area metric corresponds to the front-end projection of the impact of the memory node. The write time is also considered. Even if the use of non-volatility makes this consideration far less critical, it is still of interest to improve the programming time, in order to enable fast reprogramming.

#### 3.2.5.2 Compact Modeling

A behavioral compact model has been used to allow fast prediction of the behaviour of PCM cells [90]. The model is based on two parts. The static part reproduces the current/voltage characteristics and the dynamic part allows the programming of the cell. This compact model is made of five modules: a thermal module to derive the heating of the chalcogenide material, a temperature analysis module to detect if any phase change occurs, a time analysis module to detect whether the material can become amorphous, a storage module for the data retention and a module to derive the resistance. Figure 38 shows the architecture of this model.



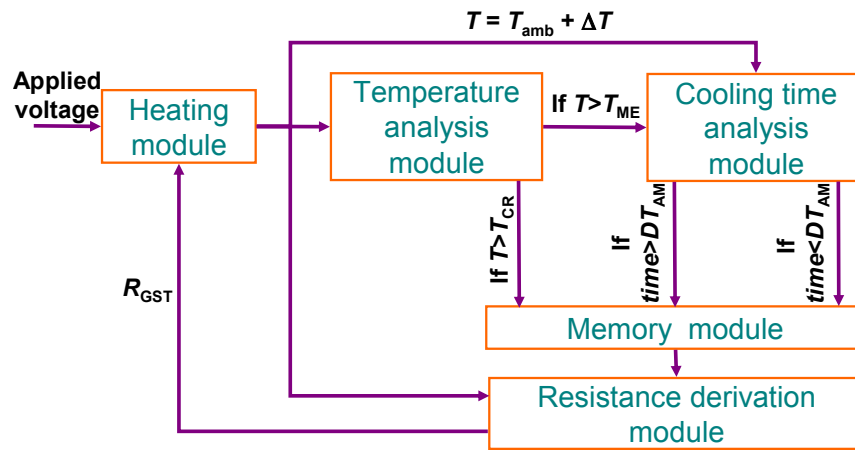


Figure 38. Synopsis of the behavioural model with the five modules.

In the figure,  $T$  represents the temperature of the cell,  $T_{amb}$  the ambient temperature and  $\Delta T$  the heating.  $T_{ME}$  is the melting temperature and  $T_{CR}$  is the crystalline one.  $DT_{AM}$  is the maximum time of cooling for the cell to become amorphous, while  $R_{GST}$  is the resistance of the material. Beyond a given voltage, the resistivity of the chalcogenide material decreases with temperature increase. The amorphous resistance can be derived with or without snap back (In the reset state, the device voltage increases non-linearly with current to the threshold voltage, then suddenly drops). All these modules are unified thanks to smooth functions. [90]. To achieve the model, some assumptions were taken into account: the threshold voltage is assumed constant (for a given ambient temperature after the drift end [91]) and there is no partial set and reset state. Moreover, two access resistances are taken into account. The model has been compared to measured data from [73] and the results are displayed in figure 39.

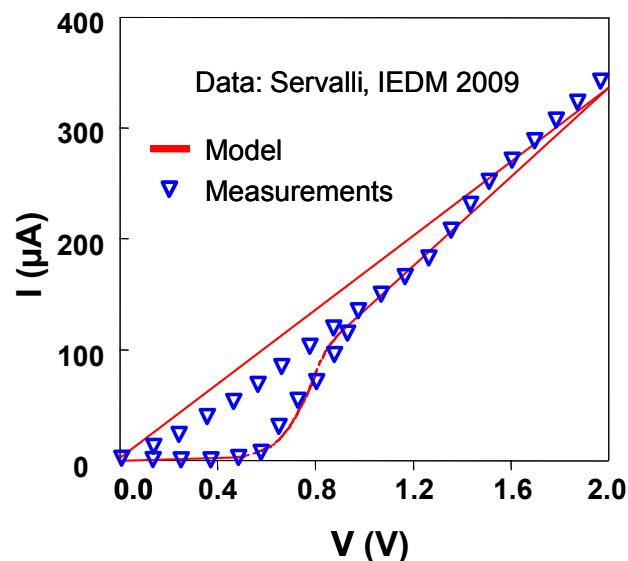


Figure 39. Comparison between the model and measurements of I-V characteristics of the amorphous and crystalline phases of a PCM cell.

### 3.2.5.3 Transient Simulations

Transient simulations have been performed in order to validate the global behavior of the structures. In addition, they also allow the validation of the compact model in circuit application contexts.

Figure 40 shows the electrical results of a typical memory use case. Here,  $V_{dd}$  is equal to 500mV and  $V_{ss}$  is tied to 0V. In region A, the elementary node is initially configured to apply a logic '1' signal to the output node. At that point, the PCMs are configured according to figure 29. The memory node is then reprogrammed in the region B, by sequentially applying a RESET pulse to the PCM that is connected to  $V_{dd}$  and a SET pulse to the other PCM, while the programming transistor is selected by setting *Prog\_Enable* high. This leads the structure

to flip its memory content. In the final read operation C, logic '0' is sensed at the output node.

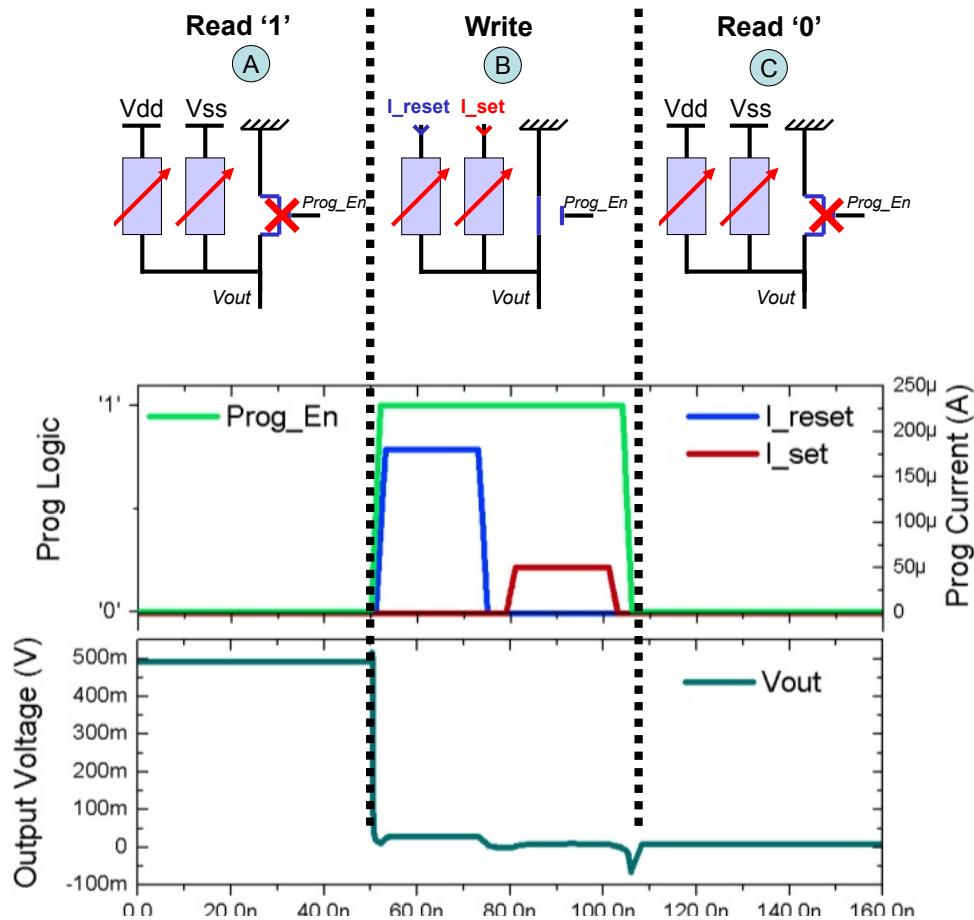


Figure 40. Binary memory node transient simulation

Figure 41 describes the simulation of a unique routing device. This routing device is built by a single resistive memory. Firstly, the resistance of the memory is high, leading to open switch behaviour. In region A, we observe that  $V_{out}$  does not follow the  $V_{in}$  signal. In fact,  $V_{out}$  is slightly impacted by  $V_{in}$ , due to the small leakage current flowing through the memory in its high resistive state. In region B, the memory is reprogrammed. Thus, the Programming structures are activated, in order to validate all the outbound programming paths from the switchbox. Then, the programming signal is applied to the resistance, in order to change its resistance. By applying the  $I_{reset}$  signal, the memory is expected to be in its low resistance state. This is confirmed in region C. Indeed, it is now clear that the  $V_{out}$  signal follows the  $V_{in}$  signal. The small impact on the signal is due to the RC circuit between the output node and the memory on resistance. This signal deformation is corrected by the output buffers that are placed at the switchbox outputs.

#### 3.2.5.4 Performance Estimation

Table IV shows some characterization results in terms of area and write time for the proposed solution compared to traditional FPGA memory nodes. Considering that all these elements are driving an equal load, we omitted the pass transistor at the output node in our simulations. Thus, the SRAM cell is considered to be a 5-transistor (5T) structure. The Flash topology is implemented by 2 Flash transistors [35]. The Magnetic RAM implementation is realized by an unbalanced flip flop as proposed in [37]. In fact, this allows efficient separation of the programming path from the data path.

We see that the proposed PCM cell is the most compact solution, even with the impact of the programming current on the access transistor. This advantage is due to the reduction of the memory front-end footprint to only one transistor, compared to 5 for the SRAM cell, and compared to 2 for the Flash solution (one pull-up transistor coupled to a floating gate transistor). We should remark that PCMs offer a significant reduction in writing time for non-

volatile memory technologies, as well as a consequent reduction in writing energy. In our context, it is possible to reduce the area by a factor of 1.5 and the writing time by a factor of 16.6 and the programming energy by a factor of 500 compared to an equivalent flash technology. However, the programming energy is about 10 times larger than in the flash technology. Compared to another emerging equivalent non-volatile magnetic resistive technology, we can observe that the proposed structure improves the area by a factor of 3.8. This is due to the chosen structure for the magnetic resistive node, which is based on an area-hungry unbalanced flip flop. Nevertheless, it is important to highlight that the PCM-based node is slightly slower for reconfiguration, with a 33% difference as compared to magnetic memory. Finally, while the writing time is slower, the required programming energy is reduced by a factor of 50 as compared to the magnetic technology, which requires energy to create the programming magnetic field. We should therefore consider that this result is strongly dependent on the technology used. In this work, we considered Thermally-Assisted-Switching (TAS) MRAM technology. While this technology has the same level of maturity as PCM, it also requires a large programming energy (between 100pJ and 150pJ per cell). Most advanced writing schemes, such as Spin Transfer Torque [92, 93] allow further reduction of writing energy to the range of 3pJ per cell.

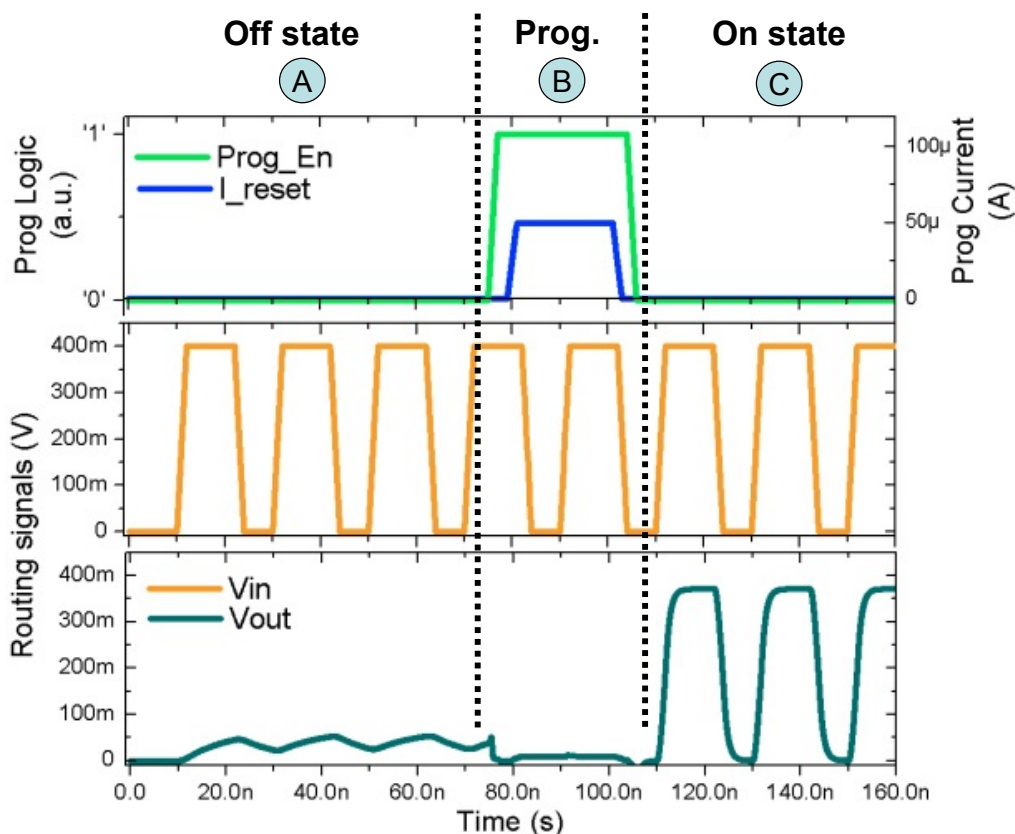


Figure 41. Elementary routing element transient validation simulation

Table IV. Detailed technology performance evaluation

	Cell elements	Area (F <sup>2</sup> )	Write time (ns) [6]	Programming energy (pJ) [6]
<i>SRAM</i>	5T	115	0.3	$7 \cdot 10^{-4}$
<i>Flash cell</i>	2T	46	1000	0.1
<i>MRAM UFF</i>	5T2R	115	40	300
<i>PCM cell</i>	1T2R	30	60	12
<i>Flash vs. PCM</i>	-	x 1.5	x 16.6	x 0.08
<i>MRAM vs. PCM</i>	-	x 3.8	x 0.6	x 25

Table V shows some characterization results in terms of area and writing time for the solution

and conventional FPGA memory nodes. The SRAM-based node is formed by four 4-input multiplexers, programmed by 5T SRAMs. The Flash-based solution is taken from [35], while the MRAM-based solution is taken from [37]. In this work, 2 flash transistors are used to replace a combination of one SRAM and one pass gate. The structure of the switchbox is considered to be the envisaged PCM structure. We should note that a GeTe technology [116] with a very low on-resistance value is envisaged, as it corresponds to the requirements for a high performance switch.

We see that the proposed PCM switchbox is the most compact solution, even with the impact of the programming current on the access transistor. This advantage is due to the reduction of the memory front-end footprint to only metal lines, compared to 5 transistors for the SRAM cell and 2 transistors for the Flash solution (one pull-up transistor coupled to a floating gate transistor). It is also worth noticing that PCMs offer a significant reduction in write time versus flash technology. In our context, it is possible to reduce the area by 3.4 and the write time by 16.6 as compared to an equivalent Flash technology. Another important metric is the data path resistance. In the context of a logic signal propagating through the switchbox, we observe that the on-resistance of PCM is very low compared to the other solutions. In fact, all the other solutions still use a MOS transistor to realize the data path switch. Hence, since the on-resistance is decreased by a factor of 182, this means that a reduction of the propagation time is expected. From the point of view of the off-resistance however, we observe the opposite, since the off-resistance follows the same trend as the on-resistance. This means that the off-resistance decreases by a factor of 144, and since off-resistance impacts directly on switch leakage, this decrease is likely to worsen static power dissipation and logic levels. However, this could be mitigated by the fact that the  $R_{on}/R_{off}$  ratio has not decreased - on the contrary, PCM cells demonstrate an improved ratio with a value of around 1/20000, with respect to the value of 1/15824 obtained for the other solutions. The ideal solution will require a low  $R_{on}$  value, while the  $R_{on}/R_{off}$  ratio should be increased with respect to CMOS.

**Table V. Technology performance evaluation (2x2 switchbox)**

	<b>Cell elements</b>	<b>Area (F<sup>2</sup>)</b>	<b>Write time (ns) [6]</b>	<b>On data path resistance (<math>\Omega</math>)</b>	<b>Off data path resistance (M<math>\Omega</math>)</b>
<i>SRAM</i>	72T	2576	2.4	9100	144
<i>Flash cell</i>	48T	1104	24000	9100	144
<i>PCM cell</i>	8T24R	321	1440	50	1
<i>MRAM UFF</i>	168T48R	4448	960	9100	144
<i>Flash vs. PCM</i>	-	x 3.4	x 16.6	x 182	x 0.007
<i>MRAM vs. PCM</i>	-	x 14	x 0.6	x 182	x 0.007

### 3.2.5.5 Discussion

The presented performance metrics make the solutions built around resistive memories of high interesting for the purposes of reconfigurable applications. Indeed, we showed that it is possible to create a compact configuration memory node that can store a logic level in two resistances. This logic node improves the size as compared to flash memories by 1.5x. While this 33% reduction in area reduction is significant, it is worth noticing that the limitation is due to the programming transistor. Indeed, the selection transistor size is dictated by the level of current to be driven to the cell being programmed. While the required programming current is large, the programming energy remains the lowest over the benched technologies. Other technologies might be envisaged in order to further reduce the programming current. First of all, other PCM technologies, such as GeTe or GeTeC [116], require smaller currents for programming compared to standard GST. It is then possible to migrate to other resistive memory technologies. For example, OxRAM requires less current overall for their programming, and will lead to more compact configuration nodes [117].

Concerning the routing part, we propose to introduce the PCM directly into the logic data path. This is of interest from two points of view. First, we obtain an area reduction of a factor of 3.4 as compared to flash. However, the most interesting advantage comes from the technology: the ReRAM technology is shown to have a much lower on-resistance than CMOS switches, which should lead to significant reductions in propagation delay in complex logic circuits. While the well-known GeSbTe-based PCM technology has a “quite high”  $R_{on}$  value, it is attractive to look at new PCM materials or new technologies. In particular, special alloys such as GeTeC and GeTe [116] have demonstrated a very small on-value, well-suited to the requirements of the routing structures. Nevertheless, it will be necessary also to study the reliability of the proposed switchboxes as well as the data retention time in the structure. Resistive memories are switched by a controlled current or voltage applied to or through them. Nevertheless, the switching behavior is generally more complex, as far as timing and environment (e.g. temperature) must also be considered. Unpredictable signals are flowing in the data path. Unpredictable means that the signals depend on the application. Thus, the memories introduced in the flow maybe switched to unwanted states leading to an unreliable structure.

### 3.3 Proposal 2: Monolithic 3-D Integration Process

In the previous section, we looked at a technology, which allows a passive resistive memory device to be embedded into the back-end levels. While this approach is advantageous for memory or routing structures, the elements are only passive. It then appears interesting to assess the interest of 3-D integration, in order to diversify the functionality of above IC devices. In this section, we propose the use of a monolithic 3-D integration technology to stack two active layers. This represents a first step towards the use of active devices in a 3-D scheme.

#### 3.3.1 Introduction

For several decades, the semiconductor industry has invented new approaches to increase integration density and transistor performance, the main vector for this being MOSFET scaling.

In this context, the IC integration in three dimensions appears to be a promising alternative path to scaling, and to some extent would avoid the huge investments required by scaling. While the concept is not entirely new [94], the development of the technology has witnessed significant growth over the last decade. In particular, the technology of *Through-Silicon-Vias* (TSVs) is currently the reference for 3-D technology processes. A TSV could be defined as a large via built across the substrate, in order to contact the active front-end to the reverse-side of the chip. Bumps are then used to contact the dies between them, as shown in figure 42.

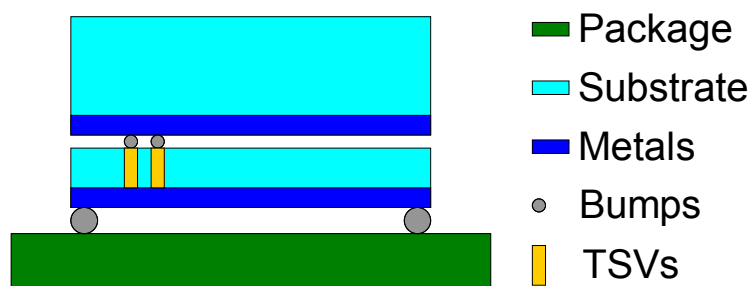


Figure 42. TSVs-based die stacking cross-sectional view

While such integration is challenging due to TSV density requirements (leading to a double specification of TSV aspect ratio and wafer thickness), the process is mature enough for industrial applications [95]. However, there are still some hurdles to face. Table VI depicts the principal characteristics of TSV processes [96]. It is worth noticing that TSVs are quite area-hungry with diameter up to  $5\mu\text{m}$  and pitch up to  $10\mu\text{m}$ . This means that the connection density is quite poor (from 400 to 10000 TSVs/ $\text{mm}^2$ ) and this results in a loss of active size.

Thus, designers are limited to high level interconnect, such as memory/core communication. However, this segregation is of great interest for performance, since processes will be tuned to optimize each layer to a given class of application (low power, general purpose ...)

Table VI. Density survey of TSV technologies [96]

	Diameter ( $\mu\text{m}$ )	Pitch ( $\mu\text{m}$ )	Density (TSVs/ $\text{mm}^2$ )	CMOS 65-nm gate equivalent
Low Density TSVs	20	50	400	1250
High Density TSVs	5	10	10000	50

In a reconfigurable application, a large number of interconnections are required if separation between memory and logic is to be envisaged. This means that other integration processes should be used, to overcome the limitations of connection density. In an FPGA, connections between memory and logic are done at gate level. We estimate the required density at about 500 000 3-D contacts/ $\text{mm}^2$ . Thus, instead of processing the layers separately and stacking them *a posteriori*, it is possible to use monolithic sequential integration. In a monolithic integration, the circuit is processed from the bottom to the top. This means that the stacked layers are build by a set of technological steps above the already processed stack, as illustrated in figure 43.

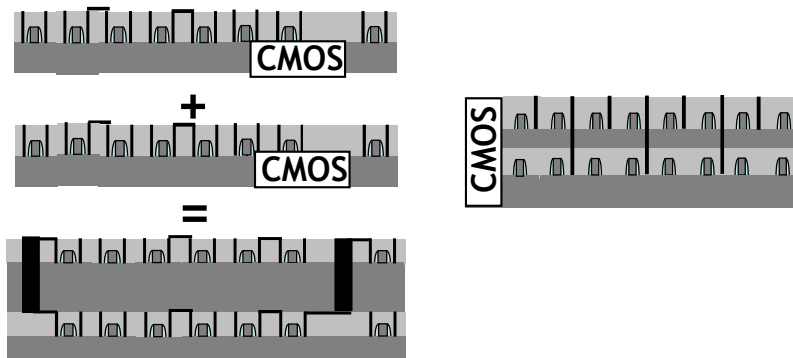
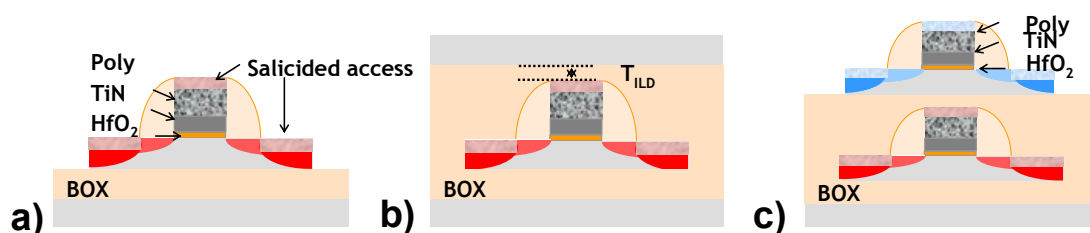


Figure 43. Cross-sectional view of 3-D parallel integration (left) and 3-D monolithic sequential integration (right)

Such an integration scheme is promising for several reasons. Firstly, it enables a much higher via integration density. In fact, the contacts use a planar scheme step, as opposed to TSVs. Thus, it is possible to obtain alignment accuracy in the range of lithographical precision (around 10-nm), whereas TSV alignment accuracy is in the range of  $1\mu\text{m}$ . However, monolithic integration gives rise to other issues. In particular, conversely to parallel integration, where the layers are processed separately, it is necessary for the fabrication of the top transistors not to degrade the performance of the bottom layer transistors. This means that the top layer thermal budget is limited and that processes have to be thought through in consequence.

### 3.3.2 Technological Assumptions

As briefly introduced above, the monolithic 3-D integration process integrates at least two layers of active silicon. These layers are processed sequentially, as shown in figure 44. The principal process steps are (a) the realization of the bottom transistor, (b) the deposit of the top film silicon and (c) the realization of the top transistor. Several constraints appear at each step.





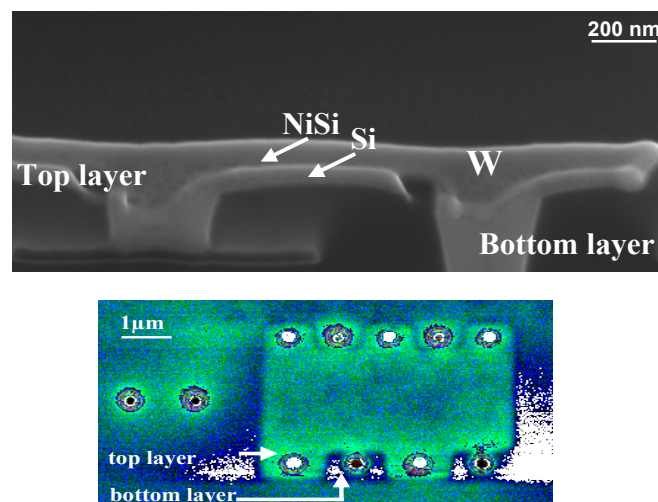
**Figure 44.** Cross-sectional view of 3-D monolithic steps – a) optimized bottom FDSOI process b) high quality top film deposition c) low temperature top FDSOI process

The first transistor layer must be optimized in order to improve its robustness to temperature. It is worth noting that the second active transistor layer will have an impact on the first layer and thus the thermal robustness of the bottom layer is critical. The maximal thermal budget for subsequent steps is limited by salicidation. In fact, standard NiSi dewetting occurs in only 3 minutes at 600°C. To obtain a low resistance access with processes around 600°C, a salicide stabilization procedure has been demonstrated in [97]. This salicidation approach is based on an optimized W-NiPtSi-F. With this technology, the salicide is stable at 600°C, which is considered to be the maximal thermal budget for all the subsequent steps.

The second step corresponds to the realization of a high quality top film. While seed window techniques were first used in the literature [98], they lead to crystalline defects, poor thickness control and density loss. It is thus better to use molecular bonding [97], where a blanket Silicon-On-Insulator wafer is transferred on top of the processed MOS wafer. This step plays a major role in alignment accuracy, since the alignment occurs after bonding, conversely to parallel integration [99].

A low temperature Fully Depleted Silicon-On-Insulator is then processed on the top film. Due to the limited thermal budget (600°C), it is not possible to perform thermal activation of dopants (around 1000°C-1100°C). The process uses a solid phase epitaxial re-growth [100], which consists of pre-amorphization of the top film, followed by dopant implantation and finishing by a re-crystallization at 600°C.

Finally, 3-D contacts are realized using the same contact techniques as in standard processes. Only one lithography step is required for all contacts. Figure 45 illustrates the contact realization, where W (tungsten) plugs exist between the top and bottom layers. The process has been validated for several simple cases, such as SRAM memories or inverter cells [100, 101].



**Figure 45.** Contacts between metal layer and top/bottom layers [97]

### 3.3.3 Elementary Blocks

In order to evaluate the opportunity of the monolithic 3-D integration process, elementary nodes for configurable logic have been designed. These elementary blocks have been implemented in the test chip SNOW1<sup>4</sup>, and are currently under fabrication. The blocks that have been chosen are a top memory, a 3-D Look-Up table and a 3-D Cross point. The global idea is to place the configuration memory just above the circuit that requires it.

<sup>4</sup> Silicon Nanoelectronics On 300-mm Wafer (SNOW1) is a CEA LETI's internal mask set sent to mask shop on 09/2010. This mask set is multi-project and multi-technologies. Addresses technologies are: FDSOI, Ultra-Thin BOX FDSOI, 3-D monolithic, 3-D TSV, 3-D NEMS, OxRAM, Flash and implemented circuits range from simple technological test circuits to complete design.

### 3.3.3.1 3-D Memory

A configuration memory as used traditionally in reconfigurable systems is depicted in figure 46. Figure 46-a shows the implemented schematic, which is built around a latch. The latch is programmed by applying a strong signal (i.e. strengthen driven current) to one branch of the memory. The strength of the signal depends on the latch size and must be stronger compared to that generated by the internal latch feedback inverter. The programming signal is applied by the inverter connected to the SRAM. In order to ensure its sizing, the dimensions are increased by three to four times the size of the feedback inverter of the latch. It is also remarkable that the structure uses two cascaded tri-state inverters before the latch, to allow electrical insulation from the input node and consequently memorization. The third state of the inverters is controlled by an inverted clocking scheme, which makes this architecture edge triggered. Figure 46-b presents the layout, which uses regular drawing techniques. Conversely to standard cell design, regular layout technique means that polysilicon lines must be drawn with a constant pitch. This means that dummy polysilicon lines are added to match this requirement. In addition, contacts and active regions must be aligned and also pitched regularly. This high level of regularity in the layout, particularly in critical layers (i.e. the layers defining the front-end devices), is expected to improve the reliability of the technology. In this layout, it is worth pointing out that most of the area is used for the programming tri-state buffer.

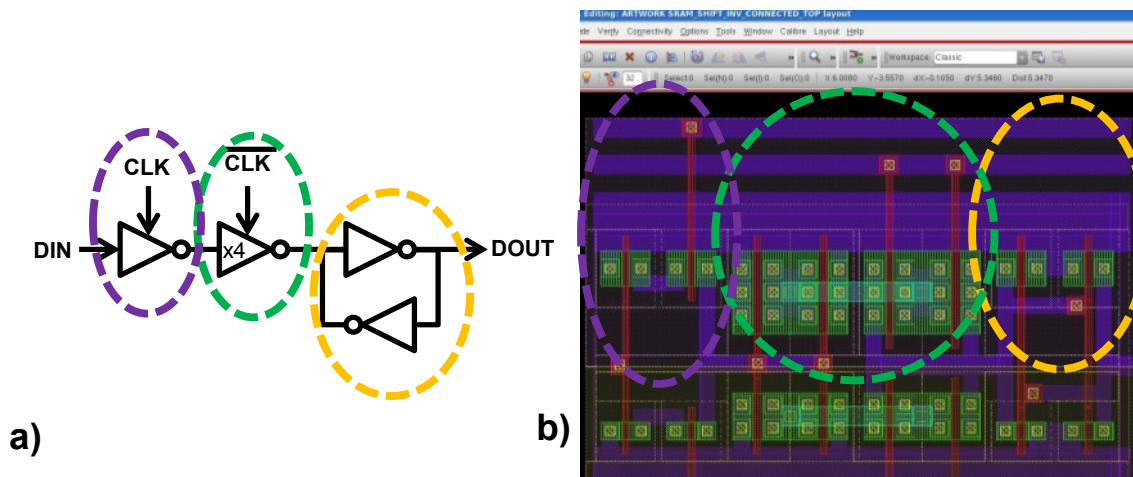


Figure 46. 3-D memory for reconfigurable logic – a) schematic b) layout on top levels

### 3.3.3.2 LUT Impact

Since we are able to achieve configuration memories implemented in the top layer, it is quite straightforward to develop a two layer 3-D look-up table. In figure 47, the layout and the schematic of a 1-bit Look-Up Table are depicted. While the structure is too simple for any architectural considerations, it is useful for the purposes of technological demonstration.

Two configuration nodes are thus placed above a multiplexer circuit. The dimensions of the multiplexer are largely relaxed, due to limitations on the metal layers available for the technological demonstrator.

In this way, we could consider that only a small amount of size is used for the structure, in addition to the multiplexer. Furthermore, while we observe a strict separation between the configuration memory part and the data path multiplexer, it is possible to optimize the technological process to reach specific properties. Especially, it is envisaged to have a low-leakage – low-power process for memories, while the active data paths are created using high performance transistors. This makes the implementation highly promising for circuits requiring large configuration memory.

### 3.3.3.3 Routing Structure Impact

Using the same technological assumptions, we propose a routing cross point. The cross point is built around a 2-transistor pass-gate, which is driven by a static memory, as illustrated in



figure 48. The memories and data path will be split and placed on different levels. Memories are placed on top of the cells that directly use their information, and thus a large reduction in the number of wires is observed. Two different structures have been realized.

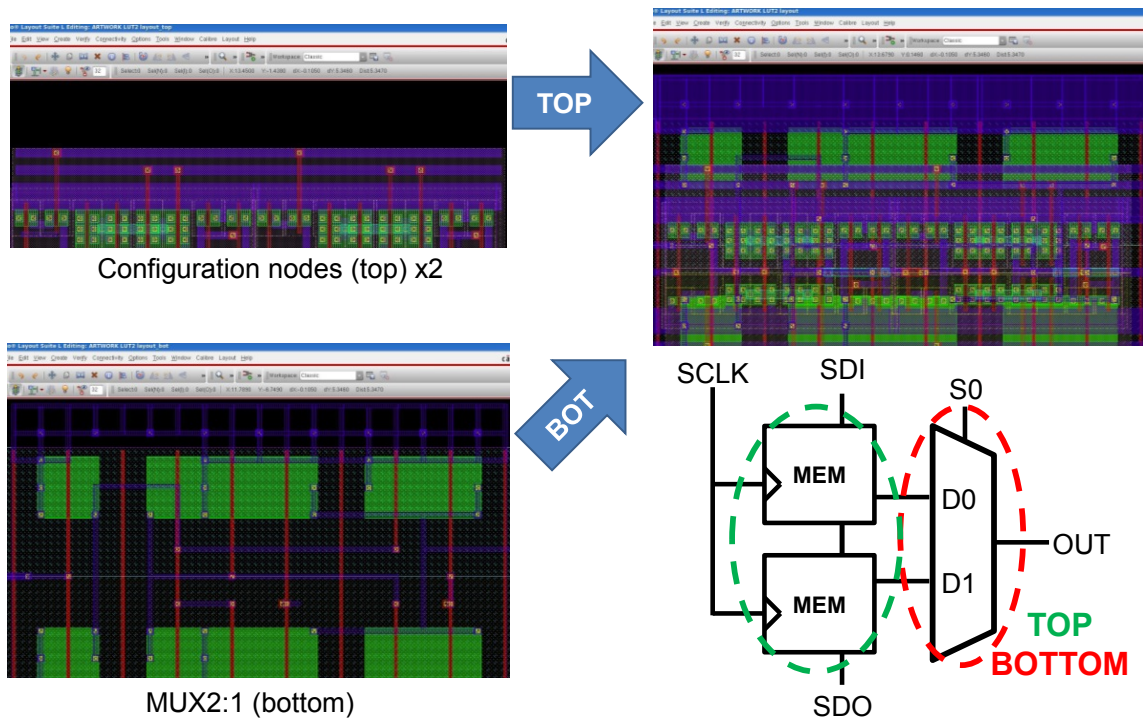


Figure 47. Two layer monolithic 3-D based 1bit Look-Up Table – layout and schematic

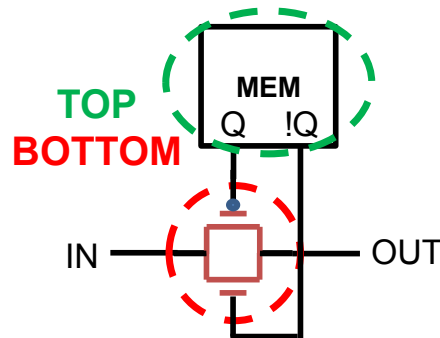


Figure 48. Cross point schematic

Figure 49 shows the cross point using the configuration node presented above. It is worth noticing that, in this circuit and conversely to the LUT, the memory dominates the dimensions of the cell. Indeed, the pass-gate, even when considering a large size transistor, has low requirements in terms of the metal layers, and thus its implementation below the memory is kept at the front-end. We also point out the H-shape of the pass-gate, due to the low number of metal layers in the technological demonstrator. This H-shape is not necessary when the number of available metal layers is greater than 2.

The above implementation uses a configuration node, which could be programmed serially. This is obviously of high interest for FPGAs, since it allows memory cascading and thus simplifies the programming circuitry. Nevertheless, in the context of separation between logic and configuration memory in a 3-D process, it is of interest to use a standalone memory approach. Hence, a typical SRAM organization could be envisaged on the top layer with word lines and bit lines crossing through the circuit, while the FPGA data path is found below. Figure 50 shows a possible implementation of an SRAM memory with a pass-gate in a very compact 3-D implementation.

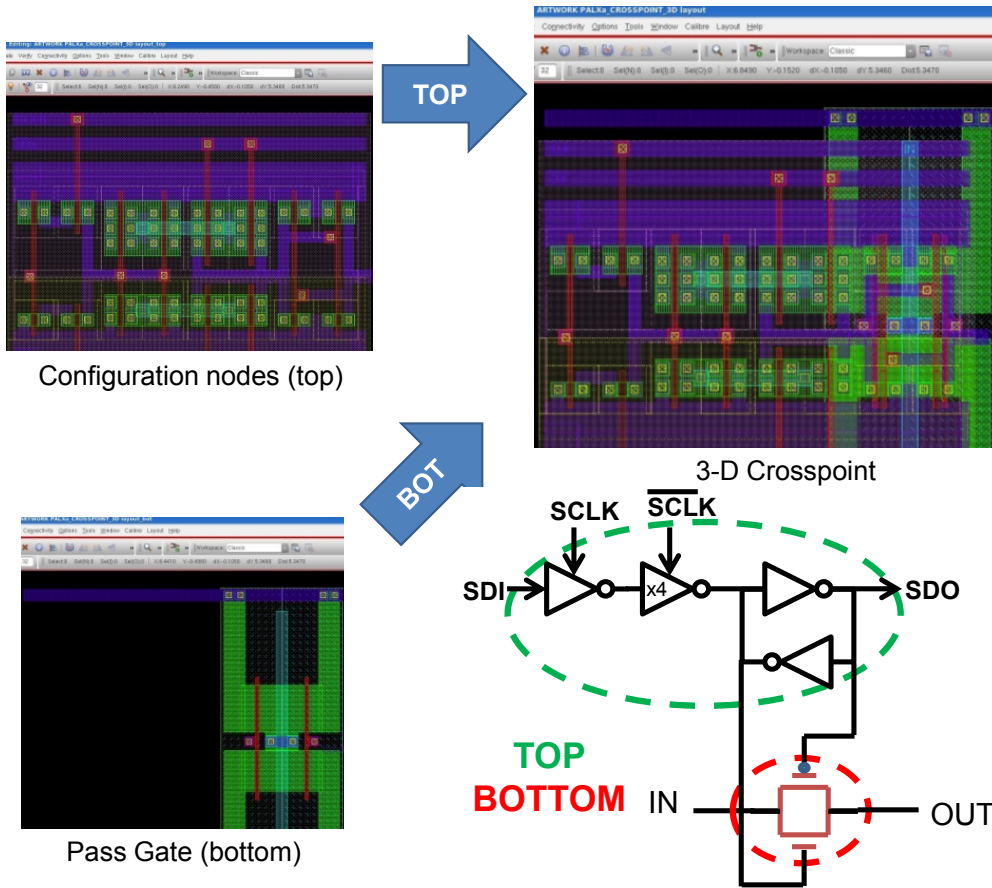


Figure 49. Two layer Configurable cross points (Configuration node based) – layout and schematic

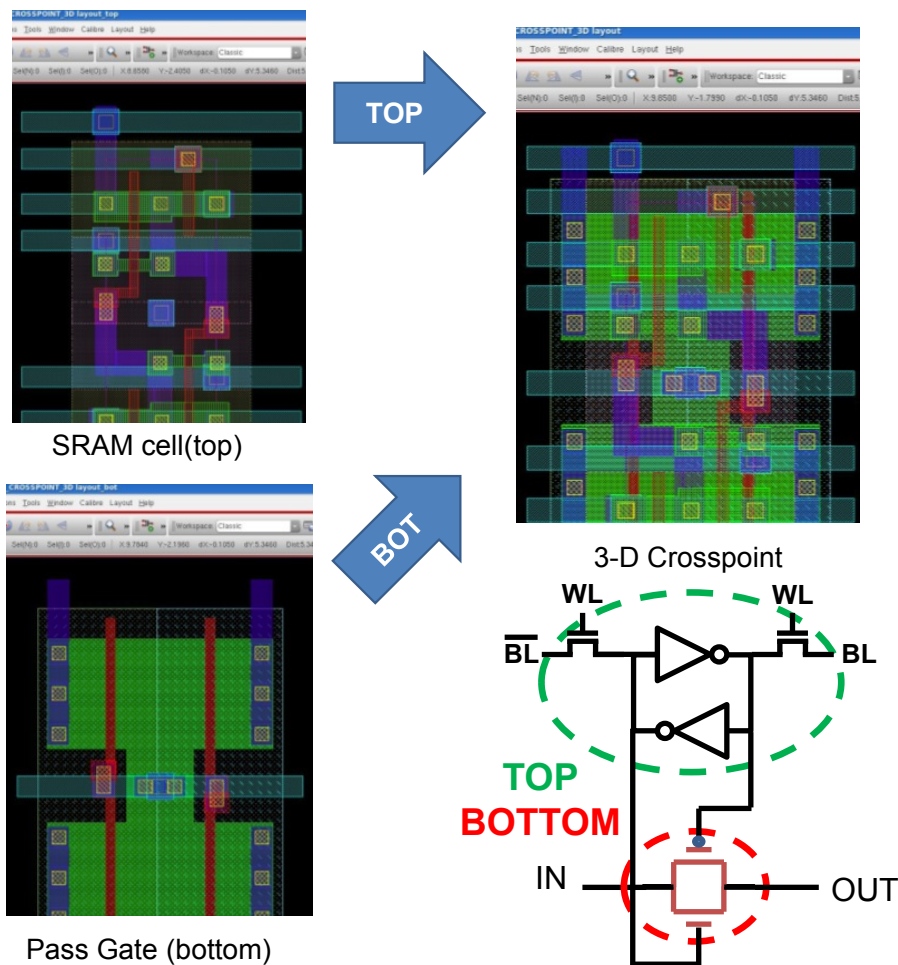


Figure 50. Two layer Configurable cross points (SRAM cell based) – layout and schematic

### 3.3.4 Performance Characterization

#### 3.3.4.1 Methodology

The presented elementary blocks have been designed on a real silicon test chip. At the time of writing, the test chip is in fabrication, so obviously silicon measurements are not available. We will thus only consider metrics extracted from layout and electrical simulations.

The simulations are realized using the Synopsis HSPICE electrical simulator. This simulator allows a fast and accurate simulation for the LETI's homemade model.

The metrics will be area, delay and power consumption. The area is extracted considering the layouts. The test chip uses two different set of design rules for front-end and back-end. In order to demonstrate the capabilities of the technology with a relatively low cost, the front-end can be scaled down to the 16-nm node, while the back-end uses a 65-nm lithographic node. Nevertheless, in order to be compared fairly with existing technology, all the dimensions used in the layouts use an equivalent 65-nm lithographic node. The delay is assessed by electrical simulation and broken down into the intrinsic delay and the load delay (expressed as the  $K_{load}$  factor). The power consumption is extracted by electrical simulation with an FO4 load. The circuit is operated at the gate's maximum achievable frequency with the FO4 load. This means that rise and fall time are chosen in accordance to  $5\tau$  charge and discharge. The related frequency obviously depends on the circuit. At this frequency, we swept all possible input vector combinations and we averaged the power consumption.

In order to compare the performance of the 3-D FDSOI, we compare all metrics to a standard industrial low power 65-nm process. In order to differentiate the gain due to the FDSOI technology and that due to the 3-D implementation, we also perform the evaluation of the equivalent circuit in 2-D FDSOI. We should mention that the electrical model is the same for both bottom and top transistors and that parasitic post-layout extraction is not available in the design kit as used. Thus, comparisons in terms of performance and power will not differ from 2-D and 3-D implementation of the FDSOI cells.

#### 3.3.4.2 SRAM Configuration Performance

Table VII shows the performance evaluation for the configuration node. In this evaluation, we highlight that the name of 3-D FDSOI is misleading. In fact, since we expect to use the configuration memory node with other circuits, we are benching only the 2-D implementation. Nevertheless, the so-called 3-D comes from the fact that this implementation is done on the top silicon.

Table VII. Evaluation of SRAM Shifter performance

SRAM Shifter 65-nm node	Area ( $\mu\text{m}^2$ )	Intrinsic delay (ps)	$K_{Load}$ ( $\text{ps}\cdot\text{fF}^{-1}$ )	Static average power (nW)
<i>2-D LP Bulk</i>	13.52	18.95	6.06	6.34
<i>2-D FDSOI (bottom)</i>	10.64	6.30	2.70	1.19
<i>3-D FDSOI (top)</i>	10.97	6.30	2.70	1.19
<i>2-D Bulk vs. 3-D FDSOI</i>	x1.80	x 3.00	x 2.24	x 5.33
<i>2-D FDSOI vs. 3-D FDSOI</i>	x 1.03	-	-	-

The general conclusion that we can draw is that the technology improves the performance for all the metrics compared to a 2-D low power bulk process. Area is improved by 1.8x with respect to 2-D bulk. While the lithographic node remains the same, the FDSOI technology leads to smaller transistors for the same performance. Furthermore, the SOI transistors do not require any bulk patterning, such as guard rings. Comparisons between the top and bottom implementations of the same circuit show that the area is slightly larger for the top

implementation. This is due to the lithography rules, which are slightly relaxed for top layers in comparison to the bottom layers. In terms of electrical performance, we observe a gain of 3x in intrinsic delay, 2.24x in the load factor and 5.33x in static power. This is a direct result coming from the FDSOI technology. In fact, the technology is especially suited for low power operations, while maintaining a high performance quality.

#### 3.3.4.3 LUT Performance

Table VIII shows the performance comparisons for a very simple 2-input LUT test case. The considered circuit corresponds to that presented in figure 47. Concerning the area, we obtain a 2.03x improvement in regards to 2-D bulk. This figure is mainly due to the stacked integration of the memory on top of the multiplexer. It is also worth mentioning that the 2-D FDSOI implementation is larger than the 2-D bulk one. This result is counter-intuitive, since the lithographic node used is the same in both cases. However, we should note that the 2-D FDSOI layout has been carried out following regular layout techniques, as well as some relaxed technological demonstrator rules. This clearly leads to larger cell implementations than in an equivalent non-regular bulk process.

In terms of performance and power, an improvement of 1.62x in intrinsic delay can be observed, as can figures of 6.11x in load influence factor and 2x in dynamic power at 1GHz. The contributing reasons for these good numbers are twofold. As already stated, the technology is compliant to low-power high-performance circuits. Nevertheless, the improvements are also due to 3-D integration. It is possible to group the low-power optimized circuits on one layer, while performance-optimized blocks are grouped on another. This allows specific process optimization, but it also allows the dimensions of some transistors to be relaxed. For example, in this context of LUT, it is possible to size the multiplexer to be quite large, while the memories are placed on top. This strategy leads to a high-performance multiplexer, placed under the necessarily area-hungry memories.

**Table VIII. Evaluation of Look-Up Table performance**

<b>LUT2 65-nm node</b>	<b>Area (<math>\mu\text{m}^2</math>)</b>	<b>Intrinsic delay (ps)</b>	<b><math>K_{\text{Load}}</math> (<math>\text{ps}\cdot\text{fF}^{-1}</math>)</b>	<b>Average power at 1 GHz (<math>\mu\text{W}</math>)</b>
<i>2-D LP Bulk</i>	64.48	85.08	36.29	46.90
<i>2-D FDSOI</i>	122.57	52.52	5.94	23.23
<i>3-D FDSOI</i>	31.7	52.52	5.94	23.23
<i>2-D Bulk vs. 3-D FDSOI</i>	x 2.03	x 1.62	x 6.11	x 2.02
<i>2-D FDSOI vs. 3-D FDSOI</i>	x 3.87	-	-	-

#### 3.3.4.4 Cross point Performance

Table IX shows the performance evaluation of 3-D cross points. As previously described, two different cross point schemes are studied. The difference between the two cross points resides in the configuration. One is based on a shifter configuration node and the other is based on an SRAM.

In terms of area, we should note that the SRAM-based cell is the most compact one, with an improvement factor of 2.93x compared to 2-D bulk. This is obtained mainly due to the compactness of SRAM cells, where each transistor and layout is optimized to obtain the best density. Nonetheless, due to its greater size, the configuration-node based cross point improves the area by “only” 1.48x. It is worth noticing that this cross point implementation embeds all the circuitry required to program the node in a shift register manner. Thus, it appears obvious that this implementation is larger compared to the fully optimized SRAM-based implementation.

In terms of performance, the intrinsic delay is shortened by 1.6x and 3.1x for the configuration-node based implementation and the SRAM-based implementation, respectively.

The load factor is almost the same with a gain of 1.06x. The difference comes mainly from the pass-gate cell. It is in fact possible to use larger transistors in the SRAM-based implementation; indeed, since the cell is more compact, the shape of its layout facilitates its placement and efficient connection of the bottom pass-gate.

Finally, in terms of dynamic power, we observe a 3.1x improvement at 2GHz for the configuration-node based cell, and 2.1x for the SRAM-cell, compared to bulk. These numbers come mainly from the use of FDSOI, which is an intrinsic low power technology. The difference between the two memory structures is due to pass-gate sizing. Since most of the contribution to power originates in the data path, these results can again be attributed to the sizing of the pass-gate. In the SRAM-based case, the transistor is larger and leads to larger power consumption during the signal drive.

**Table IX. Evaluation of Cross point performance**

<b>65-nm node</b>	<b>Area (<math>\mu\text{m}^2</math>)</b>	<b>Intrinsic delay (ps)</b>	<b><math>K_{\text{Load}}</math> (<math>\text{ps}\cdot\text{fF}^{-1}</math>)</b>	<b>Average power at 2 GHz (<math>\mu\text{W}</math>)</b>
<i>2-D LP Bulk</i>	15.83	0.44	0.71	0.12
<i>3-D FDSOI Configuration node</i>	10.66	0.28	0.67	0.039
<i>3-D FDSOI SRAM cell</i>	5.40	0.14	0.66	0.058
<i>2-D Bulk vs. 3-D FDSOI Conf. node</i>	x 1.48	x 1.6	x 1.06	x 3.1
<i>2-D Bulk vs. 3-D FDSOI SRAM cell</i>	x 2.93	x 3.1	x 1.07	x 2.1

#### 3.3.4.5 Comments

The described performance levels have shown a significant advantage of this technology for reconfigurable circuits. In general, the FDSOI technology is highly suitable for power reduction in electronic circuits. In particular, it improves the FET properties, which leads to reduction in leakage current, while the performance is increased. Nevertheless, the ability to stack several layers with a high alignment accuracy and high via density also allows to achieve high improvements in terms of area. In particular, we have shown a reduction in area by a factor of 2x for a 2-bit LUT test case circuit. The remaining area occupation is composed essentially of the bottom multiplexer. Furthermore, we can observe that the implemented circuits are following very strict (conservative) design rules, which are required by the technologist at the current level of technological maturity. In the future, it is clear that rules will be more aggressive, and will lead to further improvements in the figures.

### 3.4 Proposal 3: Vertical Silicon Nanowire FET Process

In the previous part of this chapter, we have investigated various technologies, which allow respectively a passive resistive element to be embedded in 3-D and to stack several layers of active devices. Nevertheless, these solutions should not really be considered “true” 3-D solutions since they cannot place active devices in a real 3-D shape, i.e. with a vertical rather than planar transistor. In this section, we will assess such a technology, which is able to build vertical FETs in the Back-End metallic layers.

#### 3.4.1 Introduction

As previously discussed, the 3-D integration of transistors is an attractive solution to pursue the increase of circuit performance, while limiting the cost, as opposed to the continued single use of scaling. Stacking technologies, whether the traditional sequential technology (where wafers are processed separately and then stacked) or advanced monolithic integration (where transistors are processed step by step on the same wafer) only deal with stacks of planar

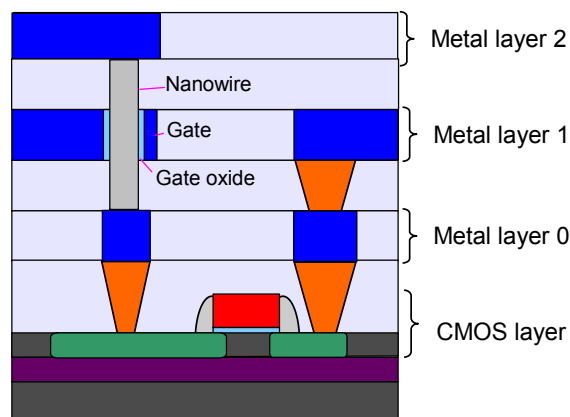


transistors. The transistor itself (or more specifically, the transistor channel) does not exploit the vertical direction in these approaches.

Meanwhile, semiconducting nanowires have recently attracted considerable attention. To further miniaturize the transistor while still maintaining control over power consumption, alternative transistor geometries have been considered [106]. With their unique electrical and optical properties, they offer interesting perspectives for basic research as well as for technology. A variety of technical applications, such as nanowires as parts of sensors [109], and electronic [108, 171] and photonic devices [111] have already been demonstrated. In particular, electronic applications are increasingly coming into focus, as ongoing miniaturization in microelectronics demands new innovative solutions. Typically, silicon nanowire transistors have a horizontal, planar layout with either top or back gate geometry [110]. However, the amount of energy and time required to align and integrate these nanowire components into high-density planar circuits remains a significant hurdle for widespread application. More advanced works show a *Gate-All-Around* (GAA) organization in a planar topology [108]. In-place growth of vertically aligned nanowires, on the other hand, would in principle significantly reduce the processing and assembly costs of nanowire-based device fabrication, while opening up opportunities for "true" 3-D. Some research works have demonstrated the possibility of fabricating transistors directly between two metal lines, within the back-end levels [112, 113]. These works make it possible to realize computation directly in the metal levels, through programmable vias, as well as potentially complete complementary logic functions.

### 3.4.2 Technological Assumptions

Recent studies have demonstrated the possibility to grow single crystalline silicon nanowires on a metallic line, into a CMOS compatible process [129]. This work represents a great opportunity to build FET devices in the interconnect levels [112, 113]. We propose to co-integrate standard CMOS with vertical Nanowires Field Effect Transistors. The cross-sectional view is shown in figure 51.



**Figure 51. Cross sectional schematic showing a BEOL FET and standard CMOS FET co-integration**

First of all, standard transistors are processed using the specified technology, which could be very versatile, such as Bulk, Silicon-on-Insulator, Fully Depleted SOI, Thin Box FDSOI, among others. Then, silicon nanowires can be grown in a CVD reactor using the VLS mechanism. Even on a metallic line, they have a single crystalline structure and semiconducting properties. Taking advantage of this, and respecting low temperature processes under 400°C, it is possible to make vertical transistors between two interconnecting lines. After etching a hole through the oxide to the metallic bottom line, a catalyst can be deposited at the bottom. Nanowires can be grown from the metallic line using the oxide hole as template and a deposited metallic catalyst. Using diborane or phosphine, nanowire can be doped to form P-N junctions. Nanowires for p-MOS and n-MOS should be grown during two distinct sequences comprising template formation and growth. After growth, a chemical etching can be used to remove a part of the oxide template. A multilayer gate stack can be achieved thanks to ALD and CVD deposit of the dielectric ( $\text{Al}_2\text{O}_3$ ,  $\text{HfO}_2$  ...), and the metal

gate (TiN ...) respectively. An oxide can then be deposited before performing a CMP step on the top. Isotropic etching allows the removal of a part of the metal gate and defining the gate length. The space left by the metal gate can be filled by oxide deposition. The top contact is achieved by top line formation using a conventional damascene process.

### 3.4.3 Vertical 3-D Logic

Since the described process could be used to integrate transistors into the metallic layers of a standard circuit, we can envisage adding some functionality, typically in relation to the back-end layers directly above the circuit.

#### 3.4.3.1 Smart Vias

The association of vertical transistors and reconfigurable circuits leads immediately to the concept of “smart” reconfigurable vias. Here, the connection between a metal line and another is done by connecting the line to a transistor-drain and the other to a transistor-source. Thus, the connection is only controlled by the gate. This requires the line to be connected to the front-end and requires the use of several layers of vias to contact the back-end lines to the transistors. Here, based on the 3-D FET technology, figure 52 shows an implementation of a controlled contact between two metal lines controlled by a combinational logic function.

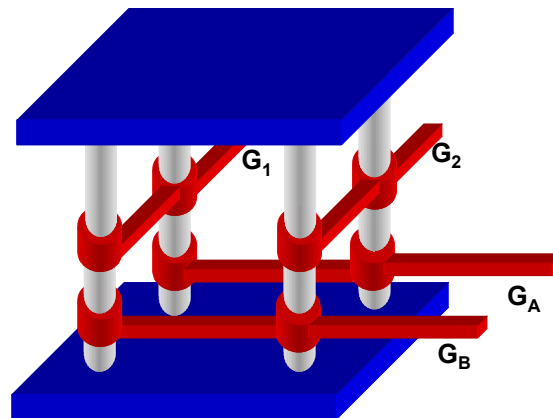


Figure 52. Boolean logic with vertical FETs

In the presented illustration, four vertical wires are shown. Two transistors per wire are built in series. Four control gates drive the wires in a meshed pattern, which are used to create a combinational equation for the via conduction. The functionality is thus increased. The considered metal lines are connected when all transistors of any vertical branch are conducting. This corresponds to the following Boolean condition:

$$(G_A + G_B).(G_1 + G_2) = 1$$

#### 3.4.3.2 Logic Gates

From the previously presented fully configurable via pattern, complementary logic cells can be built by including the dual branch using *p*-type transistors. Based on the 3-D FET technology, we propose a fully back-end implementation of standard logic gates.

As an illustration, figure 53 shows a NOT gate built with vertical transistors. The circuit is formed of two vertical transistors of different doping types. An *n*-type transistor is built between the ground and the output lines, while a *p*-type transistor is placed between the power supply and the output lines. The peculiarity of the structure is that output lines and the power lines are separated by a metal layer. This metal layer is used to realize the transistor gates. It is worth noticing that different nanowire diameters could be used in the transistor channel, in order to size the structure according to design specifications. The proposed layout can be extended to any complementary logic gate.

We show in figure 54 the implementation of a NAND gate. Various actual layouts could be envisaged. Figure 54-a shows an implementation where the serial transistors are shared on the same nanowire. While this solution is compact, it is worth mentioning that it also increases the vertical dimensions of the gate, because of the occupation of two levels of metal. Another

solution is depicted in figure 54-b. In this implementation, we only have one transistor per nanowire, i.e. between two metal lines. This allows the dimensions of the cell to be controlled, by analogy with the 2-D standard cells, where only one level of metal is used for the interconnect.

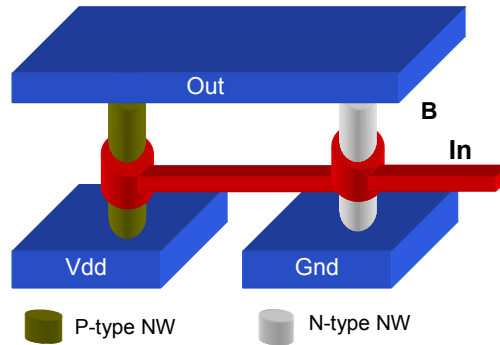


Figure 53. NOT function implemented in a 3-D standard cell

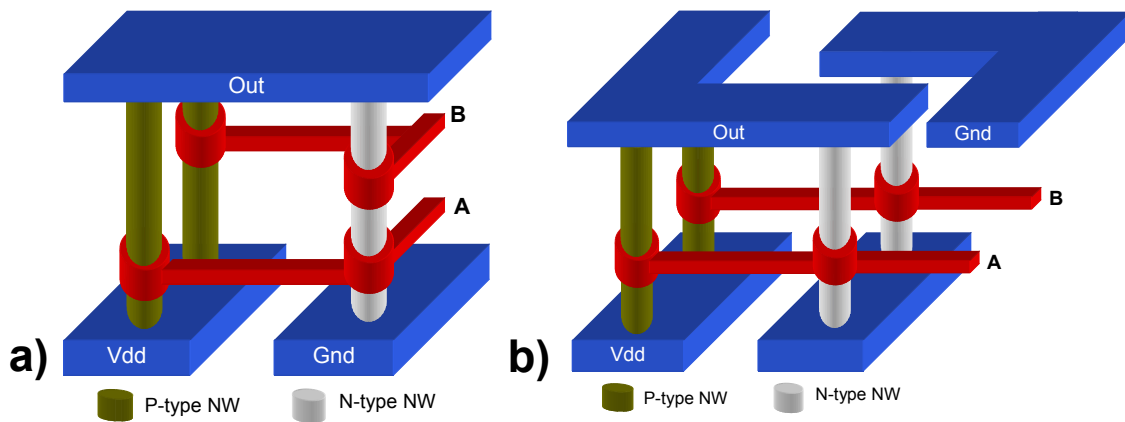


Figure 54. NAND function implemented in a 3-D standard cell – a) layout on two gate levels b) layout with only one gate per nanowire

Back-end standard cells appear to be of use for routing application. Signal lines could be buffered directly in the back-end, thus avoiding connecting front-end buffer cells merely for routing purposes. This makes sense for long signal lines, as well as for clock signals.

#### 3.4.3.3 SRAMs

Since we have shown that we could embed any logic function in a real 3-D approach, it then appears quite straightforward to build memory elements. Figure 55 shows a possible organization for a 5T-based SRAM cell. Using the same approach as for the NAND gate, the layout is realized only between two metallic layers. The interest of these elements will be twofold. Firstly, they will be used as configuration memories for the reconfigurable logic front-end. As already mentioned, positioning memory resources above the circuit is of high interest for reconfigurable logic compactness and performance. Furthermore, it is also of use to have memory elements which are realized in the same technology as the smart vias presented above. These nodes will allow the configuration of the smart vias to be stored very close to their final use, and will thus lead to extremely compact cross point nodes entirely implemented in the back-end levels, as shown in the top view in figure 56. While the cross point uses a standard CMOS architecture as shown in figure 56-a, its implementation is based on the use of smart vias and 3-D SRAM cells. In fact, four 3-D SRAM cells are used to store the configuration as close as possible to the logic circuit which uses the information. Two other SRAM cells implemented in 2-D are also used. Thus, it is possible to mix the different technologies to obtain a structural trade-off between the requirements on the front-end area and on the routing ability. It is actually not advisable to embed all the logic in the back-end, but only to place a part of it in the most efficient way.



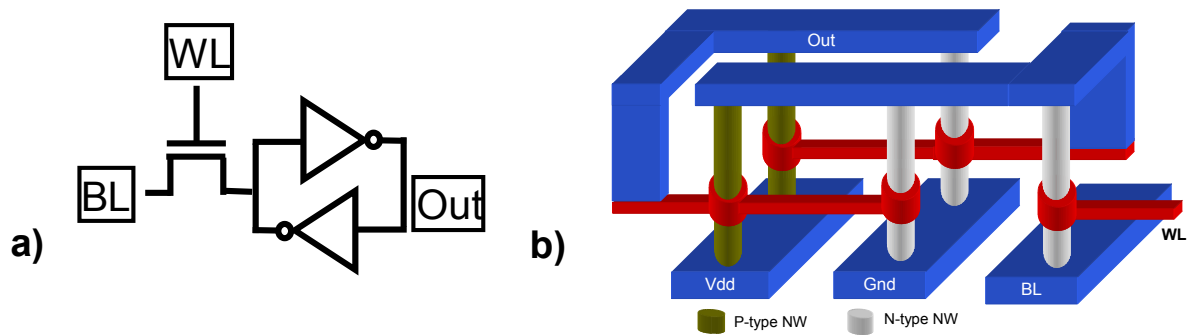


Figure 55. 5T SRAM cell implemented in a 3-D standard cell – a) schematic b) cross view

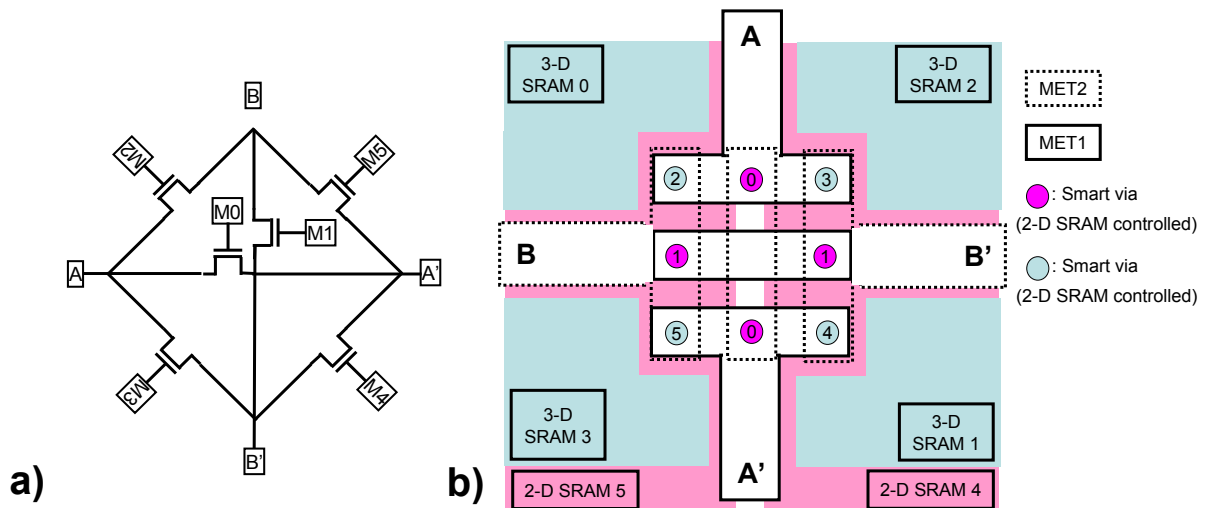


Figure 56. Cross point implemented using 3-D smart vias and 3-D SRAMs – a) schematic b) top view

### 3.4.4 Performance Characterization

#### 3.4.4.1 Methodology

In order to evaluate the impact of the back-end implementation of standard cells, we investigate the area, delay and power consumption of elementary gates. Then, we compare our logic gates to the equivalent cells taken from an industrial 65-nm CMOS bulk process design kit.

Since the considered technology is not mature enough, there are no compact models available. The evaluation of 3D-BE standard cells have been performed thanks to TCAD ATLAS [114] simulations. TCAD simulations are computationally expensive, but they are flexible and allow new fundamental devices to be modeled efficiently. These characteristics mean that while they are widely used for physics and device simulation, they are rarely used for circuit design. Nevertheless, they represent new opportunities for design methodologies working close to the technology. In our original approach, we will use TCAD simulations, not only to simulate the performances of the unique device, but also to simulate standard circuit and to extract its performance numbers.

Vertical nanowire transistor structures have been modeled according to the proposed device architecture (Figure 57). The circuit netlists are then input to the TCAD simulator and transient simulations are performed. From the transient simulation results, the performances metrics are then evaluated. The complete methodology is described in figure 58.

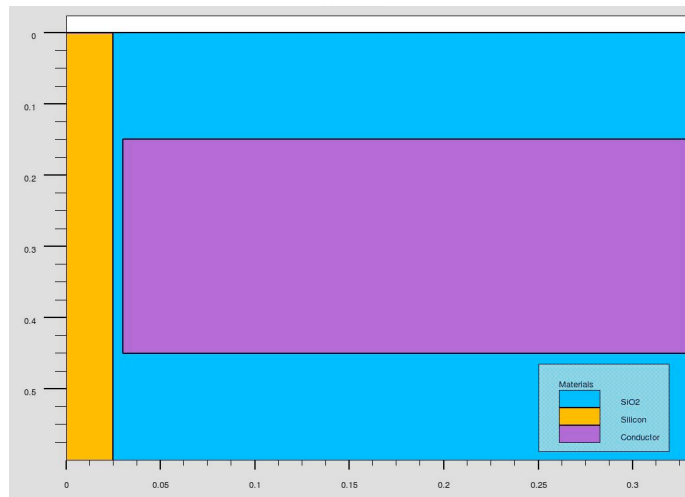


Figure 57. TCAD cylindrical representation of the vertical NWFET structure (scales are in micron)

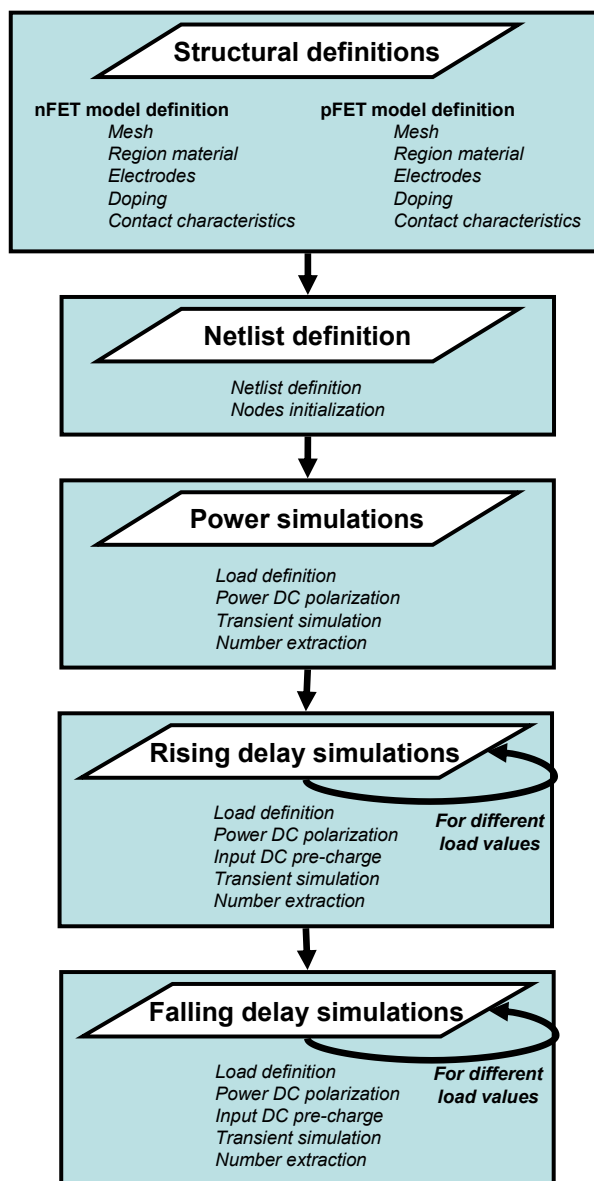


Figure 58. TCAD circuit characterization methodology flow chart

### 3.4.4.2 Device Characterization

In order to evaluate the technology at the lowest level, technological simulations have been used to extract the I-V curves of the elementary vertical transistors. Figure 59 depicts the I-V curves of an n-type and a p-type vertical FET. These curves are extracted for a NWFET with a 50-nm diameter, a 300-nm gate length and 5-nm oxide thickness. It is worth noticing that

the devices present a very good  $I_{on}/I_{off}$  ratio in the range of  $10^8$  and a low  $I_{off}$  current less than 0.1pA. These excellent properties can be explained by the structure of the transistor. First of all, the transistor is a Gate-All-Around structure, which means that the active zone must be considered all around the nanowire. This makes the device highly controllable, while the dimensions are maintained very compact. The structurally optimal electrostatic control of the device leads to the low  $I_{off}$ . Furthermore, the structure uses the vertical dimension to increase drastically the dimensions of the active region, while the front-end impact is maintained very small. Thus, it is possible to realize large FETs, with excellent electrical properties and very compact front-end projection.

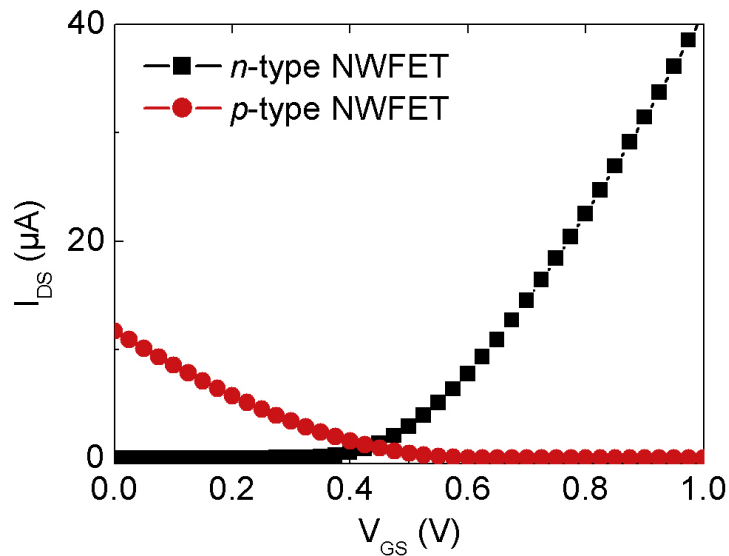


Figure 59. I-V characteristics for vertical nanowire FET

The estimated performance levels of the elementary transistor lead to the realization of high performance switches. For routing applications, it is of high interest to achieve pass transistors which exhibit a low  $I_{off}$  current, while the  $I_{on}$  current is in the range of several  $\mu\text{A}$ . It is clear that the higher the  $I_{on}$  current, the lower the  $R_{on}$  resistance and consequent switch impact on the overall circuit performance.

#### 3.4.4.3 NOT Gate Characterization

While it is quite standard in TCAD simulation to extract the I-V curve of a device from the simulation, it is also of high interest to evaluate, from a transient simulation point of view, the characteristics of simple circuits. Table X summarizes the performance results of a NOT gate and compares it to a standard bulk cell. Such a comparison is helpful to evaluate the properties of the technology not only for the device but also for its ability to improve circuit performance.

Table X. Simulation results summary

	Area ( $\mu\text{m}^2$ )	Intrinsic delay (ps)	$K_{Load}$ ( $\text{ps}\cdot\text{fF}^{-1}$ )	Leakage power (pW)
NOT (MOS 65nm)	1.6	17.5	3.8	40.1
NOT (3-D BE)	0.05	6.9	4.7	2.8
3-D BE vs. CMOS	x 31.2	x 2.5	x 0.8	x 14.5

We see that the proposed 3-D implementation of a NOT gate clearly improves the area by a factor of 31.2x. These results come from the use of the vertical direction to implement the transistors. Delay and leakage power are also improved by a factor of 2.5x and 14.5x respectively. This can be explained by the good performance levels of the gate-all-around control of transistors. Indeed, the good electrostatic control of the channel helps to have a well defined off state, and thus to obtain a very low  $I_{off}$  current, while the  $I_{on}$  current remains high (the  $I_{on}/I_{off}$  ratio is high). Thus, since these characteristics are directly involved in the power

consumption and delay estimations, the figures show clear improvement over the equivalent MOS technology.

#### 3.4.4.4 Discussion

These characteristics show some promising performances for the technology. We have seen that the elementary device can be sized efficiently, thanks to the use of the third dimension. It is thus possible to obtain high performance devices, while the impacted area on the front-end is maintained very low. This opens the way towards high performance switching between two metal lines, which is of high interest for routing resources in reconfigurable circuits such as FPGAs. Furthermore, we have seen that it is possible to implement logic functions, using complementary logic with this technology, and that the obtained circuit performance gain is also significant with respect to CMOS. While the area is improved again by the use of the third dimension, we also obtain high performance figures, due to good device performance levels. It is thus of great interest to examine the way such logic could embed elementary active functions, as required for the routing circuits and clocks. For example, it is possible to embed the signal buffers as well as high performance clock tree buffers along the line which requires this functionality. This avoids inefficient multiple vertical communication schemes such as those seen today between the active transistor front-end and low-RC metal back end.

### 3.5 Global Comparisons and Discussions

Several technologies have been investigated in this chapter in order to realize a part of functionality of the FPGAs into the back-end levels. All these techniques integrate an above-IC device. While they are all promising for routing resource improvement, it is interesting to compare them using the same template. To do so, we use a 4-input LUT, which is the elementary block used in modern FPGAs to perform the computation.

Table XI compares the performance of the structure realized with the various technologies. The assumptions for the realization differ from the technology, since we need to consider the programming circuitry for a fair comparison. The bulk implementation serves as a reference. The storage elements are implemented as SRAM cascaded in shift registers, while the data path is composed of multiplexers. The ReRAM implementation uses multiplexers for the data path, but their data inputs are driven by the configuration nodes presented in 3.2.3.1. The programming structure has been incorporated into the evaluation, at least to ensure the fair selection between the power lines. The monolithic 3-D FDSOI implementation uses multiplexers realized on the bottom silicon layer, while the SRAM shifted configuration memories are built on top. Finally, the vertical NWFET implementation uses standard 2-D multiplexers to realize the data path, while the configurations are placed above the data paths using the circuit proposed in 3.4.3.3.

In terms of area, the most compact solution is obtained for the vertical NWFET technology. Indeed, the improvement of 10.2x is the consequence of the full vertical integration of the memory circuits. While the circuit is placed in the third dimension, its volume remains constant while the front-end projection (i.e. the area occupies by the drawn layout) is very low. In second place comes the ReRAM implementation with a gain of 7.3x. In this implementation, all the configuration memories are placed above the circuit, but programming access transistors are still in front-end silicon. Finally, the 3-D FDSOI improves the structure by 5.7x, thanks to the 2-stack repartition of transistors.

In terms of performance, the best figures are obtained with 3-D FDSOI. In this case, the data path technology is moved towards enhanced silicon FDSOI. This leads to an improvement in the intrinsic delay of a factor of 4.4x, and of the load factor by 1.9x. For the other technologies, the performance levels are of the same order of magnitude as CMOS. This is due to the fact that the data path remains based on CMOS multiplexers and is thus the same between the different implementations. Nevertheless, it is worth noticing that the intrinsic

delay is improved by 1.5x. This is because the data paths differ by one gate. Thus, the propagation delay is reduced by one intrinsic delay.

Finally, in terms of power consumption, only trends have been extracted, since an accurate estimation would require several simulations that are not possible to the unavailability of compact models for the vertical NWFET. Based on the previous results, we can remark that the ReRAM solution will remain of the same order of magnitude as the bulk solution. Indeed, while the data paths remain the same, the leaky SRAM circuits have been replaced by other potentially leaky memory nodes, such that the figures can be expected to be constant. Concerning 3-D FDSOI, power might be reduced by at least 2, thanks to the low power characteristics of FDSOI. Finally, we have seen that vertical NWFET are large and possess good electrical properties. This provides them with a strong suitability for low power, and thus the trend can be estimated to be in the range of 10x.

**Table XI. Technology comparison (4-input LUT test case)**

	<b>Area (<math>\mu\text{m}^2</math>)</b>	<b>Intrinsic delay (ps)</b>	<b><math>K_{\text{Load}}</math> (ps.fF<sup>-1</sup>)</b>	<b>Average power gain trend</b>
<i>Bulk MOS 65nm</i>	547	465.6	11.2	reference
<i>ReRAM</i>	74.5	310.4	11.2	=
<i>3-D FDSOI</i>	95.1	105	5.94	x 2
<i>Vertical NWFET</i>	53.8	310.4	11.2	x 10
<i>ReRAM vs. CMOS</i>	x 7.3	x 1.5	x 1	-
<i>3-D FDSOI vs. CMOS</i>	x 5.7	x 4.4	x 1.9	-
<i>Vertical FET vs. CMOS</i>	x 10.2	x 1.5	x 1	-

### 3.6 Conclusion

In this chapter, we have investigated how disruptive technologies can be used to build enhanced basic logic circuits for routing and configuration. These elements are fundamental for FPGA technology since they represent the largest amount of required area (over 80%) and constitute the bottleneck for performance improvement

Globally, the improvement of these circuits is based on the use of 3-D integration technologies. In fact, we studied three different technologies, which respectively allow (i) passive reconfigurable elements to be placed in the back-end layers, (ii) several layers of active silicon to be stacked with a high alignment accuracy and a high via density and finally (iii) "true" 3-D transistors to be embedded in the metal layer in a vertical arrangement.

Resistive memories, and especially phase-change memories, have been initially considered. This technology is able to place a passive non-volatile resistive material in the back-end layer. This node can be reprogrammed between two stable resistive states. We propose a simple logic circuit based on 2 resistive nodes and 1 transistor in order to create a configuration node for reconfigurable logic. This node leads to an improvement of 1.5x in terms of area compared to its flash counterpart. Furthermore, we notice that the on-resistance of a resistive memory is very low compared to a CMOS transistor. We thus propose a switchbox structure that places the resistance directly in the logic data path. The solution is compact compared to the flash equivalent circuit with a gain of 3.4x in area, and we expect further improvements at the system level due to the reduction in the path resistance.

A monolithic 3-D FDSOI integration process was then considered to embed active devices into the back-end layers. Such a technology is of particular interest for separating memory resources from the logic data paths in a reconfigurable device. In this context, it is possible to optimize the circuit technologies separately towards low-power or high-performance operation. Monolithic integration leads to a high via density, which allows fine grain architecture partitioning over the layers. To demonstrate the gain of the monolithic 3-D

integration, simple test structures were designed. We demonstrated a simple LUT structure, which improves the area by 2x, the delay by 1.6x and the power by 2x compared to its 2-D CMOS counterpart. We also implemented a 3-D cross point using 2 different configuration circuits. We demonstrated that this circuit yields an improvement of 1.5x in area, 1.6x in performance and 3.1x in power. This shows the benefits of the FDSOI technology when coupled to monolithic 3-D integration, when compared to standard CMOS bulk.

Vertical NWFET has been finally assessed. This technology could be seen as the ultimate evolution to 3-D, since it allows the integration of vertically oriented transistors. Thus, the designer can go beyond a stacked design of 2-D devices, and can distribute the active devices within the back end layer with a very small impact on the front-end. This potential leads to the development of a smart back-end and the generalization of 3-D computing. In order to evaluate the interest of the technology, we performed a TCAD evaluation for a simple circuit, and showed an improvement of 31x over CMOS in terms of area, while the delay is improved by 2.5x and the power by 14x.

To conclude this chapter, we have assessed three different technologies with a view to improving the efficiency of elementary logic blocks. These logic blocks have been chosen, in order to address the field of reconfigurable logic devices. Hence, we expect to improve accordingly the performance of test case FPGA architecture. Nevertheless, this kind of evaluation must be conducted fairly, thanks to the use of generic benchmarking tools and a set of well-known-application circuits. This evaluation will be the focus of the next chapter.



## **CHAPTER 4** *Architectural Impact of 3-D Configuration and Routing Schemes*

---

### **Abstract**

In this chapter, the architectural impact of the 3-D enhanced memories and routing resources were carefully studied. The traditional FPGA architecture was enhanced by the technologies presented in the previous chapter. The envisaged technologies move devices in 3-D. Devices can be passive (e.g. resistive phase-change memories) or active (e.g. monolithic 3-D integration or vertical NWFET). Performance estimations were carried out by benchmarking simulations of the improved FPGA architecture. The benchmarking tool is based on standard tools and tuned according to the technological parameters.

We showed that, implemented in FPGAs, the resistive configuration memory node, coupled to the routing structure, yields a delay reduction up to 51%, thanks to the reduction of dimensions and low on-resistance of PCMs. This result was also reached by the vertical NWFET technology, because of the ability to size a large transistor vertically without a large impact on the projected area. In this case, the critical path delay may be reduced up to 49% compared to the traditional scaled MOS.

Regarding the area metric, the best improvement was reached by the vertical NWFET technology with an improvement of about 46%. Vertical NWFET technology allowed moving all the peripheral circuits above the IC. By opposition, the PCM technology leads to a much tighter area improvement of 13%. Indeed, this technology requires a large programming transistor per node.

Among the different technologies, we should remark that 3-D monolithic integration process yields in an area improvement of 21% on average and in a delay improvement of 22% on average. Such a technology represents a good trade-off process for short term micro-electronics evolutions.



In the preceding chapter, we proposed several approaches to enhance the performance of routing and memory structures in standard FPGA architectures, using 3-D integration techniques, where part of the circuits is stacked above the conventional silicon layer. The use of resistive memories allows passive memory devices to be placed in the back-end layers. For active devices, we explored the use of monolithic 3-D integration and vertical nanowire FETs, to respectively stack 2-D FDSOI devices in a 3-D scheme, and orient the active device channels in the third dimension. We analyzed the improvement of routing and memory blocks in terms of footprint, write time, data path resistance, etc. While these alternative circuits have been shown of interest in terms of their intrinsic properties, it is also mandatory to examine how they improve performance in a complete environment. Hence, in this chapter, we will focus on real-life circuit benchmarking, using FPGA architectures enhanced with the proposed approaches. After describing the motivation for architectural evaluation, we will describe the benchmarking tool flow that allows the evaluation. Then, for all the previously introduced proposals, we present a specific organization and the results of the benchmarking.

## **4.1 Motivation and Global Methodology**

Research into emerging technologies is generally focused on devices and simple circuits. Nevertheless, it is necessary to pursue the analysis and assess the amount by which the use of a new technology can improve a complete application test circuit. The cost actually required developing, stabilizing and ramping up a new process can only be justified only if circuit performance improvement is significant. Estimating such performance gain at an early stage in technology development is thus a critical step to channel work on emerging technologies. Thus, as introduced in section 1.5, a methodology based on the evaluation of emerging technology through a complete application context is presented.

In this work, we investigate reconfigurable logic applications, which are well-suited to disruptive technology analyses, due to the high architectural regularity (mirrored in the regular arrangements of emerging technologies), and high flexibility, an important factor for fault-tolerant circuits in the context of unreliable devices.

In this section, we propose an assessment of the impact of the elementary blocks designed earlier on the complete FPGA architectural scheme. Indeed, while the individual block performance metrics are improved, it is necessary to check their impact at the system level.

Thus, we will consider the conventional island-style FPGA scheme, according to its description in section 2.1.1.2. In particular, the logic part of the architecture is formed by CLBs. The CLBs are built by 10 BLEs of 4-input LUTs. 22 inputs connect the CLB to routing lines. Hence, we will replace the MOS routing and memory elements by the proposed structures.

In order to evaluate the performance metrics, it is necessary to set up a benchmarking toolflow. Part of the task of benchmarking is to “program” the reconfigurable structure with a set of well-known circuits, which constitute a representative sample of most application tasks the circuit has to perform. The other part of the benchmarking task is to evaluate and compare several metrics for each circuit. The most commonly used metrics are: the area used for logic and for routing, and the critical path delay through the routing structures. Others can also be added, such as power consumption, scalability, robustness to name a few.

## **4.2 Benchmarking Tool for FPGA-like Architectures**

A conventional FPGA benchmarking flow is described in figure 60. The toolflow has been developed by the Toronto University [185], and is built around the T-VPACK and VPR5.0 tools. T-VPACK handles the packing of the logic blocks, while VPR5.0 handles the place and route of logic blocks in the island FPGA scheme. The input to T-VPACK is a netlist of LUTs and FFs. Thus, the logic packing operation consists of grouping the LUTs and FFs into BLEs, and then in packing several BLEs into CLBs. This operation should take into account all

timing constraints. In particular, inter-CLB routing paths are faster than intra-CLB ones, such that it is possible to find an optimal packing arrangement that allows the connection delay to be globally improved. The netlist of CLBs, output from T-VPACK, is input into the VPR tool. VPR starts by placing the CLBs onto the island-style structure. The optimal position of the blocks over the grid is found using a simulated annealing algorithm [17]. Then, connections between CLBs are found using a pathfinder algorithm [17]. Obviously, all these steps use timing-based metrics. Finally, VPR outputs several metrics such as the total area used by the circuit and the critical path delay. To ensure the connection between the VPR flow and standard BLIF benchmarks [204], the ABC tool performs the logic optimization, the technology mapping onto a set of LUTs, and finally outputs the netlist of LUTs and FFs that is used by the flow.

This toolflow was initially developed for architectural exploration in the context of LUT-based FPGAs. It is therefore not possible to target an architecture that is not based on LUTs. Nevertheless, it is worth pointing out that the VPR tool is highly versatile in terms of architectural description, expressed in XML. The XML file consists of several sections that specifically describe part of the FPGA area. Hence, it is possible to specify the properties of the segment switches (i.e. the transistors that are used for routing), in great detail (albeit with a large number of parameters). Thus, the architecture can be adapted quite easily to advanced routing technologies, using this architecture fine tuning.

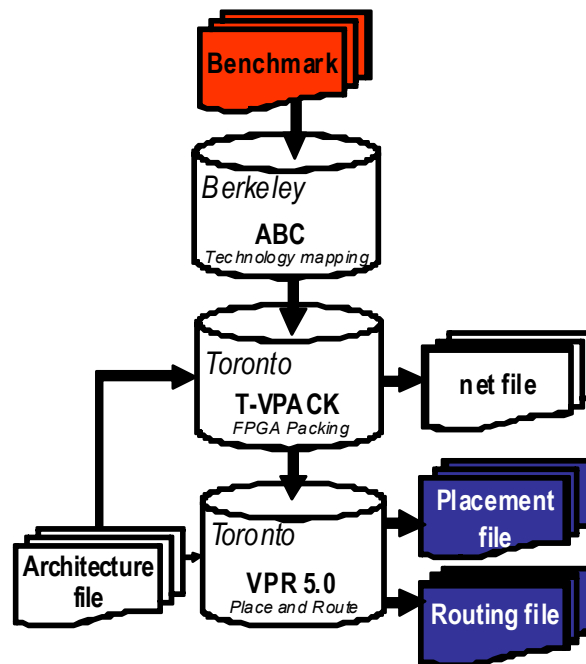


Figure 60. FPGA benchmarking flow diagram

### 4.3 Resistive Memory-based FPGA Performance

In the previous chapter, we proposed two circuits using resistive memories: the configuration memory node and the switchboxes. In this section, we will assess their impact on the FPGA architecture. First, we will draw a global view of the improved FPGA structure and then will discuss the benchmarking results.

#### 4.3.1 PCM-FPGA Architecture

The potential improvement of the FPGA scheme is twofold. Firstly, concerning the configuration nodes, we propose to replace the area-hungry SRAM by a smaller and more efficient circuit. While this replacement is expected to reduce the size of the complete FPGA, it is also of high interest to take advantage of the low on-resistance of ReRAM. A switchbox implementation has been proposed, which uses resistive memories to replace the pass switches.

The PCM-based FPGA will use both of these proposals. Figure 61 gives an illustration of the structure. Configuration nodes will be used to feed the multiplexer inputs of each LUT. Furthermore, since such a configuration node is able to place a fixed logic level on a logic gate, it is also intended to be storage mechanism for routing multiplexers, which can be found at every level. A multiplexer is used in the BLE, in order to choose between the registered or the unregistered version of the LUT signal. Multiplexers are also used at the CLB level for the local interconnect. Finally, routing multiplexers are used in the connection boxes to select the signal to be connected to a CLB input. For all of these routing BLEs, the proposed configuration node will be used instead of an SRAM. Finally, all switchboxes will be replaced by the proposed PCM-based switchbox.

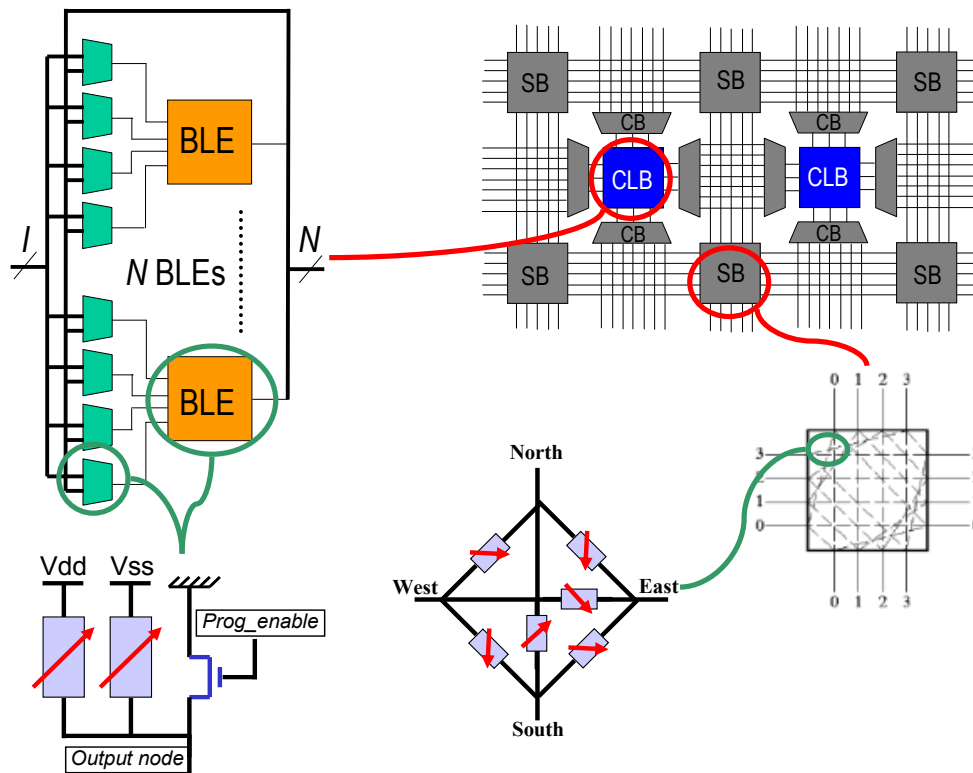


Figure 61. PCM-based FPGA organization

### 4.3.2 Methodology

To perform the evaluation, we used the standard FPGA benchmarking flow. A set of logic circuits taken from the MCNC benchmark [204] were first synthesized using the ABC tool [183]. We then performed the technology mapping with a library of 4-input LUTs ( $K=4$ ), also using ABC. We subsequently performed the logic packing of the mapped circuit into CLBs with  $N=10$  BLEs per CLB and  $I=22$  external inputs using T-VPACK [186]. Finally, the placement and routing were carried out using VPR [186]. We synthesized the considered benchmark twice. The first (conventional) design was based on SRAM-based LUTs and MUXs in a 65-nm bulk CMOS process, using a pass-gate design. In the second (experimental) design, we replaced the SRAM cells in the LUTs and routing MUXs by PCMs. It is worth pointing out that, in routing topologies, PCMs are placed on the data path. Contrary to CMOS logic, which requires signal restoration, output buffers become superfluous with PCM technology.

### 4.3.3 Simulation of Large Circuits

In order to study the impact of each contribution, we analyzed two variants: one with an FPGA scheme improved by only the routing structures, and a second with a complete improved arrangement.

### 4.3.3.1 Impact on the Routing Structures

In this first variant, PCMs have been introduced only in the logic data path of the routing structures. We mapped the benchmarks onto both PCM- and SRAM-based FPGAs, and ran simulations to estimate the FPGA delay, as shown in figure 62. These benchmarks show a delay reduction ranging from 38% to 51%, with 44% on average. Hence the main benefits of using PCMs instead of SRAM cells are the compact area of the cell and the lower internal resistance of data paths. As a matter of fact, we extracted from our design kit the internal resistance of a pass-gate cell which is of the order of the on-resistance of an n-type transistor (9.1k $\Omega$ ); while the experimental results show that PCMs have a lower on-resistance that has been reported close to 4 k $\Omega$  [115]. This makes the PCM-based switchboxes potentially faster than the SRAM-based counterparts.

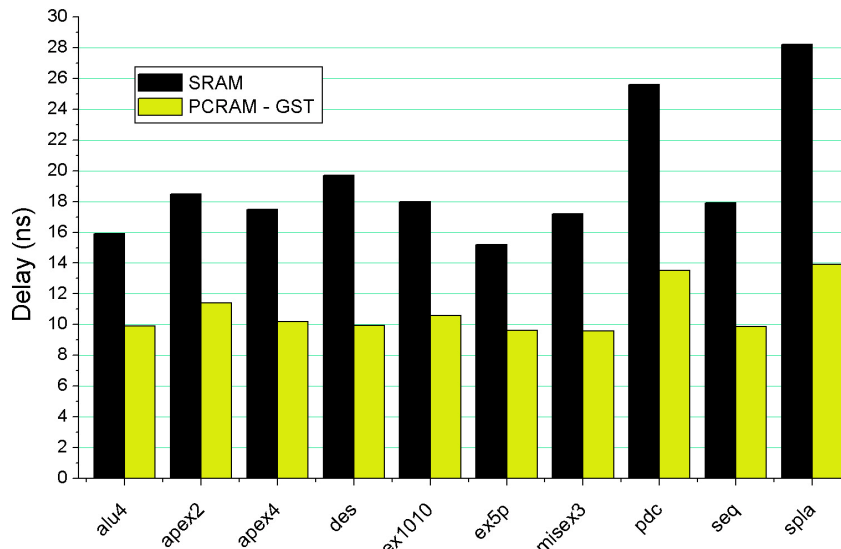


Figure 62. Delay estimation for FPGAs synthesized with GST-PCM- and SRAM-based Switchboxes

### 4.3.3.2 Impact on the Configuration Memories

We showed that the introduction of PCMs directly into the logic data path leads to an improvement in delay. In this section, we also add the memory node to the evaluation. Since routing is more compact than that of standard SRAM, we expect an area reduction. We mapped the benchmark onto PCM- and SRAM-based FPGAs, and ran simulations to estimate the FPGA area and delay. With respect to the previous evaluation, we have added more benchmarks in order to achieve a better test circuit diversity. We measured an area saving of about 13% on all our samples, due to the decrease in area allocated to memory. The delay estimation is depicted in figure 63, and shows a delay reduction ranging from 22% to 51%, with 40% on average. Here, the main benefits of using PCM- instead of CMOS-based cells are the compact area of the cell and its lower internal resistance. Again, we extracted from our design kit the internal resistance of a CMOS cell, of the order of the on-resistance of an n-type transistor (9.1 k $\Omega$ ); while the measurements show that our PCMs have an on-resistance of 3.7 k $\Omega$  [115]. This makes the PCM-based FPGA potentially faster than the SRAM-based counterparts, given the lower resistive data path through the memory and the associated pass-gate MUX. Further, the compact area of the cell allows for a reduction of the CLB size and consequently a lower wire delay. The delay reduction due to smaller overall area is however less significant than the delay reduction due to faster routing elements.

## 4.3.4 Impact of Technologies

We demonstrated with the previous simulations that the combination of the switchbox design and the GST technology results in a significant delay improvement. In the following, we showcase the impact of the PCM technology type on the delay improvement. We remind the reader that other materials exist, which may replace the GST as a phase-change material. These include GeTe and GeTeC $\alpha$ %. Besides the difference in RESET current and time, the on-resistance depends on the chosen material and it has an impact on the routing path

resistance. We simulated the delay of the FPGA benchmark with various resistance values corresponding to different materials. We notice that the delay improvement is linear with the decrease of the on-resistance. However, the delay sensitivity is low: a decrease in on-resistance by 2 orders of magnitude from about 4000  $\Omega$  to about 50  $\Omega$  causes a delay reduction of only 5%. The reduction of the on-resistance of the switches decreases the inter-CLB routing delay. In this situation, the intra-CLB routing delay becomes the dominant contribution and a larger reduction of the external switches on-resistance has no drastic impact.

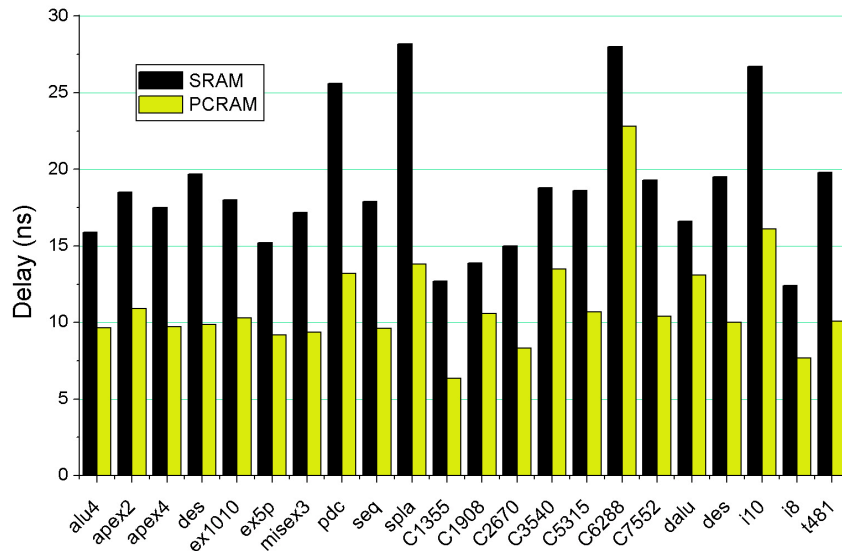


Figure 63. Delay estimation for FPGAs synthesized with PCM- and SRAM-based LUTs and MUXs

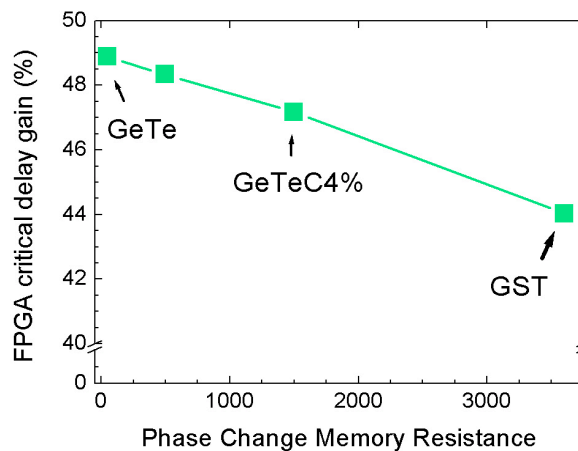


Figure 64. Variation of the FPGA critical path delay with the PCM on-resistance (delay averaged over the whole benchmark set)

### 4.3.5 Discussion

In this section, we analyzed various performance metrics of a PCM-based FPGA. Our first observation is that the highest gain in terms of delay comes from the improvement of the routing structure. In fact, the reduction of the switch on-resistance yields an improvement of 40% on average. The configuration memory improvement is less significant and consists essentially of a reduction in terms of CLB area. This reduction appears quite low, at only 13% on average.

These results demonstrate significant benefits through the use of this technology. Indeed, it is worth pointing out that the technology is back-end of line compatible with CMOS technology. This makes it remarkable for its low-cost properties, while giving an improvement of 40% in delay.

We must also consider the limitation of the technology. Phase-Change memories require only a small amount of energy for their programming, which makes them power efficient.

However, since the phase-change mechanism is based on the heat control of the material, the current density must be high. It is thus necessary to drive a large current into the memory node through a wide programming transistor, which remains costly in terms of silicon area. Furthermore, regarding the design of the configuration memory node, we should notice that there is inherent leakage, due to the off-resistance of the material, which must be kept as high as possible.

In order to retain the advantages of the modulated resistance, and mitigate the limitations, it is possible to consider other kind of ReRAM technologies. For example, it is worth noticing that Oxide-based memories could be programmed with a smaller current [86], while Conductive-bridge memories demonstrate a high  $R_{off}$  value [87].

## 4.4 3-D Monolithic Integrated FPGA Performances

In the preceding chapter, monolithic 3-D FDSOI basic blocks have been designed. In particular, memories were integrated on top of logic circuits. These blocks have been studied in comparison to standard FPGA blocks. In this section, we will assess the performance of a complete FPGA circuit, built in a 3-D manner.

### 4.4.1 Overall FPGA Architectural View

In a CMOS FPGA, memories are found throughout the logic. Indeed, configuration memories are used to drive the LUT inputs, as well as to configure all the routing multiplexers. In the following FPGA evaluation, we will use the LUT scheme proposed in 3.3.3.2. As already stated, routing multiplexers can be found at each hierarchical level, from the BLEs to the Connection Boxes. Each routing multiplexer structure will use a top memory driving the configuration path of the multiplexer. All the switchboxes will use the scheme shown in 3.3.3.3. Figure 65 depicts a final view of the envisaged FPGA organization.

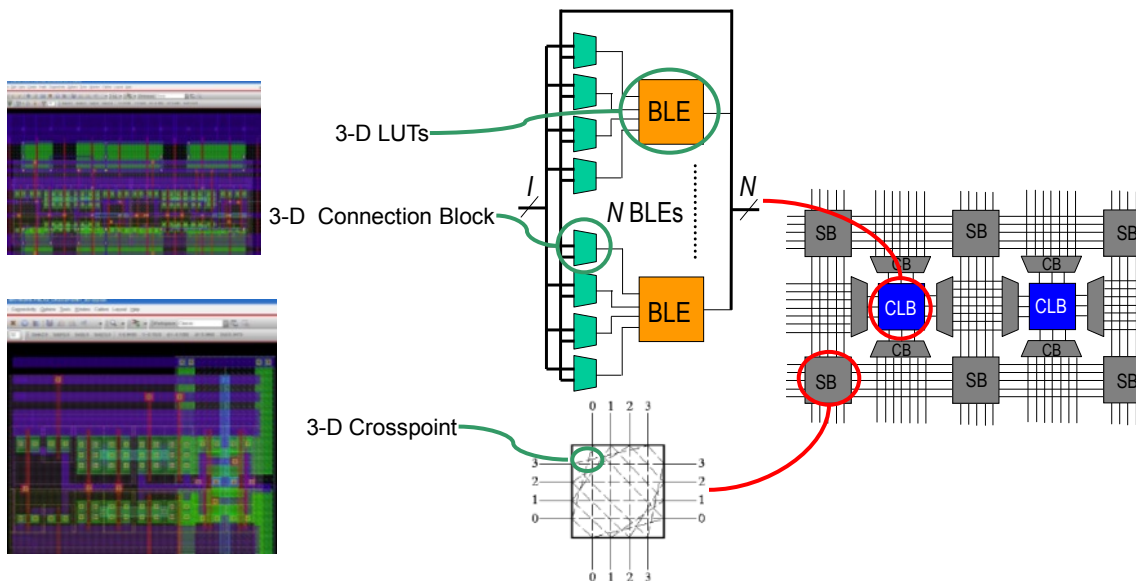


Figure 65. Monolithically 3-D integrated FPGA organization

### 4.4.2 Methodology

The methodology is in principle the same as that described in 4.3.2. The standard FPGA benchmarking flow will be used to map a set of logic circuits taken from the MCNC benchmark [204]. Optimal sizing for the FPGA will be used with  $K=4$ ,  $N=10$  and  $I=22$  [175]. As already stated, the envisaged technology may enhance the logic block organization and the routing part. Thus, parameters are adapted accordingly in the architecture description files. The technology node under consideration for the elementary block sizing and performance evaluation is a 65-nm bulk CMOS process.

### 4.4.3 Simulation of Large Circuits

#### 4.4.3.1 Monolithic 3-D Integration Impact on Global Area

We mapped the benchmark in Bulk and monolithic 3-D integrated FDSOI FPGAs, and we simulated the FPGA area. The area estimation is normalized and shown in figure 66. The benchmarks in figure 66 show an area reduction ranging from 20% to 25%, with 21% on average. The main benefit of using the 3-D integration technology is the smaller impact on memory circuits on the area. Memories circuits are placed on top of the FPGA structure, while all the data path circuits are placed above. Nevertheless, since the top transistors are directly stacked above the bottom silicon layer, the circuit suffers some loss of performance in terms of routing. Indeed, the top transistors are impeding the connections between bottom transistors and the first metal layer. To improve the quality of routing and find the optimum routing organization, it is of interest to study the best trade-off in terms of internal metal layers.

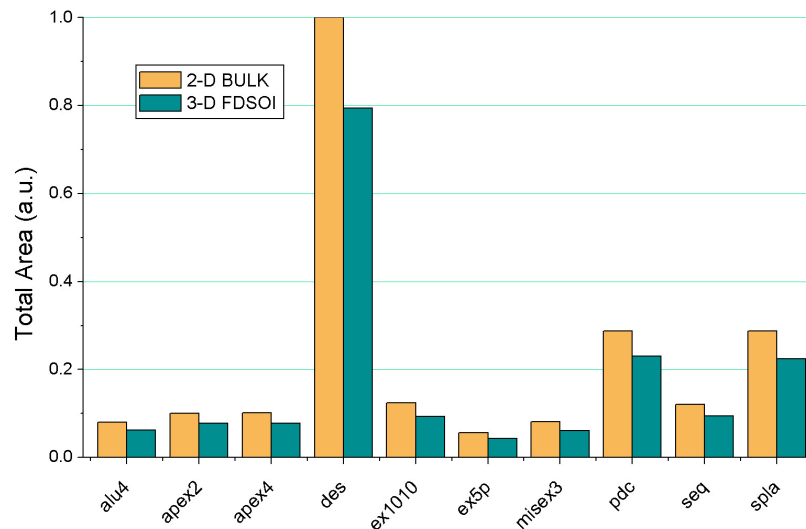


Figure 66. Area estimation for FPGAs synthesized with standard bulk circuits and monolithic 3-D integrated FDSOI circuits.

#### 4.4.3.2 Monolithic 3-D Integration Impact on Critical Path Delay

The delay estimation with respect to the benchmark circuits is shown in figure 67. We can observe a delay reduction ranging from 10% to 45%, with 22% on average.

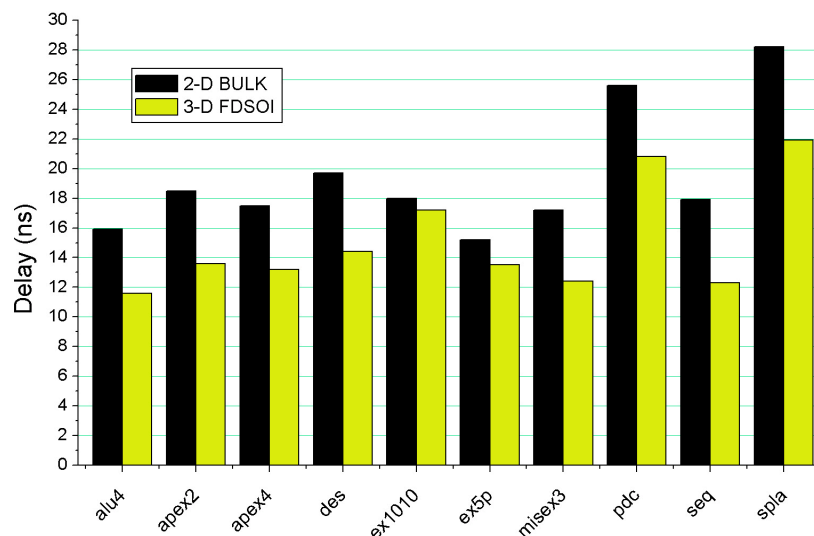


Figure 67. Delay estimation for FPGAs synthesized with standard bulk circuits and monolithic 3-D integrated FDSOI circuits.

The main benefits of using the 3-D FDSOI instead of standard Bulk are twofold. First, the intrinsic performance of FDSOI compared to bulk contributes to speed up the gate delay and thus reduce the computation time. Second, the reduction of routing area, as well as logic area,



leads naturally to a reduction of routing wires. Shorter routing wires will be driven more efficiently, besides the fact that they will be controlled by drivers with improved output conductance. This makes the 3-D integrated FPGA potentially faster than its standard Bulk counterparts.

#### 4.4.4 Discussion

Monolithic 3-D integration appears to be an interesting technology from two aspects. First, it is currently reaching industrial maturity, which makes suitable for consideration for upcoming generations of reconfigurable circuits. In particular, it allows reductions in both the area of the global circuits and the critical path delay. The intrinsic performance of the FDSOI process is able to improve the speed of logic blocks, as well as the strength of the driver lines for routing resources. Furthermore, the integration on two layers clearly shows a reduction in the dimensions. Monolithic 3-D integration is able to reach a high via density and high alignment accuracy. This means that it is possible to interconnect the various layers with the finest granularity, and thus optimize the circuit area. Nevertheless, the bottom and top layers are highly correlated. In particular, the active regions are very close, which means that significant coupling is expected between the layers. However, this coupling has never been studied and must be investigated in terms with silicon experiments and models. Furthermore, the simplest 3-D process stacks the active silicon sequentially before processing the back end layers. This means that connectivity between the bottom transistors will impact the top transistors as well, and may lead to a non-optimum stacking of the transistors. This means that the 3-D design is not a simple stacking of cells designed separately, but must be done directly at the layout step. It is thus necessary to consider creating real 3-D cells, and optimized functions using the proposed scheme. From the technology perspective, it is also possible to integrate one layer of metal between the two active layers, and in the future, it will also be possible to integrate several layers between the active regions. This will improve the layout capabilities and lead to a smaller footprint. Further investigations will be required to find the optimal number of internal metal layers required for an efficient layout scheme.

### 4.5 Vertical NWFETs-based FPGA Performances

The vertical NWFET technology opens the way towards a smart routing back-end. This is obviously of high interest for FPGA applications, where a large number of switches are required to route the signals from one side of the circuit to another. The use of vertical active devices will lead to more compact routing resources.

#### 4.5.1 Real 3-D Routing for FPGAs

The technology is intended to improve the FPGA routing resources, by pushing them into the third dimension. Routing resources are realized by several sub-circuits. Connections between segments are proposed to be realized by the smart via scheme proposed in 3.4.3.1, and which intends to realize a configuration via between two metal layers. In segment connections, we will consider the simplest case where a single transistor is used between two metal lines. Switchbox cross point nodes will be realized as proposed in 3.4.3.3. Indeed, a 6-transistor scheme is used for the cross point, while the SRAMs are realized close to the node in a back-end logic scheme. Finally, it is worth noticing that every buffer required (for signals or for clocks), will be realized directly in the back-end, using the logic gates proposed in 3.4.3.2. This integration is very relevant to our considerations, since it appears inefficient to connect a signal every time to its front-end buffer, before coming back to the metal line.

In the proposed organization, the connection boxes and logic blocks remain in the front-end levels. Thus, we expect to distribute equally the complexity between the circuits implemented in front-end and those moved in back-end.

Figure 68 gives an overview of the final FPGA scheme, based on the Vertical NWFET technology.



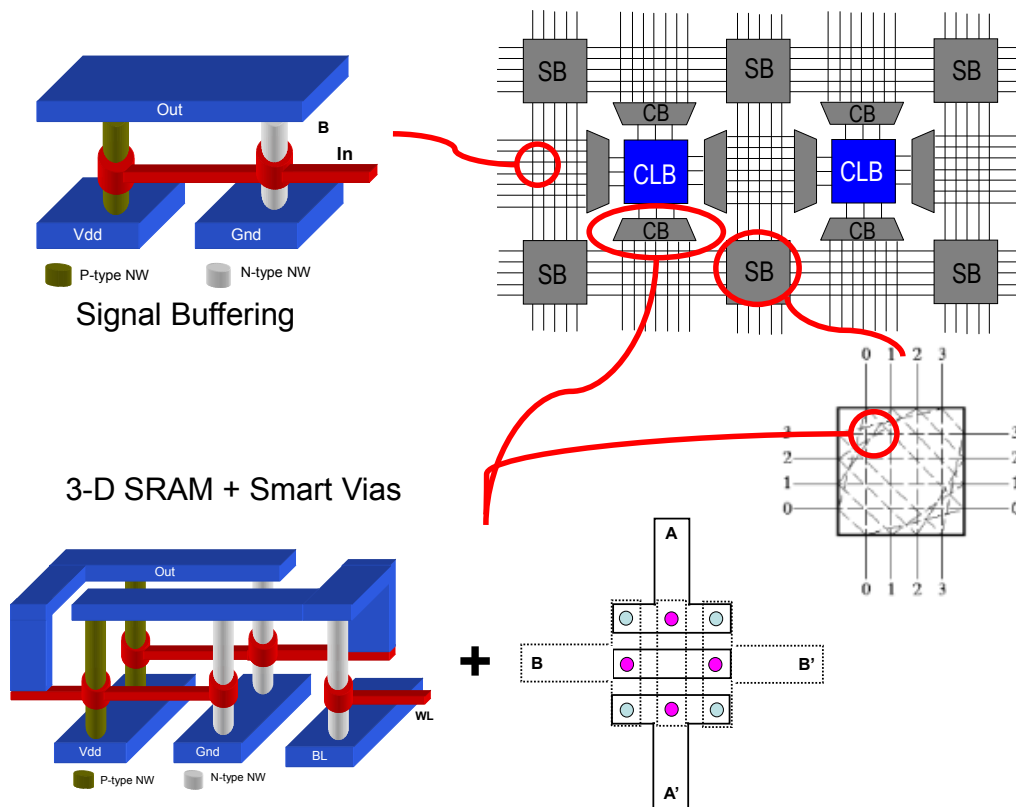


Figure 68. Vertical-NWFET-based FPGA overall organization

### 4.5.2 Methodology

In this evaluation, the methodology is again the same as that of 4.2.2 and 4.4.2. Since the technology “only” improves the routing resources, we adapt the architecture file accordingly. The performance of each elementary block in the back-end uses the numbers extracted from TCAD simulations directly.

### 4.5.3 Simulation of Large Circuits

Simulations have been carried out by mapping the benchmark to SRAM-Bulk and the vertical integrated FET scheme proposed above. We thus evaluate the area and the critical delay. We measured, relatively to 2-D CMOS FPGA, an area saving of about 46% on all our samples. This number is due to the decrease of the area allocated to routing. Indeed, the technology is able to place the pass-transistors inside the back-end layer with a very small front-end projection. The delay estimation is shown in figure 69. The benchmarks show a delay reduction ranging from 37% to 48%, with 42% on average.

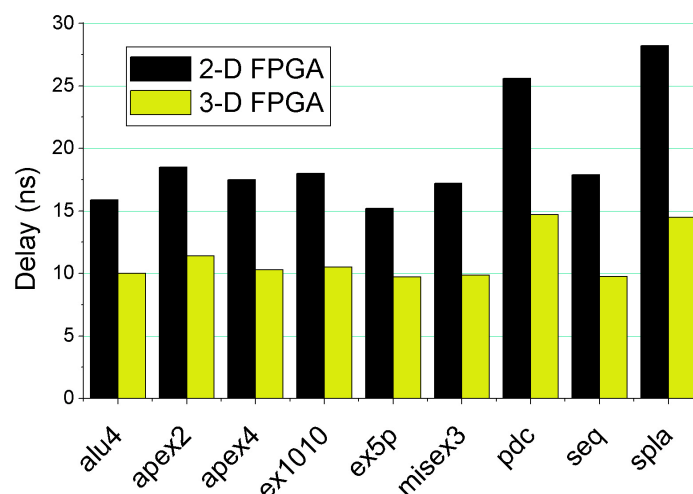


Figure 69. Delay estimation for 2-D and 3-D FPGAs

In fact, with the proposed integration process, it is possible to realize large transistors with very low front-end projection and with very good electrostatic channel control. In fact, the Gate-All-Around structure leads to the existence of a channel controllable from all points on the circumference of the wire. As such, for a small diameter  $d$ , it is possible to achieve a large effective transistor width ( $\pi d$ ).

#### 4.5.4 Discussions

The presented technology is suited for reconfigurable applications, since it yields significant reductions in circuit area and critical path delay. However, only large transistors ( $>100\text{nm}$  diameter) are currently manufacturable with this technique. While this first appears to be contrary to the general tendency of scaling, it is in fact not a constraint from a design perspective. In fact, large transistors exhibit good performance levels and can be used as high performance switches, while at the same time they can be realized with a low front-end silicon area. This leads directly to an improvement of area and delay.

The presented results are of interest both for circuit design and for technology. While the technology is not mature enough for mass production, we have used a predictive evaluation methodology to study its impact in a full application environment. The results achieved helps to conclude that it is not always necessary for technology to reach the most advanced scaled dimensions. It is thus preferable to converge through a reliable solution for mass production, in terms of variability, clean room compatibility ... Feeding back conclusions to the technologists is also a significant achievement for the methodology.

In this thesis, we focus on FPGAs. Nevertheless, it is worth noticing that this proposal should not be limited to this application only; complex 3-D ASICs can also be envisaged. Figure 70 shows how the technology could be used to stack several layers of standard cells and make them share the same space as routing layers. In this way, routing layers could then not be used entirely and include some areas where routing is not allowed. Back-end logic is a promising technology for numerous applications, but several questions remain open concerning EDA tools with 3-D technology. These questions have been heavily researched for 3-D TSVs [118]. It is potentially possible to make use of the available tools for coarse-grain circuits and extend them to the fine granularity of this approach.

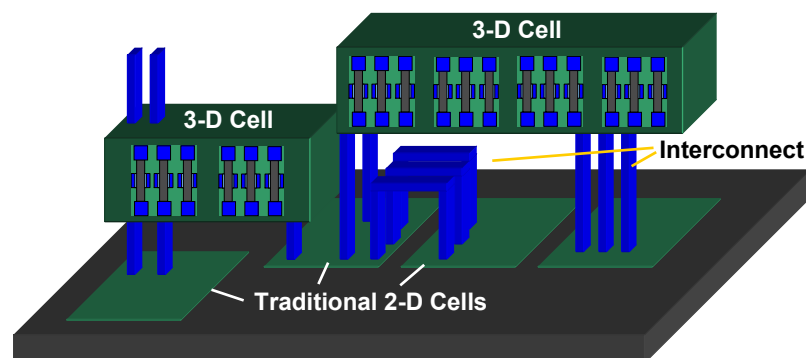


Figure 70. Illustration of the 3-D and 2-D cell co-integration

Finally, we could also expect to use other real 3-D technology to fulfill the role of a vertical switch. In this way, it is possible to use vertical *Nano-Electro Mechanical Systems* (NEMS), with nano-relays realized vertically [119]. This technology is even more promising than the FET one, thanks to the very good on/off resistances of NEMS.

## 4.6 Discussion

Table XII gives a summary of the gain observed through the use of the different technologies. We should note that the best results are obtained by the vertical NWFET technology for the area, with a gain of 46%, and by the PCM resistive memories for the critical path delay with a gain of 44%.

**Table XII. Architectural evaluation summary**

	<b>Total Area Gain</b>	<b>Critical Path Delay Gain</b>
<i>PCM</i>	13%	44%
<i>Monolithic 3-D</i>	21%	22%
<i>Vertical NWFET</i>	46%	42%

In terms of area, the vertical NWFET technology is the most compact in terms of implementation of routing resources. Indeed, the possibility to move most of the active volume into the third dimension leads to a large reduction in the front-end cell area. The impact on the whole circuit is thus of around 46%. The Monolithic 3-D integration process has a similar effect. Nevertheless, it only stacks layers of 2-D transistors. Thus, the impact of a cell is much greater than the real 3-D implantation. Furthermore, conversely to the monolithic 3-D process, the vertical FET integration is able to place the transistors above with a very small front-end impact, but it is also possible to choose the vertical placement of the transistors. This means that it is possible to address the question of metal line layout with this technique in a relatively simple way. In monolithic 3-D, some techniques can be considered, such as several intermediate metal lines, but these solutions still suffer from technological issues in terms of process temperature budget. Finally, the PCM gives only an improvement of 13%, which is due to the fact that all PCM circuitry needs to be addressed during its programming. While the PCM programming requires a large amount of current, the remaining access transistor is large, and leads to a loss in terms of gain. Other resistive memory technologies might be envisaged to overcome this limitation.

In terms of critical path delay, the best gain is achieved by the PCM implementation with a gain of 44%, while the vertical NWFET implementation is very close with 42%. Delay in FPGA data paths depends linearly on the on-resistance of the programmable switches, which means that the smaller the resistance, the better the gain in delay. The proposed PCM structure places a resistance memory directly into the data paths. Its on-state resistance is lower than standard silicon transistors, and this leads to the reduction of delay. Vertical NWFET transistors are naturally large, due to their realization in the metal layers. These large transistors exhibit good electrical properties, which makes them able to work as high performance switches, correlated to the on-state of the switch and good delay reduction. Finally, it is worth noticing that monolithic 3-D improves the delay by 22%. Even if this figure is lower than that of the two other technologies, the gain is significant. It is due to the good electrical properties of FDSOI, linked with the reduction of routing wires along the circuit.

Other metrics might be envisaged in the future. In particular, the analysis of power consumption of the whole proposed FPGA should be another comparison objective. However, this requires power models of the architecture and specific tools, such as in [120]. Nonetheless, we could imagine that a solution based on high performance complementary logic, such as vertical NWFET, will lead to the lowest power consumption. On the opposite, RRAM technology is purely passive and will contribute to leakage (impact of low  $R_{off}$ ) and resistive losses (in the data path).

## **4.7 Conclusion**

In this chapter, we explored the impact on the FPGA architecture of the technologies introduced in the previous chapter. The evaluated technologies are expected to improve the routing resources of FPGA, thanks to the use of the third dimension. In the context of a conventional FPGA architecture, we benched well-known circuit functionalities using the improved structures and compared the results with respect to traditional CMOS counterparts. We have seen that the best results are obtained by the PCM technology for the critical path delay with 44% in gain, thanks to its low on-resistance value. The vertical NWFET technology improves the area by about 46% and the critical path delay by 42% respectively.

These improvements are due to the extremely promising performance levels of this high-performance and fully 3-D technology. Finally, we should also remark that monolithic 3-D integration yields an improvement of 21% in area and 22% in critical path delay. This relatively mature technology thus represents a good trade-off for several performance metrics.



## **CHAPTER 5**    *Disruptive Logic Blocks*

---

### **Abstract**

In this chapter, emerging technologies will be used to create disruptive elements for Field Programmable Gate Arrays. We focus mainly on the combinational function blocks, in order to improve the computing performance of future reconfigurable systems. We propose to study the use of an ambipolar carbon electronics process and two different silicon nanowire crossbar processes.

Carbon electronics, and especially the Carbon Nanotube Field Effect Transistor, exhibits the property of ambipolarity, which means that  $n$ - and  $p$ -type behaviors are achievable within the same device. It thus becomes possible to obtain a device with tunable polarity, thanks to the addition of a second (polarity) gate to the device. This novel programmability of CNFETs is leveraged in a compact in-field reconfigurable logic gate and in a new approach to designing compact dynamic logic gates. We expect an improvement up to 3.1x in area, due to the cell compactness, and up to 2x in power consumption, due to the carbon-based device electrical properties.

We then propose the use of a sublithographic silicon nanowire crossbar process. This process is used to build a configurable logic cell with the same functionality as that of the carbon-based cell. This allows the two different structures to be compared. It is worth noticing that using the crossbar organization helps to compact the dimensions (up to 6x) required by the logic circuits. Nevertheless, a technological process build around a sublithographic arrangement of nanowires is highly unreliable, and its feasibility remains uncertain when considering all the access contacts. In order to correct the lack of manufacturability of the sublithographic crossbar process, we propose a variant on this crossbar process. This is realized on a modified Fully Depleted Silicon-On-Insulator process, and enables the construction of circuits in a crossbar scheme with lithographic dimensions. In this way, it is possible to build regular gates with an area reduction up to 6x. In this process however, the main hurdle is parasitic coupling. This could be tuned by the process, especially by oxide thickness. Using the correct tradeoff, the delay can be reduced by up to 1.3x, and we finally show that this crossbar scheme will surpass the performance of traditional CMOS schemes at the 16-nm node.

In the previous chapter, we focused on improvements of the routing and memory parts of the FPGA architecture. While these improvements are relevant to solve many issues in FPGAs, the proposals can still be considered only to be incremental improvements of the conventional structure. It is thus of interest to consider the limitations of FPGAs more carefully. Indeed, in an FPGA architecture, only 14% of the area is used for logic blocks, while the entire chip is occupied mainly by “peripheral circuitries” [38].

Logic blocks are the core of computation and thus the architecture could be considered to be inefficient with respect to the ratio between computing and periphery. We will therefore try to explore another direction for the architecture. In this chapter, we will work on the elementary logic block, as it is the foundation of the entire architecture.

## 5.1 Context and Objectives

As introduced, we focus in this chapter on logic blocks. The main objective is to provide new seeds for basic logic in reconfigurable architectures. These seeds will serve to build new architectural schemes in chapter 6.

To work in this direction, we intend for each block to be disruptive as compared to conventional technology. We define the disruptiveness as a gain of around 10x with a combination of area/performance/power metrics.

Two principal hypotheses/requirements must be formulated for this work.

Firstly, the disruptive blocks must be designed in an "architecture-friendly" manner, meaning that the blocks must be adapted to the architectural arrangement. This implies, for example, that a compact block which requires a large amount of extra circuitry is not suited for large applications. To illustrate this aspect, we consider diode logic, which needs signal restoration after each computing block. So even if the logic element is small compared to an equivalent MOS, the number of additional output buffers required to ensure correct logic behavior can lead to a significant reduction in the overall area gain. In the same way, logic blocks should minimize memory requirements since, similarly to the preceding discussion, additional memory leads to an increase in overhead. Further, the memory logic levels should also be suited to the logic.

Secondly, the design must be feasible from a technology point of view, and moreover be adapted to the strengths and weaknesses of the fabrication process. First of all, it seems quite fundamental that the envisaged technologies should remain compatible with standard CMOS processes. It is clear that semiconductor companies must amortize the cost of current existing processes and will not abandon their facilities and equipment. Due to this, any proposal must be compatible with clean room facilities by, for example, avoiding the use of gold as a catalyst. All technological assumptions in this thesis are thus driven by analyses from the LETI technology department.

Several emerging devices have been proposed during the last decade, which exploit various physical phenomena and properties. In this chapter, we consider two principal means of improvement of circuit performance: the use of improved device functionality and/or the use of increase in device integration density.

The first direction is based on the improvement of device functionality. Here we consider in particular *Double-Gate Carbon Nanotube Field Effect Transistors* (DG-CNFET), which can be configured to *n*-, *p*- or *off*-type devices depending on the polarity-gate voltage. Such device ambipolarity has previously been shown to be of benefit for complex logic designs [121, 122]. In this work, we refined the previous approach in terms of technology, device optimization and design. In particular, a credible technology process and device layout are proposed. Device optimization was performed through simulation in order to ensure correct hypotheses for the analysis, in close collaboration with IMS in the framework of the Nanograin project.

An extension of the back-gate use will also be presented in order to improve dynamic logic circuits.

The second direction is based on the improvement of elementary device integration density. Based on realistic assumptions, a regular process can be used to obtain high integration density, through a variety of different processes. We will examine the sublithographic integration of devices and, based on the NASIC approach [65], we will derive it to implement simple cells, rather than complex circuits. The integration of simple cells changes some paradigms of the NASIC approach, but is motivated from a technological point of view. In order to compare fairly the approach with elementary reconfigurable system gates, we use a multiplexer gate as baseline logic cell. The technological assumptions are mainly based on speculative bottom-up arrangements of sublithographic wires. Thus, in order to simplify the hypotheses, we will move in a second approach from crossbars to a sublithographic process integration. An innovative process, derived from *Fully Depleted Silicon-On-Insulator* (FDSOI), will be proposed. A layout methodology for logic cells using the technology will also be described. Finally, an optimization of the technology will be done by considering the global performance of the gate.

## 5.2 Proposal 1: Ambipolar Carbon Electronics

### 5.2.1 Introduction

Since the discovery of *Carbon Nanotubes* (CNT) in the early 1990s [123] and the first isolation of a graphene sheet in 2004 [124], carbon electronics has witnessed a growing interest. In particular, the study of the intrinsic properties of carbon and its potential application to microelectronic devices has been of high interest.

Carbon nanotubes, which usually consist of a single sheet of carbon rolled up to form a seamless tube [137], possess exceptional electrical properties such as high current carrying capability ( $>10^9 \text{A/cm}^2$ ) [138] and excellent carrier mobility ( $9000 \text{cm}^2/\text{V.s}$ ) [139]. Due to the small diameter ( $\sim 1 \text{nm}$ ), nanotubes are ideal candidates to provide one-dimensional (1-D) electrical transport. Ballistic transport, even at room temperature, has been demonstrated over short distances ( $\lambda_{\text{MFP}} \approx 700 \text{nm}$ ) [139]. It should however be noted that only carbon nanotubes built from a single layer of carbon (single-wall carbon nanotubes) can behave as semiconductors. Multiple-wall carbon nanotubes are composed of a number of coaxial single-wall nanotubes. High-performance electronics applications only use single-wall tubes, due to their small size and good semiconducting behavior.

Graphene is a monolayer of carbon atoms in a honey-comb lattice with unique electronic properties [125]. With a high saturation velocity ( $5.5 \times 10^7 \text{cm.s}^{-1}$ ) [126], graphene is considered to be a very promising candidate for high frequency applications. In addition, the ultra-thin body thickness of graphene offers ideal 2-D electronics channels for the ultimately scaled down device. Recently, cut-off frequencies in the gigahertz range have been demonstrated in top-gated graphene transistors built on exfoliated single-layer graphene sheets [127, 128], and on few-layer graphene grown on SiC substrates [130]. Epitaxial growth has been investigated in particular to produce wafer-scale, high-quality graphene [131]. Such techniques are compatible with all traditional lithographic patterning techniques to build the electronic circuits, following on from traditional circuit fabrication processes. Nevertheless, the graphene material is a semi-metal, meaning that it has no gap and thus no semiconducting behavior. One way to introduce a band gap is to confine electrons by patterning a 2D graphene sheet into a narrow ribbon ( $< 10 \text{nm}$ ), better known as a *Graphene NanoRibbon* (GNR) [132, 133]. Another way to modify the band structure of graphene is to stack two mono-layers to form a bi-layer which has a semiconducting band structure with zero band gap [134]. This band gap can be additionally tuned by creating a potential difference between the two layers [135, 136].



In addition to their promising performance levels, it is worth noticing that carbon technology is ambipolar. This means that a transistor can conduct for both positive and negative gate voltages. Traditionally, this kind of behavior is not suited for micro-electronics applications, where unipolar devices are required. While technologists are working on masking this property, it is also interesting to examine this property at the circuit level to assess opportunities with such ambipolar devices. In [141], a single carbon nanotube field effect transistor with a resistive pull-up to the power supply was used in order to implement a dynamic XOR logic gate. Further, ambipolarity can be controlled by adding a second control gate, such that it becomes possible to obtain a single functionality-improved device with a tunable polarity [142].

In this section, we will focus on carbon technology, and in particular on the control and efficient use of ambipolarity.

### 5.2.2 CNT-based Devices Properties and Opportunities

Most *Carbon Nanotubes Field Effect Transistors* (CNFETs) studied so far have adopted a back-gate top-contact geometry, as shown in figure 71-a. In this back-gate configuration, the nanotubes are dispersed or grown on a conducting substrate covered by an insulating layer. Two metal contacts are deposited on the nanotube to serve as source and drain electrodes, while the conducting substrate is the gate electrode in this three-terminal device. Figure 71-b shows the subthreshold characteristics (drain current measured as a function of gate voltage) of a back-gated CNFET for various drain voltages [142]. In this figure, it is worth noticing that, as the gate voltage increases from negative values, the current decreases, reaches a minimum value and then rises again. This exhibits an ambipolar transistor characteristic. In order to understand this result, figure 71-c depicts schematic band diagrams of a CNFET for negative and positive gate voltage respectively. Conversely to a conventional *Metal Oxide Semiconductor Field Effect Transistor* (MOSFET), switching of a CNFET is dominated by the modulation of Schottky barriers formed at the nanotube/metal contacts [143]. Since carbon nanotubes are intrinsically ambivalent, meaning that both electrons and holes can flow through the structure, the height of the Schottky barriers allows the selection of the carrier conduction. At sufficiently negative (resp. positive) gate voltage, the Schottky barrier is sufficiently thinned to enable hole (resp. electron) injection from the source (resp. drain) contact into the nanotube. In the behavior of a FET, these barriers are not desirable because they reduce the controllability of the device. Several works have been proposed to improve the contact quality. In [144], a graphitic wet layer is used to improve the contact between metal and nanotube.

In [142], another approach is taken. Instead of eliminating the ambipolarity of the devices, a dual-gate approach is introduced. The dual-gate creates pure *n*- and/or *p*-type devices with excellent off-state performance, simply by using electrostatic control on the carriers. Figure 72-a shows the scanning electron microscopy image, as well as the device cross section of a *Double-Gate CNFET* (DG-CNFET) structure. The DG-CNFET possesses an additional Al gate electrode placed underneath the nanotube between the source and drain contacts. The Al gate region is denoted as region B, and the area between the Al gate and the source/drain contacts are denoted as regions A. In the design, the Al gate is the primary gate that governs the electrostatics and the switching of the nanotube bulk channel in region B. The Schottky barriers at the nanotube/metal contacts are controlled by the Si back gate (substrate). This prevents the electrostatics in region A from being influenced by the Al gate. Figure 72-b shows the I-V characteristics of the DG-CNFET device. It is worth noticing that the same dual-gate CNFET exhibits clear *p*- and *n*-type unipolar properties for negative and positive  $V_{gs-Si}$ , respectively. In the device, electrostatic doping effects in CNFETs is utilized to eliminate ambipolar characteristics in a Schottky barrier CNFET and to obtain a bulk-switched transistor possessing a tunable polarity (*n* or *p*), steep subthreshold swing and excellent off-state performance. This *p*-FET and *n*-FET behaviour of the dual-gate CNFET can be understood by the schematic band diagrams shown in Figure 72-c.

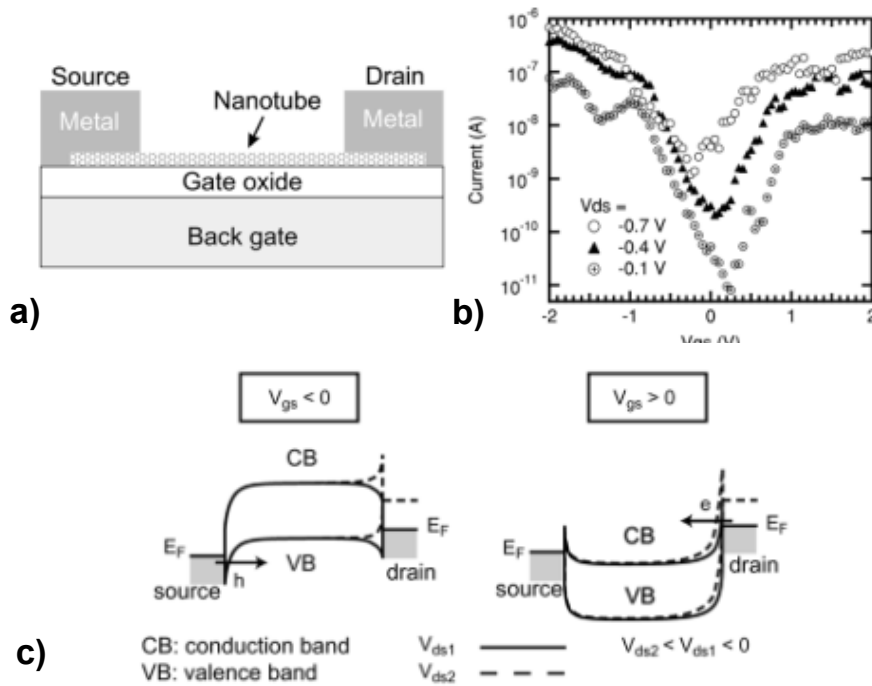


Figure 71. a) Schematics of a CNFET with a back-gate configuration [142] b) Measured characteristics of a typical CNFET (CNT  $\phi=1.4\text{nm}$ , Ti contacts and gate oxide 10-nm of  $\text{SiO}_2$ ) [142] c) Band diagrams of the structure [142]

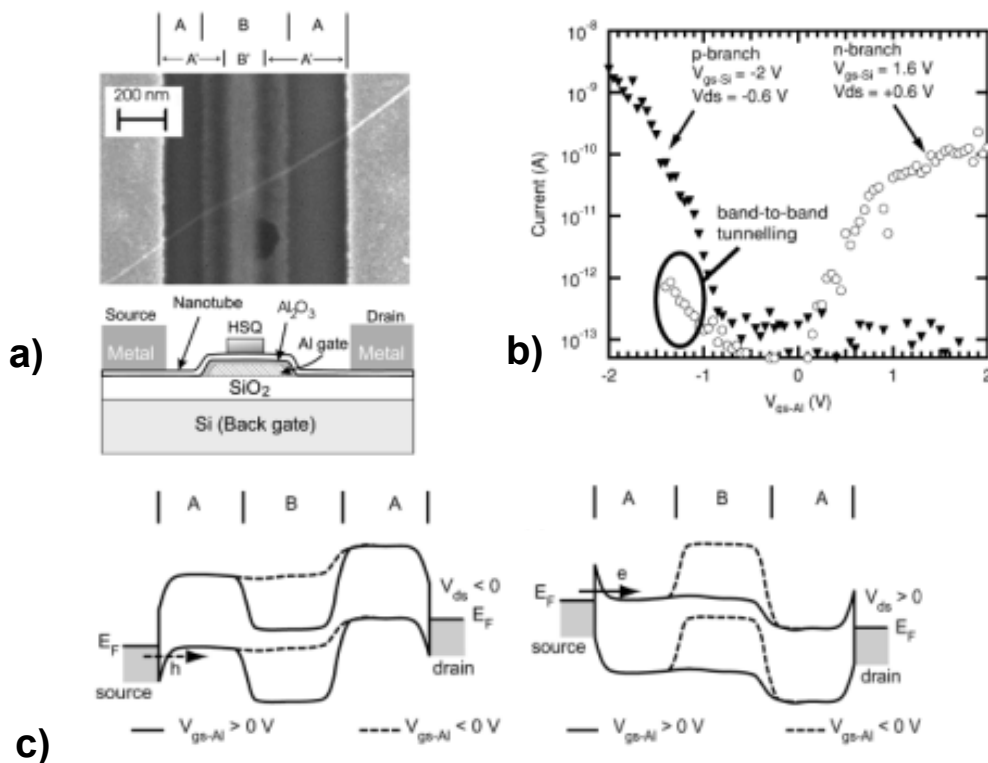


Figure 72. a) SEM image (top) and schematic cross-sectional diagram of a DG-CNTFET [142] b) Measured characteristics of a DG-CNTFET (Back-gates are polarized to shown p- and n-type unipolar behaviors) [142] c) Band diagrams of the structure (Left  $V_{gs-Si} < 0$ , and right  $V_{gs-Si} > 0$ ) [142]

For a sufficiently negative (resp. positive) Si gate voltage, the Schottky barriers are thin enough to allow for hole (resp. electron) tunnelling from the metal contacts into the nanotube channel. Thus, regions A become electrostatically doped as *p*-type (resp. *n*-type), resulting in a *p/i/p* (resp. *n/i/n*) band profile that allows only hole (resp. electron) transport in the nanotube channel. The dual-gate CNFET is switched *on* and *off* by varying the Al gate voltage that alters the barrier height for carrier transport across region B. In this configuration, regions A

serve as extended source and drain, and the device operates similarly to a conventional MOSFET through bulk-switching in region B.

### 5.2.3 Technological Assumptions and Device Modeling

In [142], the authors used an integration process based on a generalized and common back-gate electrode. This process is simple for research and characterization purposes. Nevertheless, it requires a global and shared back-gate for all the devices. However for circuit design, the principal advantage of the device is its unique in-field reconfigurability, meaning that each back-gate needs an individual control. Control individuality is required by several circuits, such as [52, 121]. Here, we will firstly describe the technological background compatible with the requirement of individual back-gate access. Then, we will discuss parameter fitting for the compact model, and process optimization required for design.

#### 5.2.3.1 Process Flow and Layout Proposal

The proposed process flow is based on *Silicon-On-Insulator* (SOI) wafers. In this process, it is possible to build silicon mesas (i.e. islands of silicon surrounded by oxide) to realize the DG-CNFETs back-gate. The individual control requirement will thus be guaranteed. The final device is shown in figure 73.

The expected technology process is based on a realistic and CMOS-compatible process flow. Starting with a SOI wafer, where the silicon-on-insulator layer will be used as back-gate,  $N^{++}$  doping (or NiSi salicidation) can be realized to ensure a good conductivity of the back electrodes.  $\text{SiO}_2$  is subsequently deposited as back-gate oxide, and intrinsic carbon nanotubes are transferred on top of this [145]. Then, the gate oxide ( $\text{HfO}_2$ ) and the metal (Al) of the top gate are deposited and patterned. Then, the active area and the back gate are defined by  $\text{SiO}_2$  and Si etching respectively. Subsequently, the metal is sputtered onto the contacts to drain, source and both gates.

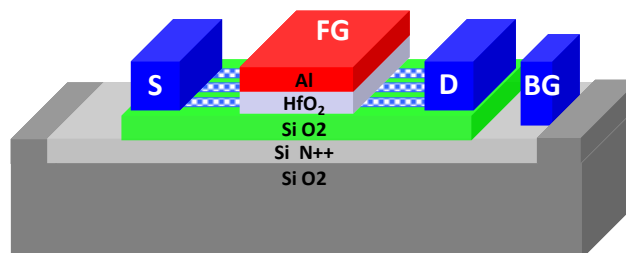


Figure 73. DG-CNFET device schematic using the proposed process-flow and showing the source (S), drain (D), front- (FG) and back-gate (BG) contacts

Using the previously defined process flow, we consider the layout requirements for our DG-CNFET device. The expected layout is shown in figure 74. This layout corresponds to parameterized-cell requirements and follows a standard industrial back-end rule set.

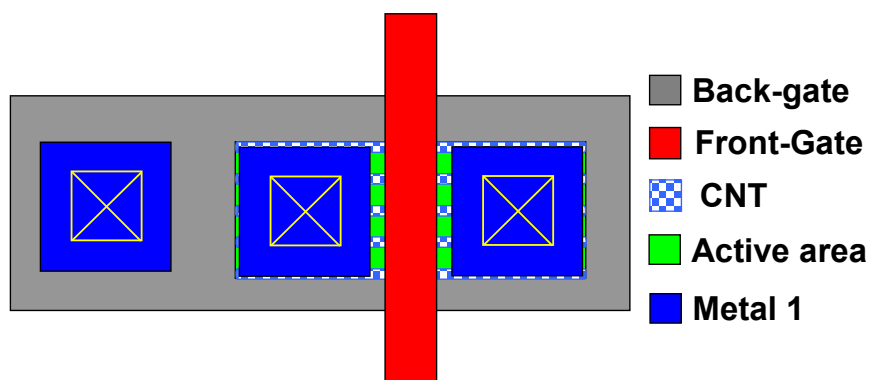


Figure 74. DG-CNFET device layout

### 5.2.3.2 *DG-CNFET Compact Model*

A physical compact model of the device is described in [146]. This model has been developed in the context of the Nanograin project by IMS. This physical model is made of three different regions: source access, inner part (underneath the front gate) and drain access.

In this structure, four energy barriers appear in the device. At the metal to source (or drain) access junction, two Schottky barriers appear, while at the source (or drain) access to the inner part junction, the barrier is more conventional and is a *pn*-junction type. Depending on the work function difference between the metal contact and the nanotube, carriers at the metal-nanotube interface encounter different heights of energy barrier: Carriers with energies above the Schottky barrier height reach the channel by thermionic emission. On the other hand, the probability of carriers with energies below the Schottky barrier height reaching the channel is defined by a transmission function describing the tunnel effect which can be calculated from Wentzel-Kramers-Brillouin (WKB) approximation.

To overcome the complexity of the WKB expression for compact modeling, an approximation based on works from [147] is applied. This effective barrier height model is described in [146]. The electron (hole) current is calculated through the Landauer equation, by integrating the energy from the dominating barrier to infinity. The dominating barrier position depends on the applied bias. In fact, the electron current is limited by three barriers: (i) the Schottky barrier from the source, (ii) the Schottky barrier from the drain and (iii) the conduction (valence) band of the inner part. The analytical expression of the drain current is given in [146]. In this model, several other features are included. On the one hand, a band-to-band tunneling model has been developed for MOSFET-like CNFETs in [148] and has been validated through Non Equilibrium Green Function simulation. On the other hand, charges have been modeled according to the ballistic assumption and the analytical expression of charge in each region is given in [142]. The potential calculation inside the device is given in [146].

Finally, it is worth noticing that the parasitic devices have been taken into account in the model. They have been calculated using the chosen layout (5.2.3.1).

### 5.2.3.3 *Process Tuning*

The I-V curve of the DG-CNFET device, shown in figure 72, demonstrates clear *p*- and *n*-type unipolar behaviour. Once the process flow is fixed, the following step is to optimize the process materials in such a way that a symmetric operation is guaranteed. The meaning of a symmetric operation is threefold. First, the  $I_{on}$  current of the *p*- and *n*-type devices should be within the same range (10-20% difference). Second, the back gate voltage range, which has been reported in the literature to be from -2 to +2V, should be scaled down to  $-V_{dd} - +V_{dd}$  (with  $V_{dd}$  typically between 0.7 to 1V) in order to be compatible with the front gate voltage range. Third, the  $I_{off}$  current must be sufficiently low to control the leakage of the device. Hence, we will optimize the process flow by band-gap engineering, in order to obtain the I-V shape, presented in figure 75.

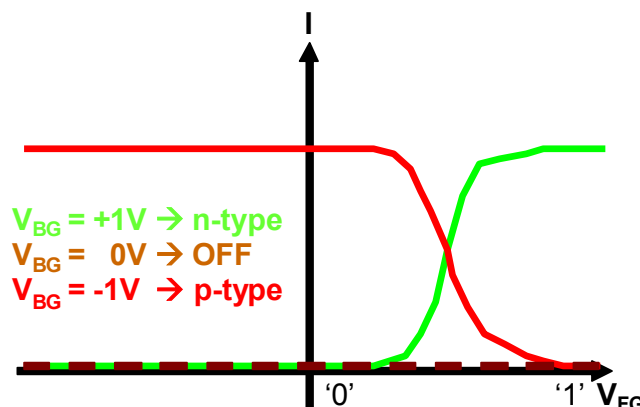


Figure 75. Specified I-V curve of a tuned DG-CNTFET

This optimization has been carried out by IMS by tuning the work function of the top and back gates. The gate materials (doped Si for the back gate and aluminum for the front gate), the gate oxides (SiO<sub>2</sub> for the back gate and HfO<sub>2</sub> for the front gate) and their thicknesses have thus been chosen in order to meet the conditions on the work function. Figure 76 demonstrates the operation as *p*- and *n*-type devices, after tuning, for a back gate voltage set to -1V and 1V respectively.

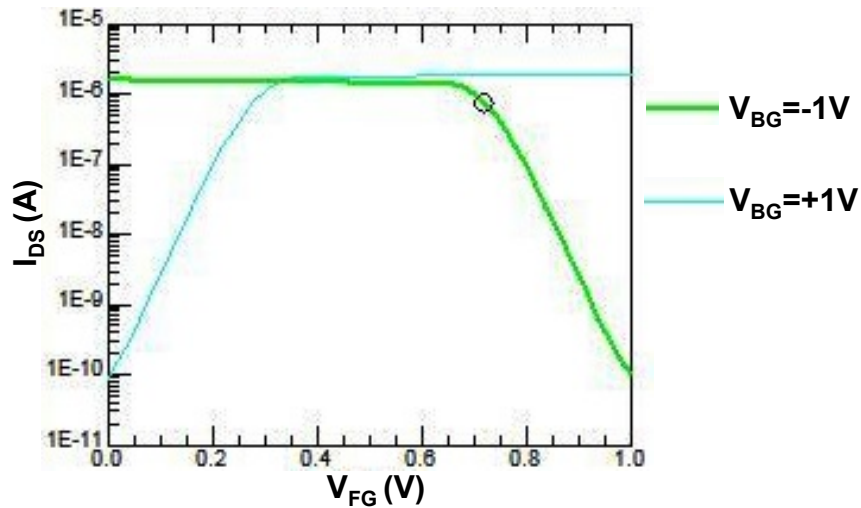


Figure 76. Simulated I-V curve of a tuned DG-CNTFET [IMS]

#### 5.2.4 Functionality improvement : In-Field Reconfigurability

The in-field reconfigurability property of the DG-CNTFET increases the functionality of circuits at the device level. It is thus possible to build compact logic cells using the in-field reconfigurability property. Such a cell was proposed by INL in the Nanograin project. In [121], a dynamic reconfigurable logic cell is presented, and the schematic is shown in figure 77-a. The cell is based on dynamic logic and is composed of seven DG-CNTFETs organized into two logic stages: logic function and follower/inverter. A four-phase clock signal set is used to perform the logic operation, and consists of two precharge inputs (PC<sub>1</sub>, PC<sub>2</sub>) and two evaluation inputs (EV<sub>1</sub>, EV<sub>2</sub>). The signals are non-overlapping as in classical CMOS dynamic logic gates. The polarities (*n*-type / *p*-type) of DG devices T<sub>1</sub>, T<sub>2</sub> and T<sub>3</sub> are controlled by the corresponding back-gate bias voltages V<sub>bA</sub>, V<sub>bB</sub> and V<sub>bC</sub>, as previously explained. The cell may thus be configured to one of fourteen basic binary operation modes, as shown in figure 77-b. In this cell, there are seven inputs and one output:

- two Boolean data inputs A and B (The logic levels are represented by the supply voltage values V<sub>dd</sub> and Gnd);
- three control inputs to configure the circuit according to figure 77-b (Back-gate bias voltages are according to the tuned control values);
- The clocking signals PC<sub>1</sub>, PC<sub>2</sub>, EV<sub>1</sub> and EV<sub>2</sub>;
- The circuit output Y.

Figure 78 illustrates how this logic gate works. When V<sub>bA</sub> = V<sub>bB</sub> = V<sub>bC</sub> = 1V, CNTFETs T<sub>1</sub>, T<sub>2</sub> and T<sub>3</sub> are all configured as *n*-type FETs. When PC<sub>1</sub> is enabled, the first stage is pre-charged, and the voltage of the internal node C (V<sub>c</sub>) is discharged to 0 V. If for example either of the data inputs A or B = logic '1', then when EV<sub>1</sub> is enabled, the first layer evaluates its output such that the internal node C is set to logic '1'. Then PC<sub>2</sub> is enabled (pre-charge of the second stage), and the output Y is charged to logic '1'; and when EV<sub>2</sub> is enabled, the output Y is evaluated to logic '0'. In fact in this configuration, the only situation where C is not set to logic '1' and Y therefore evaluates to logic '1' (since T<sub>3</sub> is off) is when both A and B = logic '0'. This means that for V<sub>bA</sub> = V<sub>bB</sub> = V<sub>bC</sub> = 1 V, the cell is configured as a NOR gate.

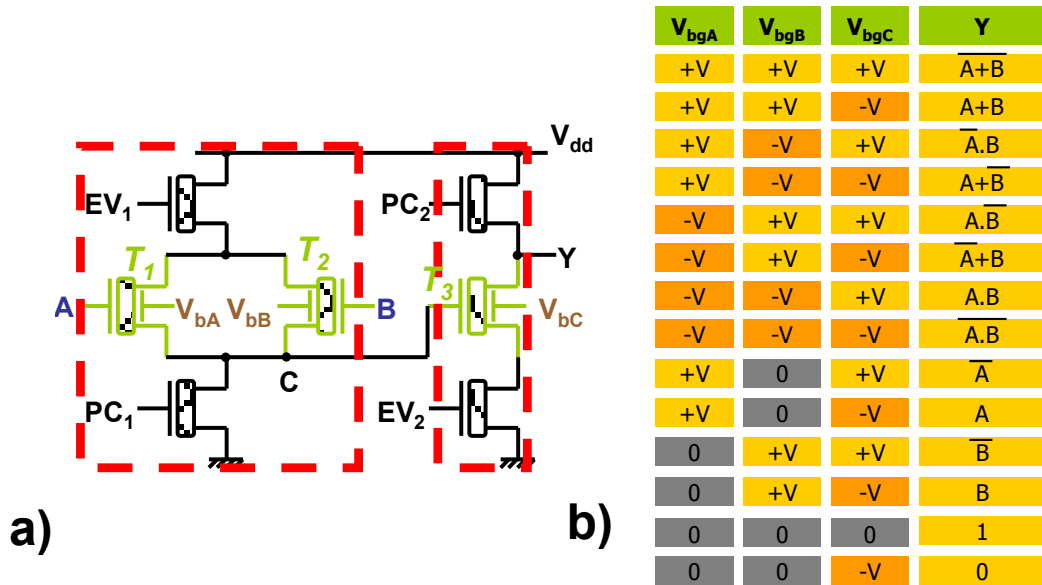


Figure 77. a) Transistor level schematic and b) configuration table for CNT reconfigurable cell [122]

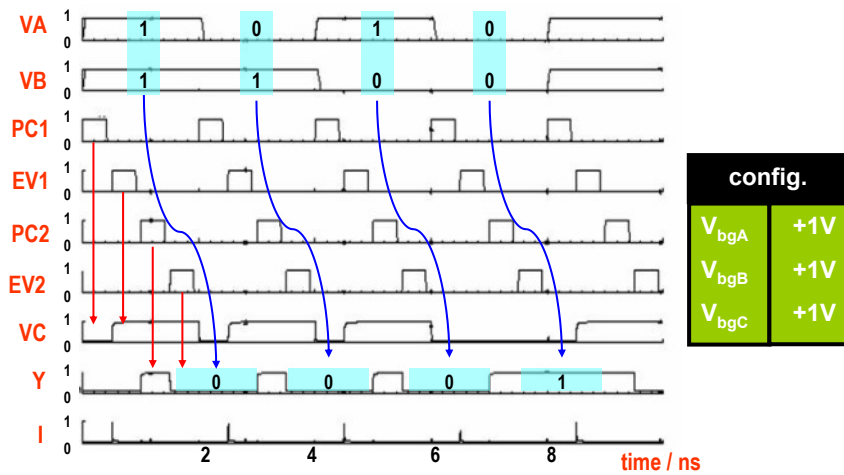


Figure 78. Reconfigurable logic gate waveforms in NOR configuration [122]

### 5.2.5 Performance Evaluation

The cell presented above promises circuit compactness for reconfigurable applications. We will now assess the performance of the in-field reconfigurable gate and compare it to the equivalent CMOS counterpart.

#### 5.2.5.1 DG-CNFET Evaluation Methodology

In order to evaluate the performance of the DG-CNFET reconfigurable logic cell, we investigate the area, delay and power consumption. The area will be extracted considering its layout in a 22-nm lithographic node. The choice of this extrapolated node is due to the fact that this solution has been estimated for middle-term. Since the logic cell is dynamically clocked, the maximal performance is controlled by the maximal clock frequency reachable. The performance of DG-CNFET cell was estimated in [150] through the use of a simple RC-model. However, the extracted performance metrics did not take the DG-CNTFET intrinsic parasitics into account. Hence, parasitic values were extracted from the layout and the RC model has been corrected accordingly. Simulations have also been conducted using different load values.

The power consumption numbers are extracted from [150].The power consumption is computed by electrical simulations with a *fan-out-of-four* (FO4) load.

### 5.2.5.2 Silicon CMOS Performance Evaluation Methodology

In order to evaluate in an objective manner the performance of the advanced DG-CNFET technology, we will consider the use of an advanced lithography node of 22-nm (i.e. the expected equivalent silicon technology node that will be available when CNT technology will be mature enough for industry). To perform this evaluation, we will extract the metric values from a well-established 65-nm technology and extrapolate it, as explained hereafter.

Considering reconfigurable logic architectures, the cell is compared with a standard LUT4:1 as used in actual FPGA systems. A standard MUX4:1 is considered to implement the equivalent CMOS circuit. The performance metrics of the MUX4:1 will be used to illustrate the methodology.

- **Area**

Area is measured using the layout view, and then extrapolated to the target 22nm technology node. The scaling factor used is

$$K_{Area}(x \text{ nm} \rightarrow y \text{ nm}) = \frac{x}{y}$$

In the 65nm→22nm case (i.e. 3 technological generations following a scaling factor between generations of  $\sqrt{2}$ ):

$$K_{Area} \approx 2.95 \approx (\sqrt{2})^3$$

and area is scaled down according to

$$K_{Area}^{-2}$$

Hence, since the area of a CMOS MUX cell in 65nm is  $10.4\mu\text{m}^2$ , the equivalent cell in 22nm will be  $1.19\mu\text{m}^2$ .

- **Performance**

To evaluate the maximal achievable performance, the intrinsic delay and load influence factor of the operator must be extracted with respect to the parasitic devices.

Estimation of the delay in CMOS MUX is done using a 65nm standard cell, and the data are then extrapolated to the 22-nm node. The scaling factor used, based on data taken from the ITRS [6], is:

$$K_{Delay}(x \text{ nm} \rightarrow y \text{ nm}) = \frac{ITRS_{On-Chip \text{ Local Clock}}(x \text{ nm})}{ITRS_{On-Chip \text{ Local Clock}}(y \text{ nm})}$$

In the 65nm→22nm case, the factor used is

$$K_{Delay} \approx 1.95,$$

and delay is scaled with

$$K_{Delay}^{-1}.$$

Figure 79 shows the estimated performance of 65-nm and 22-nm CMOS MUX.

- **Power consumption**

The predictive evaluation of power consumption evaluation is a difficult issue, due to its dependency on the target load, the associated architecture, configuration details and frequency. An FO4 load has been chosen for the purposes of evaluation. The procedure for the evaluation of CMOS needs extrapolation.



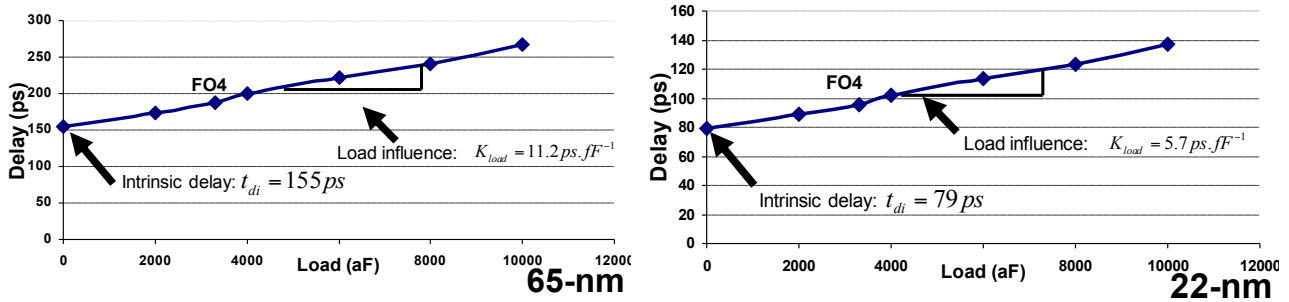


Figure 79. Delay analysis of a CMOS MUX in 65-nm and 22-nm

First, an evaluation is performed, through simulation, using a 65-nm standard cell at a frequency compatible with this technology and normalized with respect to the frequency.

$$P_{Hz}(x \text{ nm}) = \frac{P_{Eval}(x \text{ nm})}{f(x \text{ nm})}$$

Then, this normalized power figure is scaled to the extrapolated node.

$$P_{Hz}(y \text{ nm}) = \frac{P_{Hz}(x \text{ nm})}{K_{Power}(x \text{ nm} \rightarrow y \text{ nm})}$$

The determination of the scaling factor is quite complex, due to the lack of data on power evolution. In fact, power consumption of a circuit can be split into 3 contributions: the dynamic power ( $\alpha.f.C.V_{dd}^2$ ), the static power ( $V_{dd}.I_{leak}$ ) and the short-circuit power. These contributions follow different trends which do not evolve according to the Moore's Law anymore. Nevertheless, while it is clearly impossible to consider that this factor is close to the area scaling factor for such an advanced technology node, some degree of correlation can be found with the frequency scaling factor, evaluated using the ITRS. We thus (somewhat arbitrarily) assume that

$$K_{Power} \approx K_{Delay} \approx 1.95$$

Finally, denormalization to the targeted frequency is computed

$$P(y \text{ nm}) = P_{Hz}(x \text{ nm}).f(y \text{ nm})$$

### 5.2.5.3 Simulation Results

Table XIII summarizes this comparative study. The estimated size of a CMOS MUX at the 22-nm node is  $1.19 \mu\text{m}^2$ , while the DG-CNFET cell area is  $0.39 \mu\text{m}^2$ . This gain of 3.1x can be explained by the significant reduction in the number of transistors, from 26 for a standard MUX to only 7 with the in-field reconfigurable cell.

Table XIII. Global evaluation of the DG-CNFET cell performances

	Area ( $\mu\text{m}^2$ )	Intrinsic delay (ps)	$K_{Load}$ (ps.fF <sup>-1</sup> )	Average power at 4 GHz ( $\mu\text{W}$ )
<i>MUX MOS</i>	1.191	79	5.7	3.56
<i>DG-CNFET</i>	0.39	149	124.5	1.78
<i>Gain</i>	x 3.1	x 0.53	x 0.04	x 2

The timing evaluation of the cell is depicted in figure 80, with all the parasitics considered. For this metric, the DG-CNFET cell demonstrates lower performance values than the CMOS equivalent cell. This is counter-intuitive, considering the intrinsic properties of carbon electronics: (i) that the good electrical transport of carbon should increase the device speed, and (ii) that the size reduction of the elementary cell should reduce the parasitics. While both factors should lead to an improvement in the performance metrics, the comparison in fact considers a static (CMOS) logic cell and a dynamic (DG-CNFET) logic cell. The comparison is thus not completely valid; however it is justified by the fact that dynamic logic is used to



compact the cell and to compete with standard FPGAs. Nevertheless, this does lead to an increase time delay, due to the four clock phases used to charge and discharge the capacitive nodes.

This could be improved by using transistors with much higher drive strength. In fact, in the simulation, the DG-CNFET transistor is built with a single nanotube as the channel material. In this context, it seems reasonable to envisage the use of multiple CNTs per transistor, where each CNT contributes to the transistor channel. Using multi-channel transistors will result in a higher  $I_{on}$  value, which accelerates the charge and discharge of all the gated stages (with of course corresponding impact on power consumption figures).

Table XIV shows the power estimation figures for the CMOS MUX and the in-field reconfigurable cell, for all achievable functions. We can observe that the cell improves the power by an average of 2x compared to CMOS. Regarding these two opposing results, we must consider the *Power-Delay Product* (PDP). The PDP is improved by only 6% compared to CMOS technology. In fact, the power consumption is essentially due to the dynamic power contribution. Hence, it depends mainly on the frequency and the load (and parasitic capacitances). Due to the lithographic size of the transistor, the various capacitances of the devices are mostly the same, and this leads to an equivalent PDP.

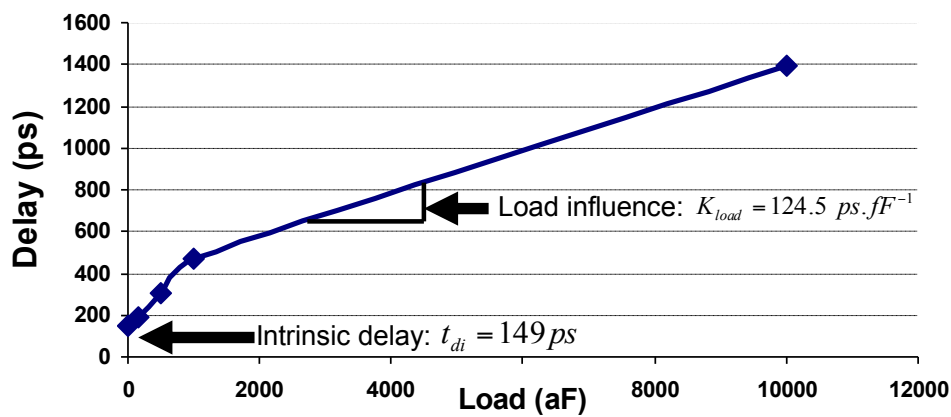


Figure 80. Delay analysis of DG-CNTFET Reconfigurable Cell

Table XIV. Power consumption vs. function (22-nm node – 4GHz)

Power in $\mu\text{W}$	CMOS MUX	DG-CNTFET [150]	Gain
$\overline{A+B}$	3.80	1.87	2.03x
1	0.93	1.12	0.83x
$A+B$	4.04	1.85	2.18x
$\overline{A}$	3.85	1.83	2.10x
$A$	3.63	1.81	2.00x
$\overline{A.B}$	3.95	1.86	2.12x
$A+\overline{B}$	3.71	1.84	2.02x
$\overline{B}$	4.23	1.82	2.32x
$B$	4.39	1.79	2.45x
0	0.89	1.82	0.49x
$\overline{A.B}$	4.12	1.84	2.24x
$\overline{A+B}$	4.10	1.82	2.25x
$A.B$	4.01	1.84	2.18x
$\overline{A.B}$	4.19	1.82	2.30x
$A\oplus B$	7.30	-	-
$\overline{A\oplus B}$	6.87	-	-
<i>Average</i> (without Exclusive functions)	3.56	1.78	2x

### 5.2.6 Dynamic Logic Circuit Improvement

The DG-CNFET ambivalence property has been exploited in previous work to add some reconfigurability to logic cells, while transistor count is maintained low. In this section, we propose to extend the use of the back-gate reconfigurability and its associated states (*n*, *p*, *off*) to improve the structure of dynamic logic cells.

#### 5.2.6.1 Concept

Compact logic cells can be realized using dynamic-logic [151]. This allows, as in conventional CMOS technology, to reduce the number of transistors (and therefore the complexity) by a factor of almost two for each function, with respect to static-logic implementation. This reduction can be explained by the elimination of the *p*-type pull-up network and replacement around the pull-down network by two clocked transistors for precharge and evaluation. The generic scheme for dynamic logic is shown in figure 81-a.

The DG-CNFET *off*-state is used to reduce complexity. The *off*-state corresponds to a transistor that is never *on* whatever the state of the front-gate is. We will use this property to merge the evaluation transistor directly with the function paths.

The programmable polarity allows us to eliminate all input inversion circuitry. Effectively, the polarity of the transistor might be chosen in order to select whether the front-gate signal should be considered directly or as its complementary value. The generic structure is shown in figure 81-b. The information of interest is the Pull-up (resp. Pull-down for the case of the inverted function) network with the evaluation clock merged on the back gates. The function path structure remains the same as in traditional CMOS, and consists of several DG-CNFETs. The DG-CNFETs could be placed in serial branches (for AND functions), parallel branches (for OR functions) or any combination of serial and parallel branches. The evaluation clock signal controls the back-gate of at least one DG-CNFET for each path from  $V_{dd}$  to the Output node. The considered DG-CNFETs receive non-complemented inputs on their front gates. In the case of complemented inputs, DG-CNFETs have to be configured to *p*-type using the associated voltage on their back-gate.

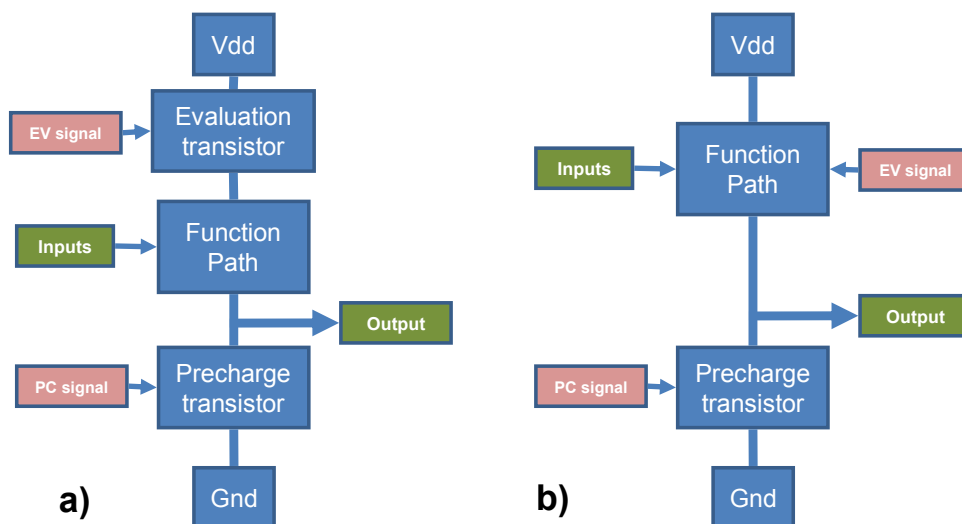


Figure 81. a) Generalized scheme for CMOS dynamic logic b) Generalized scheme for double-gate-based cells

#### 5.2.6.2 Buffer Cell Illustration

Figure 82-a shows an illustration of the concept with a buffer implementation. The cell operation is based on two DG-CNFETs. The precharge transistor (PC-gate, ground path) is configured to *n*-type, via the  $V_{dd}$  voltage applied to its back-gate. The evaluation and signal transistor (IN-gate,  $V_{dd}$  path) are combined; the evaluation clock signal is connected to the back-gate and the input signal is connected to the front-gate.

Figure 82-b shows the associated waveform for the cell. To explain the behavior of the cell, three phases are clearly identifiable:

- **Precharge:** PC=1, EV=0. The *n*-type bottom transistor is *on*, the top transistor is *off*. The output node is thus precharged to 0.
- **Sleep:** PC=0, EV=0. Both the *n*-type bottom transistor and the top transistor are *off*. Any existing charge is maintained on the output node.
- **Evaluation:** PC=0, EV=1. The *n*-type bottom transistor is *off*, and the top transistor is in the *n*-state, allowing evaluation. In this configuration, only the input signal has an influence on the output. If in=0, then the transistor is off and the output node is unchanged (from precharge state 0). If in=1, then the transistor is on and output node is charged to 1.

The resulting behavior is a buffer operator working in dynamic logic. It is worth noticing that this architecture may be used to implement all the standard logic operations: figure 82-c presents a XOR gate implementation as another example.

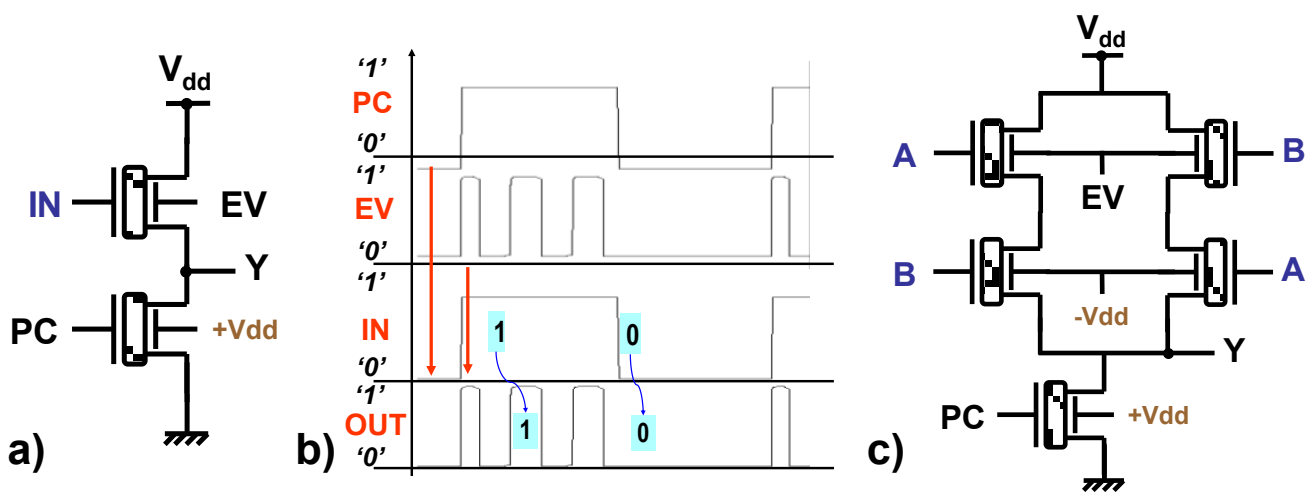


Figure 82. a) Buffer cell schematic b) associated waveform and c) XOR cell

### 5.2.6.3 Density Improvement Example

Table XV shows a comparison between the proposed approach and its CMOS counterparts in dynamic logic implementation. The proposed implementation, using the DG-CNFET device, reduces the number of transistors up to 50% compared to the dynamic MOS implementation in a standard cell ASIC approach.

Table XV. Dynamic DGCNFET-based cell transistor requirements

AREA	Static MOS cell	Dynamic MOS cell	DG-CNFET cell	Gain (DGCNFET vs. Dynamic MOS)
Buffer	4T	3T	2T	33%
Inverter	2T	3T	2T	33%
AND	6T	4T	3T	25%
OR	6T	4T	3T	25%
XOR	12T	10T	5T	50%

### 5.2.7 Discussion

All the results found above are highly promising from a design perspective. We observe that the functionality improvement of the device reduces the number of used transistors and leads to a more compact logic. Furthermore, the good expected performance levels of carbon electronics improve the power numbers, and certainly the delay also, after some structural improvements. This compactness opens the way towards ultra-fine grain computation, i.e. where the elementary logic blocks will be much smaller than conventional blocks and peripheral circuitries. This will be useful for new architectural paradigms.

Nevertheless, we should also consider that the logic compactness has a high cost in terms of clock and configuration signals: at least four non-overlapping clocks and three ternary memory nodes for the proposed cell. This makes the use of the cell quite difficult from an architectural point of view. Improvements of the cell using static logic implementations are currently under investigation by INL.

Finally, we should notice that these figures must be considered carefully. Indeed, these parameters have been taken directly from literature and have not been re-simulated using the final version of the compact model described in 5.2.3.2. Again, this evaluation work, as well as the structural validation at the circuit level, is currently in progress at INL.

### 5.3 Proposal 2: 1-D Silicon Crossbars

In the preceding chapter, we have investigated a technology that improves the intrinsic functionality of the device and saw how it can lead to a very compact logic cell. In this part, we will focus on the use of a density-improved technology. The dense integration of devices is expected to lead intrinsically to the realization of very compact logic circuits.

#### 5.3.1 Introduction

For several decades, the micro-electronics industry has been known to scale the transistor dimensions, in order to improve the performance. Nevertheless, this increase of performance is not only due to the scaling, but also due to several structural improvements.

Figure 83 shows the evolution of the silicon structures from bulk to nanowires. 1-D structures demonstrate good electrostatic controllability, which is fundamental to ensure good performance levels for scaled devices. While the trend in industry is towards the *Fully Depleted Silicon-On-Insulator* (FDSOI) process, we can note that nanowires tend to the optimal solution, merging a good electrostatic control with a low *Off* current ( $I_{off} < 1\text{nA}/\mu\text{m}$ ) [152].

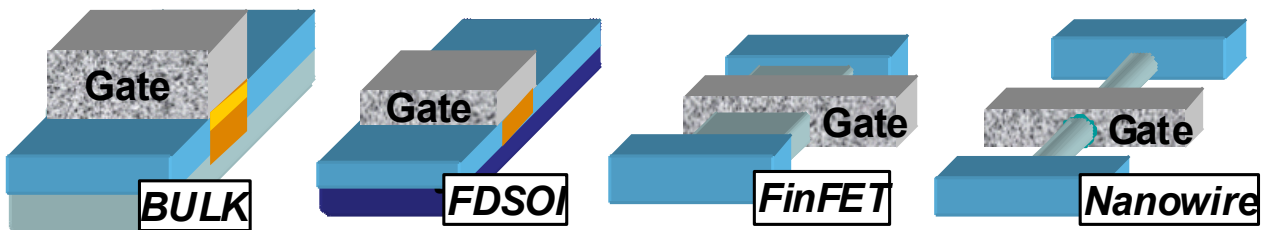
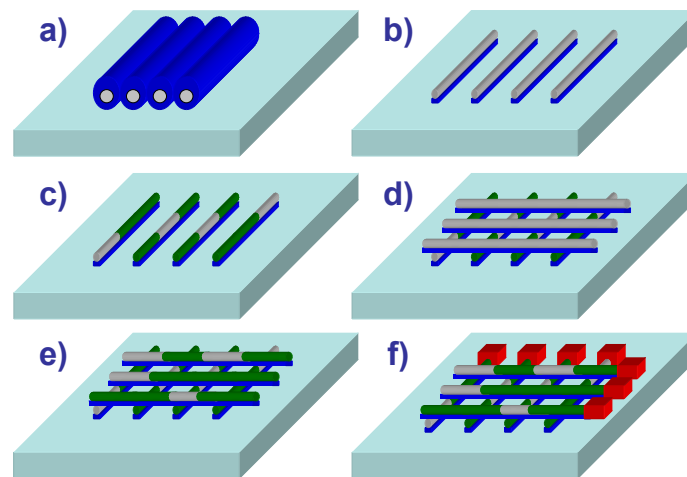


Figure 83. Silicon electronics evolution from bulk to nanowires

Furthermore, in addition to their good expected properties, we can note that the 1-D structures make the active dimensions very small. This leads to potentially high-density integrated devices, especially when a crossbar organization is proposed. Thanks to bottom-up fabrication processes, nanowires are useful to build regular crossbar structures at sublithographic scales and open the way towards enormous device density improvements over CMOS. In the literature, different active elements have been envisaged at the cross points, such as p-n junctions [153], molecular programmable switches/diodes [154] or FETs [155]. Several architectures have been proposed, based on these different devices. In [154] molecular diode-switches at crossbar intersections are proposed. This structure forms a diode logic grid, performing Programmable Logic Array. However, diode logic requires level-restoring circuitry and addressing of individual points, which leads to complex interfaces. Considering these limitations, solutions based on FET logic, as introduced in [156], are explored. The NASIC extends this concept in [157] by creating double-stage combinational logic on Crossbars of NanoWires FETs (CB-NWFETs). This structure is used as a general fabric to implement logic elements dedicated to nanoprocessors. We will use this approach to implement a logic operator with very compact dimensions.

### 5.3.2 Technological Assumptions

The crossbar can be manufactured with silicon nanowires grown by Chemical Vapor Deposition, using metallic nanoparticles as catalyst and the Vapor-Liquid-Solid mechanism [158]. This technique allows the achievement of nanowires with diameters controlled by the catalyst size [159] and diameter values around 3 nm [160]. *In situ* doping during nanowire growth can be used to obtain n-type or p-type nanowires [161]. Nanowires can subsequently be thermally oxidized to obtain a core-shell structure [162] and deposited on a substrate (Figure 84-a). The well-known Langmuir-Blodgett technique [163] can be used to align the nanowires and to obtain a first layer of nanowires with spacing controlled by the oxide thickness of the shell [162]. For better alignment when nanowires have very small diameters (i.e. below 3 nm), their lengths cannot be longer than a few micrometers to avoid gradual bending observed on nanowires with these diameters [160]. Nanowire pitch has also to be equal to (or greater than) the corresponding photolithography half pitch of a given technology node plus the nanowire radius for subsequent photolithography steps. The oxide shell should be removed except along the bottom contact between the nanowires and the substrate, in order to avoid removing the nanowires at the same time (Figure 84-b). Salicidation of some parts of the nanowires can be achieved using nickel (or platinum) physical vapour deposition, photolithography and etching to define the regions where the silicon nanowire will be the channel of transistors, annealing to form NiSi (or PtSi) regions and chemical removal of the unreacted deposited metal [155] (Figure 84-c). The second array of nanowires (crossing the first) can be obtained with the same process sequence (Figure 84-d). The oxide shell along the bottom of the second layer of nanowires serves as the gate dielectric and the NiSi (or PtSi) region serves as the metallic gate in the MOSFET structure (Figure 84-e). In the same way, some nanowire regions of the second layer can serve as the channel of the MOSFET, controlled by the nanowire metallic regions of the first layer. The process ends with a final metallization to contact all nanowires around the crossbar area by a conventional sequence of deposition, photolithography and etching steps (Figure 84-f).



**Figure 84.** Manufacturing process to build a nanowire crossbar. *a)* Deposition of silicon nanowire after growth, oxidation and Langmuir Blodgett alignment. *b)* Etching of the oxide shell. *c)* Salicidation of the nanowire regions which will not be transistor channels. *d)* Deposition, alignment and etching of the oxide shell of the second layer. *e)* Salicidation as step *c)*. *f)* metallization to contact the nanowires around the crossbar.

### 5.3.3 SiNWFET Configurable Logic Cell

#### 5.3.3.1 Nanowire Crossbar Logic

The previously presented process flow is able to build arrays of nanowires. As suggested in [156], FETs are built at the cross points. It is then possible to build elementary combinational logic functions, such as NAND, NOT or Buffer, with a dynamic-logic design style. An example of a dynamic NOT function is given in figure 85. Precharge and evaluation transistors are used, driven by two non-overlapping clock signals. The associated waveform has been obtained by simulation using an elementary model, taken from [156]. The logic

stage is first precharged to  $V_{dd}$  followed by a conditional evaluation with respect to the inputs. Thus, it is possible to build every logic function. Nanowire crossbar logic has been widely envisaged at the architectural level. In [65], the concept is generalized with the *Nanoscale Application Specific Integrated Circuit* (NASIC). The NASIC is an architecture intended to realize complex and specific designs within a lithographic crossbar framework.

### 5.3.3.2 Multiplexer Design Methodology

To ensure an equivalent behavior between this approach and the previous one, we propose the use of a 4:1 *multiplexer* (MUX). The considered configurable logic cell is designed to fulfill sixteen basic binary operations, at any given time. The selection of operation (i.e. the configuration of the cell) is made through the use of binary configuration signals. In fact, this corresponds to work with the functionality of a Look-Up Table (LUT), with no considerations on memories, since we focus on the logic part.

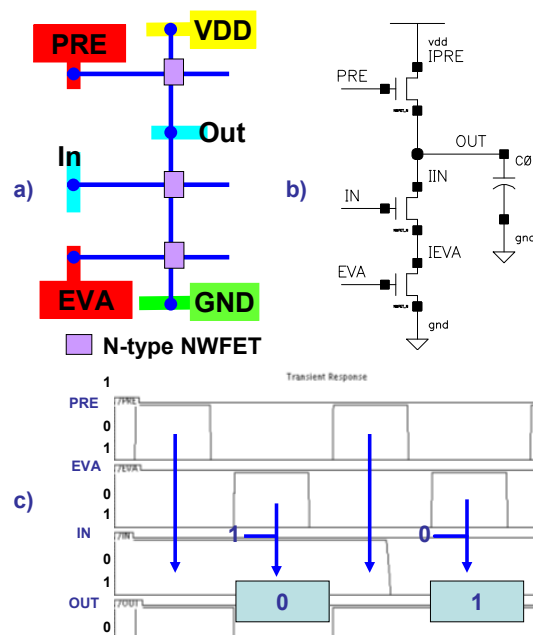


Figure 85. NOT function realized on *n*-type nanowires using dynamic logic. a) Pseudo-physical view b) Schematic view c) Associated waveform.

The result is a logic equation which can be implemented in canonical  $\sum \prod(I, \bar{I})$  form (*I* represents inputs). The MUX logic function is expressed and simplified through the Espresso logic minimizer tool [164].

A crossbar using only *n*-type NWFETs is able to implement NAND-NAND logic easily [157]. Considering the following property

$$\sum \prod(I, \bar{I}) = \overline{\overline{\overline{\prod \overline{\overline{\overline{\prod(I, \bar{I})}}}}}}}$$

any logic operator could be implemented in two NAND stages. The NASIC approach implements NAND/NAND logic for realizing complex circuit implementations, such as microprocessors (WISP-0) [165]. Due to the complexity, several NASIC tiles are connected. The connection is realized with nanowires.

In our implementation of MUX, we cannot use nanowires for interconnections between MUX and other cells. This assumption comes from the technology. In fact, it is not possible to realize wire alignment when the aspect ratio is too high (i.e. their lengths cannot be longer than a few micrometers to avoid gradual bending observed on nanowires with small diameters [160]). Thus, we expect that long connections will be implemented by micro-scale wires. In summary, the NASIC approach is used to build elementary gates in our approach. To minimize inherently area-hungry connection requirements at the micro-scale, the external contacts for complementary inputs are rendered superfluous through the use of an embedded



inverter stage which is placed before the logic. A buffer stage is also coupled in order to ensure the data synchronization required in dynamic logic.

Figure 86 shows a representation of the logic cell. The structure is formed by cascading three stages:

Input  $\rightarrow$  NOT/Buffer  $\rightarrow$  NAND  $\rightarrow$  NAND  $\rightarrow$  Output

The cell uses only two sets of dynamic clocks in order to reduce the number and complexity of clock wires. The three pipelined stages operate sequentially using these two clock phases: the NOT/Buffer and output NAND stages operate in phase, whereas the internal NAND stage operates on the opposite phase, ensuring correct charge transfer through the structure.

Due to the crossbar structure and the fabrication processes, the NAND gates are not all identical. For the first NAND stage, the inputs are y-axis nanowires, and outputs are x-axis nanowires, while the second NAND stage is exactly the opposite. Thus, due to the structure, transistors must be formed on two levels.

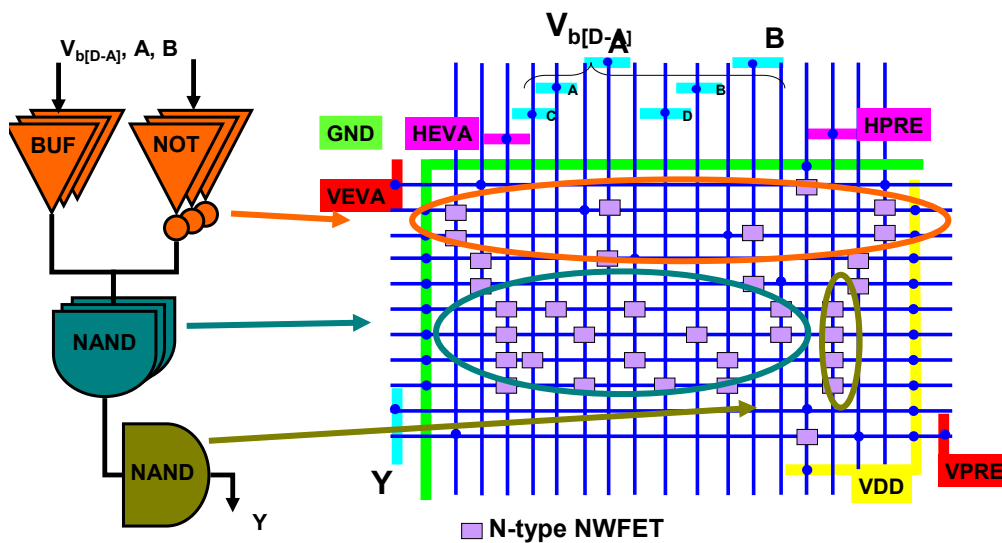


Figure 86. Representation of CB-NWFET dynamic reconfigurable logic cell

### 5.3.3.3 Layout Design

Small crossbars are connected directly to the microscale, in order to interconnect cells using typical back-end technologies. Nanowires are contacted to typical metal interconnections. It is thus possible to proceed with back-end design and subsequent Design Rule Check. An example layout, which respects an industrial back-end Design Rule set, is shown in figure 87.

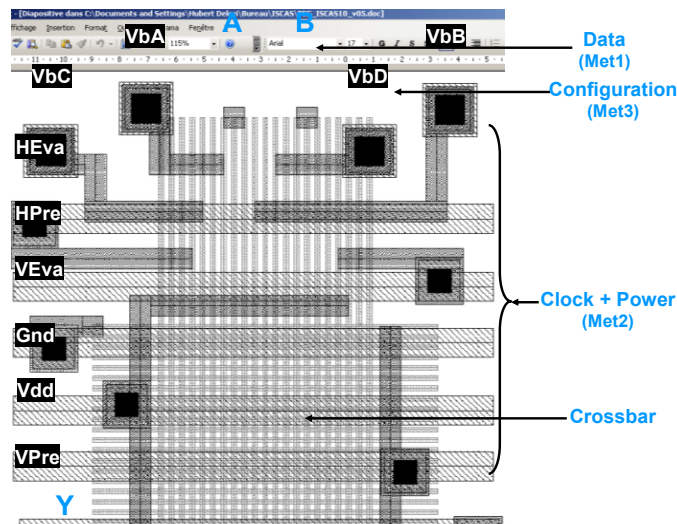


Figure 87. Metal/via layout of the dynamic logic cell

All service signals (clocks and power) are placed on horizontal lines to be shared with other cells. The configuration signals ( $V_{bA}$ ,  $V_{bB}$  ...) can be found on large vias to upper metal layers. This is of particular interest for cell addressing or back-end memory implementation. Finally, the data signals (A, B and Y) are available on lower level pads to realize local interconnections with neighboring cells.

### 5.3.4 Performance Evaluations

#### 5.3.4.1 Methodology

In order to evaluate the performance of this approach, we investigate the area, delay and power consumption. To ensure a fair comparison with the previously presented approach, we will use the same methodology and hypotheses. The area will be extracted from the proposed layout in the case of a 22-nm lithographic node. Since the NWFET cell is based on dynamic logic, the delay measurement will determine the minimum possible on-time pulse for the clock signals. Parasitic capacitance values have been extracted using the layout implementation and simulations have been conducted under various load value conditions. The elementary model of FETs presented in [156] has been used with the Eldo2008.2a simulator.

#### 5.3.4.2 Electrical Simulations

Table XVI summarizes this comparative study. Due to the sublithographic dimensions, the solution based on CB-NWFET gives an improvement in density of 4.1x, compared to the extrapolated CMOS.

Table XVI. Global evaluation of the CB-NWFET MUX performances

	Area ( $\mu\text{m}^2$ )	Intrinsic delay (ps)	$K_{\text{Load}}$ ( $\text{ps}\cdot\text{fF}^{-1}$ )	Average power at 4 GHz ( $\mu\text{W}$ )
<i>MUX MOS</i>	1.191	79	5.7	3
<i>CB-NWFET</i>	0.289	17	15.1	6.7
<i>Gain</i>	x 4.1	x 4.6	x 0.38	x 0.45

Figure 88 shows the delay analysis of the sublithographic multiplexer circuit. We can see that the intrinsic delay is close to 17ps, while the  $K_{\text{Load}}$  factor is close to  $15.1\text{ps}\cdot\text{fF}^{-1}$ . By comparison with the data extracted from CMOS technology, shown in figure 79, we can observe an improvement of 4.6x on the intrinsic delay, while the load factor is degraded by 0.38x. The intrinsic delay depends on the internal data paths. Due to the sublithographic dimensions, it has been possible to reduce the internal parasitic element values, thereby improving the delay figures. Nevertheless, while the internal delay is improved by the small dimensions, it is worth noticing that the circuits driving the output are also small. This means that the output buffer is quite inefficient for driving the load capacitance, and thus explains why the load factor is decreased conversely to CMOS equivalent. The output drivers could be improved by parallelizing the paths to increase the *on*-current.

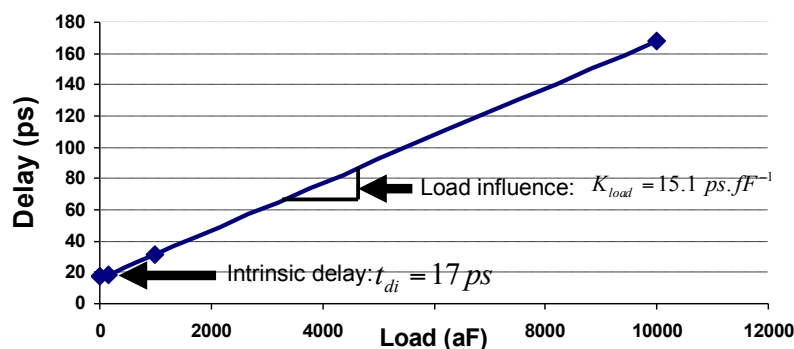


Figure 88. Delay analysis of Nanowire Reconfigurable Cell

Power consumption figures are reported in Table XVII, and are clearly worse than those of CMOS by, on average, a factor of 2.2x. It is worth noticing that the comparison is made



between a dynamic and a static logic cell, which is not a favorable case for the former. For example, in the case of function 0, only leakage current contributes to power consumption for the static cell, whereas the dynamic implementation charges and discharges nodes continuously. Another point to bear in mind is the model used. This model, with a simple variable resistance, is not favorable for NWFET results, due to a large resulting leakage current. Hence, it will be of high interest to re-run the power figure evaluation using a more realistic model, such as the one proposed in [167, 168].

**Table XVII. Detailed Power Consumption figures for sublithographic cell (22-nm – 4GHz)**

Power in $\mu\text{W}$	CMOS MUX	CB-NWFET	Gain
$\overline{A+B}$	3.80	9.50	0.40x
1	0.93	2.50	0.37x
$A+B$	4.04	4.25	0.95x
$\overline{A}$	3.85	6.87	0.56x
$A$	3.63	5.98	0.61x
$\overline{A.B}$	3.95	8.63	0.46x
$A+\overline{B}$	3.71	4.22	0.88x
$\overline{B}$	4.23	6.87	0.62x
$B$	4.39	6.89	0.64x
0	0.89	11.7	0.08x
$A.\overline{B}$	4.12	8.63	0.48x
$\overline{A+B}$	4.10	5.13	0.80x
$A.B$	4.01	9.06	0.44x
$\overline{A.B}$	4.19	4.23	0.99x
$A\oplus B$	7.30	5.99	1.22x
$\overline{A\oplus B}$	6.87	6.56	1.05x
<i>Average</i>	4.00	6.69	0.60x

#### 5.3.4.3 Discussion

As shown by the results, the solution is of high interest for area improvement and intrinsic delay reduction. The sublithographic organization helps to increase drastically the integration density of logic elements, while the sublithographic dimensions also reduce the parasitic contributions. The internal stages of the structure are fast, due to this reduction in parasitics, and the intrinsic delay is consequently improved.

Nevertheless, the contribution of sublithographic dimensions does have disadvantages. It is worth noticing that, even if the computation part is made more compact, several connections to external signals are required. This assumption is motivated by the technology and the difficulty to implement long nanowire interconnections. Due to this, metal lines follow lithographic dimensions, such that connections from the microscale to the nanoscale occupy a large part of the circuit.

Furthermore, as regards performance levels, the output drivers are built with sublithographic devices, which mean that they have a low drive capability. This can be seen from the load factor, which is high conversely to CMOS.

Finally, the main obstacle of this solution remains the technological process flow. Circuit fabrication is based on a bottom-up organization of nanowires. This process is currently in the early stages of its development, which makes the proposed solution highly long-term and speculative.

## 5.4 Proposal 3: Lithographic Crossbars

The previous proposal exploits the large density of active elements, built with a sublithographic integration process. Density-based technologies are obviously well-suited to large area reduction. Nevertheless, they are often based on highly speculative technology assumptions. In this section, we will propose a novel integration process for dense crossbar arrangements, based on lithographic FDSOI. Layout and specific methodology will be described, along with circuit performance-driven optimization of the technology.

### 5.4.1 Introduction

Previous research works addressing the crossbar architecture generally considered bottom-up nanowire fabrication techniques. The low maturity of the technology implies several limitations on the assessment of the architecture. On the one hand, bottom-up fabrication techniques yield dispersed nanowires with highly variable position and spacing. Thus, mean values with respect to the geometry have to be assumed in order to characterize the architecture. On the other hand, the ability to define two layers of perpendicular nanowires is a very challenging technological task. To date, it has been demonstrated only for metallic nanowires in a top-down approach [170], while silicon-based nanowire crossbars are dispersion-based, micrometer-scale demonstrators without any functionalization of the cross points [166]. Capacitive coupling between parallel nanowires has been addressed only in routing applications [169].

In this proposal, we will assess the questions related to the technology limitations and to the impact of the coupling effects between semi-conducting nanowires. Our approach is based on the choice of a fully characterized industrial technology. We propose a process flow to fabricate lithography-based crossbars with FETs as active devices at the cross points. Then, we electrically simulate the circuit in order to assess its performance and the potential loss due to the capacitive coupling and the resistance of the long nanowires.

### 5.4.2 Technological Assumptions

Using a *Fully-Depleted Silicon-on-Insulator* (FDSOI) process, wires can be manufactured with ultra-regular lines as demonstrated in [171]. In this section, we conceptually complete the already established process for a single layer of parallel nanowires, with a perpendicular top layer of parallel nanowires. In the proposed process, the bottom-most nanowires are defined using photo-lithography at the lithographic pitch, where their dimensions can be controlled through oxidation and etching below the lithographic limit down to 15nm [172]. Thereby, both n- and p-type dopings are allowed. On the other hand, the topmost lines are defined as *polycrystalline silicon* (poly-Si) stripes at the lithographic scale. These two perpendicular layers of parallel lines form a crossbar whereby the intersections are called cross points. In such a crossbar, the top lines can electrostatically control the nanowires underneath at the cross points in a FET fashion, when the ladders are covered by a gate oxide. Moreover, the top nanowires can form an ohmic contact to those lying underneath when a via is defined at the cross point.

Figure 89 shows the associated process flow. A P-type SOI substrate is patterned by lithography to form parallel ridges that are subsequently etched into nanowires. *Plasma Doping* (PLAD) is used to softly define N-type wires [173]. Then, the nanowires are passivated in oxide and planarized (Figure 89-a). Following this step, the passive regions, *i.e.*, the parts of the nanowires connecting every series FET, are defined by n- and p-type PLAD on the p- and n-type nanowires respectively (Figure 89-b). It is worth noticing that the implantation step is performed softly because of the small dimensions of the nanowires, such that dopant migration is limited. Moreover, the nanowires are separated by an oxide, further limiting the diffusion of dopants. This reduces the requirements on spacing, which are generally included in the design rules. This allows the smallest lithographic dimensions for all operations of patterning and doping to be reached. Then, the gate stack is defined by depositing the gate insulator, followed by the poly-silicon gate deposition and etching steps

(Figure 89-c). The poly-silicon lines carrying the gates are defined with regular parallel lines. At this level, the active devices are defined and the east-west connections between them are established through the passive parts of the nanowires, operating as resistances.

The north-south connections are composed of the poly-silicon lines and require the definition of vias between them, as well as the nanowires underneath them. The vias are defined by etching the poly-silicon lines and filling them with metal. In order to decrease the resistance of the north-south poly-silicon lines and the passive parts of the east-west silicon NW lines, it is possible to sputter a thin layer of nickel (or platinum) over the whole structure, which diffuses into the silicon and poly-silicon and forms a low resistance silicon silicide (Figure 89-d). For this reason, it is important to first etch the oxide covering the passive regions before sputtering the metal. The remaining metal after the diffusion can be removed by wet etching [156]. Finally, the contacts between the crossbar and the outer circuit are implemented through conventional metallization steps (Figure 89-e).

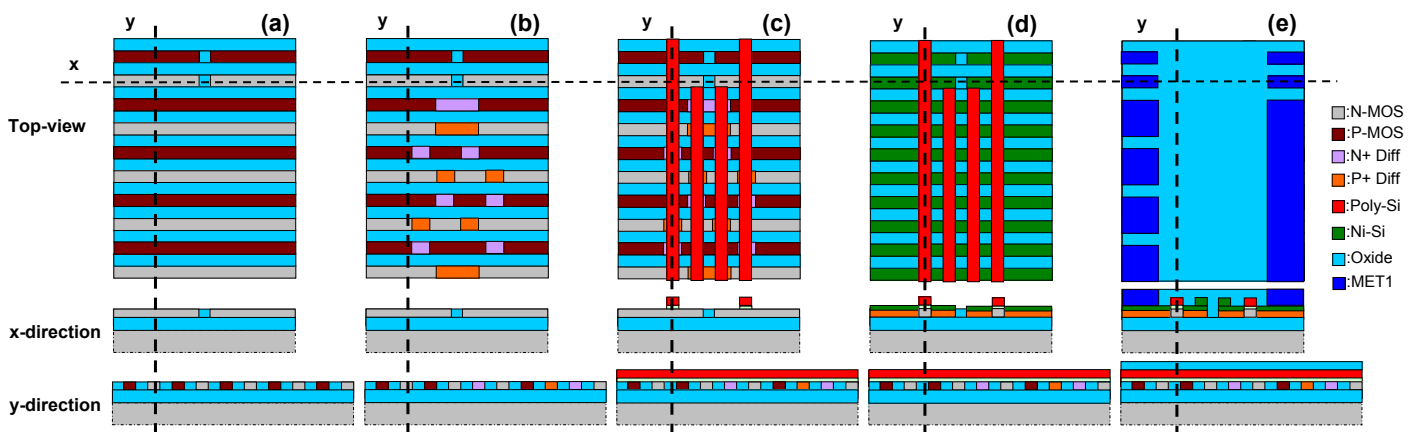


Figure 89. Manufacturing process to build a FDSOI crossbar: a) grating patterning and active regions doping. b) Passive regions definition. c) Gate deposit. d) Passive regions finalization and salicidation. e) Metallization to contact passive regions.

### 5.4.3 Logic Design Methodology

The FDSOI crossbar technology presented in the previous section can be used to implement different logic gates with a very high active area density.

#### 5.4.3.1 Inverter construction

An inverter can be fabricated for instance as depicted in figure 90 with very dense active areas. In this example, *p*-type and *n*-type transistors are realized with two separate wires, but aligned to efficiently share the gate on the same poly-Si line. It is worth noticing that the width of *p*-type transistors could be different from the width of *n*-type transistors, since separate lines are used for *n*- and *p*-Type channels. However, in a context of crossbars with ultra-regular layouts, it is preferable to keep all dimensions equal and at the lithographic limit ( $L_N=L_P=W_N=W_P=F$ ), despite the clear impact on skewed delay.

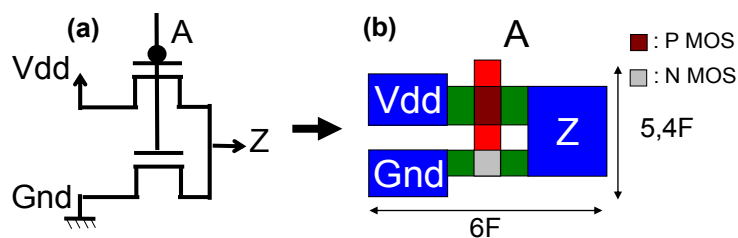


Figure 90. Crossbar inverter structure (a) equivalent circuit and (b) layout

#### 5.4.3.2 Generalization to complex logic circuits

Using the above strategy, it is possible to build any conventional logic function by folding complementary branches around the same contact pad. In the proposed technology, doping of adjacent lines is controlled only by lithography. It is possible to alternate the *n*- and *p*-type

dopings at adjacent wires or group of lines. This is useful for design purposes, when it is required to make islands of  $p$ - and  $n$ -transistors or to distribute them regularly.

- **Separate Doping Regions**

This patterning style consists of grouping  $p$ -type and  $n$ -type regions separately. This is close to the traditional representation of CMOS circuits and suitable for circuits where the output is a common pin for all branches. This situation can be found in most of the small logic circuits. An example is shown in figure 91 and illustrates a stand-alone 4-to-1 MUX. A stand-alone MUX can be realized as depicted, with the equivalent schematic in figure 91. The crossbar layout implementation is realized by associating a wire with each branch of the internal stages, whereby a branch is a sequence of transistors connecting the power line to the output buffer stage. The particularity of this multiplexer is that its data inputs control only transistor gates, thus avoiding any signal losses, while the output is buffered to keep the integrity of logic functions independent of the fan-out.

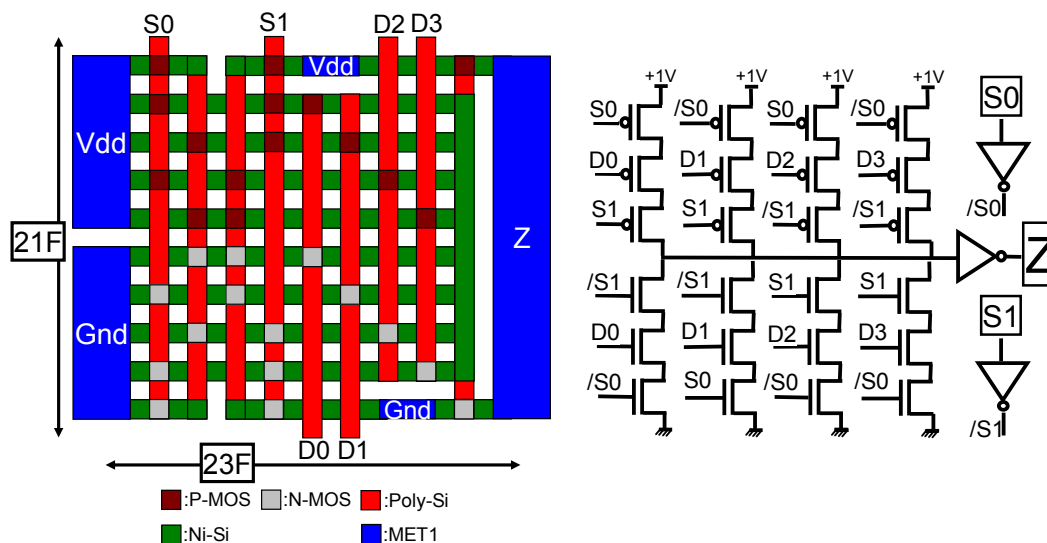


Figure 91. Stand-alone 4:1 MUX crossbar implementation (separate doping regions) and equivalent schematic

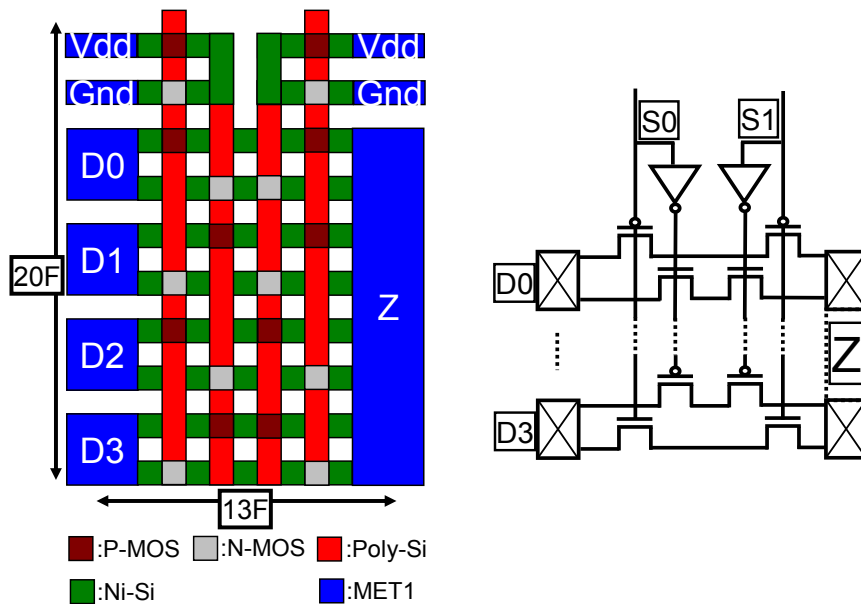
- **Alternating Doping Regions**

This patterning style consists of alternating  $p$ - and  $n$ -type regions, and is particularly suited to build pass-gates between two contact pads. This technique is used in an FPGA-adapted 4-to-1 multiplexer, as depicted in figure 92. Pass-gates have been placed between the inputs and the output of the MUX. The addressing is realized using poly-Si lines, which control all transistors. Control inputs are complemented by a primary inverter stage.

## 5.4.4 Performance Evaluation

### 5.4.4.1 Methodology

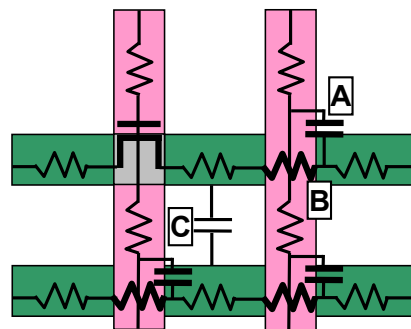
In this section, we study the impact of the crossbar technology based on the example of the designed MUX, by evaluating its area, delay and power consumption. This kind of technology can be expected as a near-term solution. In fact, all the technological processes currently exist and only the crossbar organization is fundamentally specific. This means that we can use standard industrial process data for the evaluation. The stand-alone crossbar MUX is thus compared to the equivalent MUX taken from industrial 65-nm design kit. The area is extracted considering the layouts in a 65-nm lithographic node. The delay is extracted from electrical simulation and decomposed into the intrinsic delay and  $K_{load}$  factor, which gives the load delay. The model used in the simulation is a scaled FDSOI transistor card for the PSP model. The power consumption is extracted from electrical simulations with a FO4 load. The circuit is operated at the maximum frequency achievable by the gate with the FO4 load. At this frequency, we swept all possible input vector combinations and we averaged the power consumption.



**Figure 92. FPGA-suited 4:1 MUX crossbar implementation (alternating doping regions) and equivalent schematic**

First, we simulated these metrics for the parasitic-free circuit, and then we introduced the parasitics and assessed the designed circuit under different conditions in order to minimize the impact of those parasitics.

Due to their very compact layout, a large number of parasitic devices have to be considered in the crossbar layout. A zoom on the parasitics at and around a cross point is shown in figure 93. We distinguish between two types of cross points: the *active cross points* operating as field-effect transistors, with a poly-Si line electrostatically controlling the Si line underneath it; and the *passive cross points* that are formed by a poly-Si line crossing a passivated highly doped Si line, without forming any FET. In our process proposal, the passive cross points are formed by a highly-doped conductive Si wire crossing a salicide poly-Si line. The two lines are separated solely by a thin layer of gate dielectric. Thus, the areas under the poly-Si lines, which are not used as a FET channel, result in a high parasitic capacitance between the two wires (see element A in figure 93). Moreover, we modelled the parasitic resistance through the doped conductive areas (see element B in figure 93), as well as the electrostatic inter-wire coupling using the capacitive element C in figure 93.



**Figure 93. Modeling of parasitic devices**

A parametric study was carried out to study the influence of the following design and technology parameters: *width of horizontal N-type wires* ( $F_{HN}$ ), *width of horizontal P-type wires* ( $F_{HP}$ ), *width of vertical poly-Si wires* ( $F_V$ ), *width of insulator lines between adjacent wires* ( $F_I$ ). All possible widths are multiples of the lithographic poly-Si half-pitch  $F$ .

#### 5.4.4.2 Performance Estimations with Parasitic-Free Circuits

The considered metrics were evaluated on an ideal structure by neglecting all parasitics. The results are shown in Table XVIII.

Table XVIII. Evaluation of a Parasitic-Free 4-to-1 MUX in lithographic crossbar technology

	Area ( $\mu\text{m}^2$ )	Intrinsic delay (ps)	$K_{\text{Load}}$ ( $\text{ps}\cdot\text{fF}^{-1}$ )	Average power at 650 MHz ( $\mu\text{W}$ )
<i>MOS 65nm</i>	10.4	155.2	11.2	1.47
<i>Crossbar Compact dimensions</i>	1.74	121.8	97.5	0.77
<i>Crossbar vs. CMOS</i>	x 6	x 1.3	x 0.1	x 1.9

Such a crossbar structure with its high integration density is able to improve drastically the area by a factor of 6, while the intrinsic delay and the average power are improved respectively by 1.3 and 1.9. The delay analysis is depicted in figure 94, while the full power consumption figures are shown in Table XIX. While the area gain is due to the compactness of the crossbar structure, the improvement of the other metrics is mainly due to the good isolation of the nanowires in the FDSOI technology compared to bulk. In Table XVIII, we also see that load influence on delay is higher than in bulk. This is due to the confined crossbar dimensions, which constrain transistor drive strength in  $W=L=F$  (5.4.3.1).

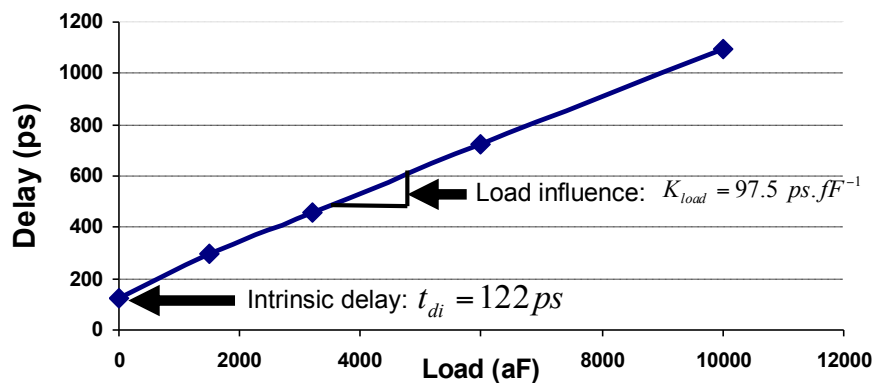


Figure 94. Delay analysis of lithographic crossbar in a parasitic-free context

Table XIX. Detailed Power Consumption figures for the sublithographic cell in a parasitic-free context (65-nm – 650MHz)

Power in $\mu\text{W}$	CMOS MUX	Crossbar lithographic	Gain
$\overline{A+B}$	1.91	0.76	2.51x
1	0.54	0.12	4.50x
$A+B$	1.86	0.77	2.41x
$\overline{A}$	1.72	0.77	2.23x
$A$	1.55	0.76	2.04x
$\overline{A}\cdot B$	1.94	0.78	2.49x
$A+\overline{B}$	1.85	0.76	2.43x
$\overline{B}$	2.13	0.76	2.80x
$B$	2.20	0.79	2.78x
0	0.52	0.12	4.33x
$A\cdot\overline{B}$	1.84	0.76	2.42x
$\overline{A+B}$	1.98	0.78	2.54x
$A\cdot B$	1.87	0.77	2.42x
$\overline{A}\cdot\overline{B}$	1.92	0.76	2.52x
$A\oplus B$	3.24	1.42	2.28x
$\overline{A\oplus B}$	3.26	1.42	2.29x
<i>Average</i>	1.89	0.77	2.46x

5.4.4.3 *Impact of the Width of Horizontal Wires Including Parasitics*

Figure 95 and figure 96 show the propagation delay in rising and falling edges for varying  $F_{HP}$  and  $F_{HN}$ , with an inverted output. When  $F_{HP}$  (resp.  $F_{HN}$ ) increases, the falling (resp. rising) propagation delay decreases and tends to a minimum value close to 480ps. When  $F_{HP}$  (resp.  $F_{HN}$ ) increases, the rising intrinsic propagation delay increases linearly by about 280ps (resp. 400ps) at each increment.  $F_{HP}$  and  $F_{HN}$  impact the crossbar structure in several ways.

When  $F_{HP}$  (resp.  $F_{HN}$ ) increases, the width of the  $p$ -type (resp.  $n$ -type) transistors increases, but so do also some parasitic devices. While the parasitic resistances (B) are reduced, the parasitic capacitances (A) at the cross points are increased. While increasing  $F_{HP}$  (resp.  $F_{HN}$ ) gives a better response at the output of the inverter stage with respect to the falling (resp. rising) time by improving the transistor drive strength, the rising (resp. falling) response, is however essentially driven by the parasitic capacitance contribution. To minimize this influence factor, we must keep the size as close as possible to the minimum  $F$ .

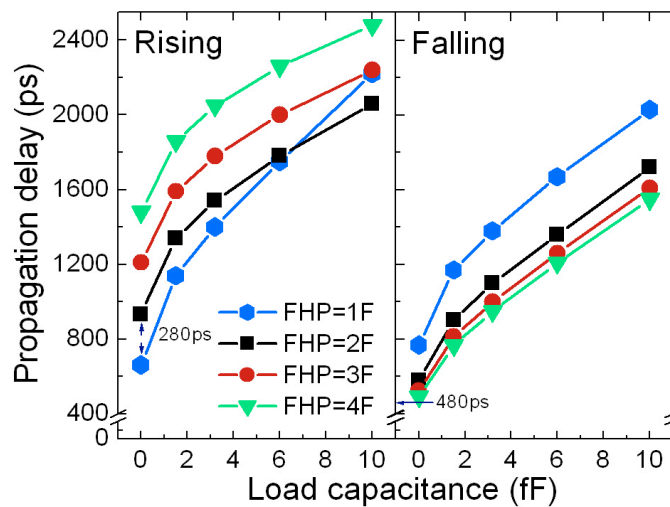


Figure 95. Influence of P-type horizontal wire size ( $F_{HP}$ ) on MUX propagation delay (Output of the inverter) ( $F_{HN} = F_V = F_I = F = 60 \text{ nm}$ )

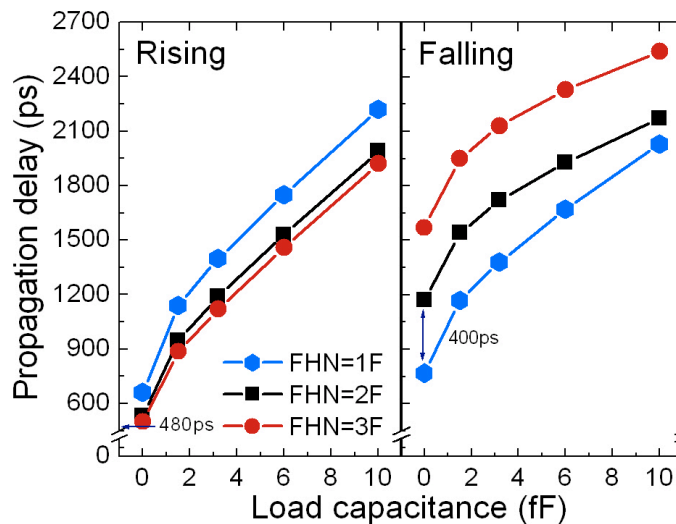


Figure 96. Influence of N-type horizontal wire size ( $F_{HN}$ ) on MUX propagation delay (Output of the inverter) ( $F_{HP} = F_V = F_I = F = 60 \text{ nm}$ )

5.4.4.4 *Impact of the Parasitic Capacitances*

Figure 97 shows the impact of inter-wire insulator size on the propagation delay. We notice that the propagation delay remains almost constant for any value of  $F_I$ . The spacing between the lines is directly related to the crosstalk parasitic capacitances between wires (C). Its influence appears negligible compared to other parasitic capacitances. Nevertheless, this result might not be adapted to aggressive scaling, and an in-depth study will be required.



Figure 98 shows the impact of gate oxide thickness ( $T_{OX}$ ) on the propagation delay. As depicted in this figure, the capacitances formed under the passive cross points appear to have the strongest parasitic contribution on the propagation delay. An optimal thickness can then be found around 6nm for a lithographic pitch of 65nm, yielding a potential reduction of the delay by a factor of 2x, if the oxide thickness is optimized. The oxide thickness ( $T_{OX}$ ) has two opposite impacts. While  $T_{OX}$  increases, the parasitic capacitances are reduced but the electrostatic effect on the FETs is degraded. In figure 10, we see that a small  $T_{OX}$  induces a large parasitic coupling, slowing the structure down. On the other hand, for very large  $T_{OX}$ , the performance is lower because of the slower FETs.

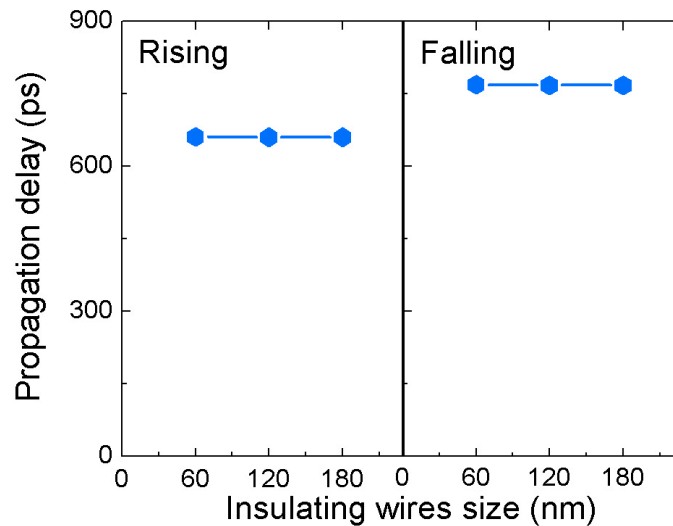


Figure 97. Influence of insulating wire size ( $F_1$ ) on MUX propagation delay ( $F_{HP} = F_{HN} = F_V = F = 60 \text{ nm}$ )

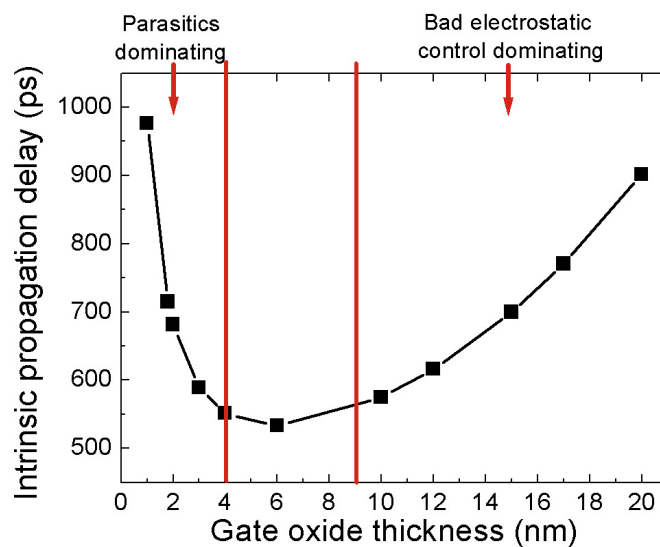


Figure 98. Influence of gate oxide thickness ( $T_{OX}$ ) on MUX propagation delay ( $F_{HP} = F_{HN} = F_V = F_I = F = 60 \text{ nm}$ )

#### 5.4.4.5 Impact of Scaling Including Parasitics

Figure 99 shows the impact of scaling on the crossbar implementation compared to CMOS. We obtained the data related to CMOS using the ITRS data [6]. Scaling plays a major role in reducing the parasitic capacitances, while the parasitic resistances are maintained constant. The propagation delay is then reduced accordingly, and thus the performance scales advantageously for crossbar implementation compared to conventional CMOS. For aggressively scaled nodes ( $< 22\text{nm}$ ), the crossbars maintain their advance in terms of area and finally surpass their CMOS counterpart in terms of performance, as shown in Table XX. Finally, we found that the average power consumption of crossbars is lower than that of the CMOS implementation by a factor of about 1.5 ( $1.47\mu\text{W}$  vs.  $0.99\mu\text{W}$  at 650 MHz).



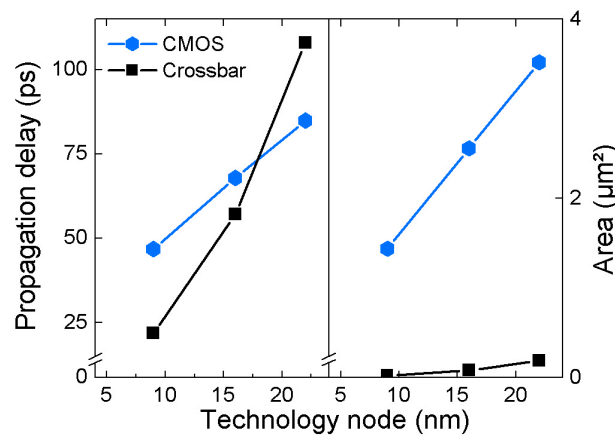


Figure 99. Impact of scaling ( $F_{HP} = F_{HN} = F_V = F_I = F$ )

Table XX. Evaluation of scaling with parasitics

	Area ( $\mu\text{m}^2$ )	Intrinsic delay (ps)
<i>MOS 65nm</i>	10.4	155.2
<i>NW Crossbar</i>	1.74	533
<i>Crossbar vs. CMOS</i>	x 6	x 0.29
<i>MOS 22nm</i>	3.52	84.9
<i>NW Crossbar</i>	0.19	107.9
<i>Crossbar vs. CMOS</i>	x 18	x 0.78
<i>MOS 9nm</i>	1.44	46.8
<i>NW Crossbar</i>	0.02	21.86
<i>Crossbar vs. CMOS</i>	x 67	x 2.14

#### 5.4.4.6 Performance Estimations including Parasitics

To summarize, we evaluated the different metrics (area, delay, power), while taking into account the presence of parasitic devices. The evaluated structures have an optimized oxide thickness  $T_{OX}$  of 6 nm. The results are shown in Table XXI, while the delay analysis results are in figure 100 and the complete power figures are in Table XXII. Even if the technology was optimized to minimize the parasites, we observe a substantial degradation in performance in terms of intrinsic delay and load factor. Nevertheless, the average power consumption of crossbars is lower than that of the CMOS implementation by a factor of about 1.5.

Table XXI. Evaluation of lithographic crossbar-based structure with parasitics

	Area ( $\mu\text{m}^2$ )	Intrinsic delay (ps)	$K_{Load}$ ( $\text{ps}\cdot\text{fF}^{-1}$ )	Average power at 650 MHz ( $\mu\text{W}$ )
<i>MOS 65nm</i>	10.4	155.2	11.2	1.47
<i>Crossbar Compact dimensions</i>	1.74	533	208	0.99
<i>Crossbar vs. CMOS</i>	x 6	x 0.29	x 0.05	x 1.48

#### 5.4.4.7 Discussion and Guidelines

The previously surveyed simulation results show different design and fabrication scenarios for a crossbar MUX, which can be considered in order to assess the impact of technology on the performance of crossbars in general. The simulations confirm the compactness of the crossbar architecture and its ability to yield up to 6x area saving for the considered MUX design. Moreover, the simulations performed on ideal structures with no capacitive coupling between the different lines also confirm the results reported in the literature [174] on the high performance of the crossbar architecture with respect to CMOS.

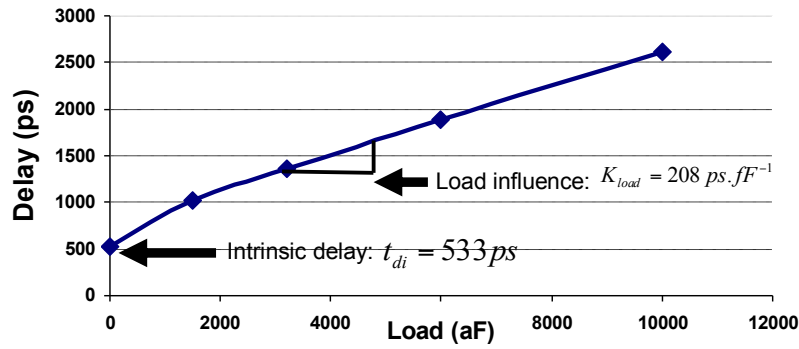


Figure 100. Delay analysis of lithographic crossbar with parasites

Table XXII. Detailed Power Consumption figures for sublithographic cell (65-nm – 650MHz)

Power in $\mu\text{W}$	CMOS MUX	Crossbar lithographic	Gain
$\overline{A+B}$	1.91	0.99	1.93x
1	0.54	0.17	3.17x
$A+B$	1.86	0.99	1.88x
$\overline{A}$	1.72	0.98	1.75x
$A$	1.55	0.98	1.58x
$\overline{A.B}$	1.94	0.99	1.96x
$A+\overline{B}$	1.85	0.99	1.87x
$\overline{B}$	2.13	1.03	2.07x
$B$	2.20	1.02	2.16x
0	0.52	0.17	3.06x
$\overline{A.B}$	1.84	0.98	1.88x
$\overline{A+B}$	1.98	0.99	2.00x
$A.B$	1.87	0.98	1.91x
$\overline{A.B}$	1.92	0.99	1.94x
$A\oplus B$	3.24	1.78	1.82x
$\overline{A\oplus B}$	3.26	1.78	1.83x
<i>Average</i>	1.89	0.99	1.92x

However, realistic estimations of the capacitive coupling between the crossbar lines show that their impact cannot be neglected, especially when they are combined with the resistance of the doped semiconducting lines. In this case, performance drops by 4.5x. This result can be understood by the fact that the lines of a crossbar perform the interconnect task. The insulator is intended to mutually shield the parallel lines (lying on the same level) and the perpendicular lines (lying on top of each other). These two opposing effects lead to a high parasitic contribution.

The largest impact of parasitic capacitances stems from those lying at the cross points, *i.e.*, those formed by the crossing lines. Their impact is dominant with respect to that of the lines lying between the parallel lines in every layer. In order to address the issue of capacitive coupling between the crossing lines, a solution can be driven from the technology side. While the insulating material between the top and bottom lines has to be optimized by using a thick or a low-k material, this insulator represents at the same time the gate oxide, which should be thin or a high-k material in order to improve the gate control over the channel. These conditions are not compatible, which is a natural result of the double nature of the crossbar: it performs both interconnect and logic with the same resources. Consequently, there is a trade-off between the speed of logic and the speed of interconnect. It is therefore necessary to accept a certain compromise in terms of delay when the insulating material between the crossing lines is chosen. Our study shows that for a  $\text{SiO}_2$  gate oxide, the trade-off is optimized for an oxide thickness of about 6nm.

The investigation we carried out in this paper is based on a MUX aiming to illustrate the general trend of crossbars. Because of the double role of this architecture (interleaved interconnect and logic), a trade-off in performance must be found. The technology must optimize the insulating oxide between the crossing layers in order to meet the optimal trade-off. It is worth noticing that the results of these investigations are valid under our technological assumptions of a lithographic and manufacturable process.

## 5.5 Global Comparisons

In this chapter, we investigated the use of three different technologies to build a compact logic computation cell. Among these three technologies, we can differentiate the first technology, which aim to increase the functionality of the elementary device such as the double-gate carbon nanotube field effect transistor, from the other technologies, which help to increase the density of elementary transistors. In the latter category, we examined a highly-speculative approach, based on the sublithographic arrangement of nanowires and a more realistic one, which is derived from an industrial Fully Depleted Silicon on Insulator process flow.

Table XXIII summarizes the metrics used to evaluate the different proposals. To compare the different technologies objectively, we will split the evaluation of future middle-term solutions (carbon electronics and sublithographic arrangements of nanowires) from that of short-term, potentially immediately manufacturable solutions (lithographic crossbar arrangement). The future solutions are compared to a scaled 22-nm node, while the short-term proposal, based on a standard industrial process, is compared with respect to a 65-nm CMOS node.

**Table XXIII. Global comparison of analyzed disruptive technology cells**

	Middle-term proposals			Short-term proposals	
	CMOS MUX 22-nm	DG-CNTFET cell	Sublithographic crossbar	CMOS MUX 65-nm	Lithographic crossbar
<i>Area (<math>\mu\text{m}^2</math>)</i>	1.19	0.39	0.29	10.4	1.74
<i>Intrinsic delay (ps)</i>	79	149	17	155.2	533
<i><math>K_{Load}</math> (ps.fF<sup>-1</sup>)</i>	5.7	124.5	15.1	11.2	208
<i>Average power (<math>\mu\text{W}</math>)</i>	4	1.78	6.7	1.47	0.99

It is worth noticing that the technologies improving the density are highly beneficial for area improvement.

Indeed, the best gain is obtained for a lithographic crossbar with a grain of 6x, compared to current CMOS technologies. Regarding the middle-term solutions, the sublithographic crossbar is also the most compact solution.

The best gain in power consumption was reached for carbon electronics, with a gain of about 2x compared to its CMOS counterpart. It appears clear that the good electron transport in carbon structures is at the origin of this.

Regarding the delay performances, the CMOS circuits are still good candidates for future generations of circuits. Nevertheless, performance optimization for the other technologies has not been carried out for this first order evaluation, and performance can certainly be improved after correct sizing and process tuning.

## 5.6 Conclusion

In this chapter, we have studied how emerging technologies could be used to reduce the size of the computation node for FPGAs. In fact, it is worth noticing that only 14% of the FPGAs

area is occupied by the logic blocks. More than 80% is used just by peripheral circuitries. Thus, the FPGA architectural scheme seems quite inefficient if we compare the circuit area used for computation to that of the others. It is thus of high interest to propose new elementary logic circuits, that may lead to alternatives to the standard FPGA scheme.

Emerging technologies have been used to reduce the size of the computation node compared to CMOS. We have surveyed two main ways to do so: *functionality*-improvement based devices and *density*-improvement based devices.

For functionality-based improvements, we surveyed the use of the DG-CNFET technology. Such a device can be configured to *n*-, *p*- or *off*-type depending on the voltage applied to the back-gate terminal of the device. We first presented an integration process for the technology and performed device optimization, in order to ensure correct hypotheses. In-field reconfigurable devices lead to compact reconfigurable logic cells. We performed the evaluation of a previously designed circuit and compared it to scaled CMOS, and showed that the reduction of the number of transistors yields a reduction of 3.1x in area. Power consumption is reduced by 2x, due to the electronic properties of carbon. Finally, the use of back-gate control is generalized to standard dynamic logic cells. We especially used the off-state to merge the clocking transistors in the function paths. Thus, we demonstrate a gain in terms of transistors of up to 50% for the XOR gate.

For the density-based improvement techniques, we surveyed the use of dense crossbar implementations of devices. The crossbar implementation expects to merge connections and active logic in a very small area. Various assumptions for crossbar implementations could be used. Firstly, we envisage the use of sublithographic techniques to realize dense arrays of nanowires that formed FET at their cross points. Sublithographic techniques are based on very speculative bottom-up arrangements of wires at the nanoscale. After describing all technological assumptions, we derived a NASIC approach to build a multiplexer intended to work and to be connected at the lithographic scale. We showed that this technology yields an improvement of 4.1x in area and 4.6x in delay, thanks to the drastic increase of density and the consequent reduction of parasitics.

Sublithographic integration processes are highly speculative and only long-term solutions. The crossbar concept is of course not limited to such an advanced technology and could be extended to lithographic processes. In this way, we derived a FDSOI process flow in order to realize a dense crossbar arrangement of transistors, with lithographic dimensions. We then proposed a layout methodology for standard circuits. After describing all the parasitic elements in the structure, we performed a circuit-driven optimization of the technology to find the best trade-off between the parasitics and the transistor performances. We showed that this solution yields an improvement of 6x in area and 1.48x in power consumption. This is due to the crossbar arrangement and the performance of FDSOI technology in terms of power. We also showed that this approach can be expected to outperform the standard CMOS scheme from the 16nm technological node onward.

To conclude, various technologies have been used with the same objective: build a compact and improved logic cell. In the light of these results, we can clearly see that a trade-off must be found depending on the required specifications in terms of delay, area or power. Nevertheless, we should notice that all the cells reduce the area figures. This may lead to the use of a new highly compact elementary logic cell for high-performance computation. This basic block will be envisaged as a new seed for architectural organization, discussed in the next chapter.



# **CHAPTER 6**    *Disruptive Architectural Proposals and Performance Analysis*

---

## **Abstract**

In this chapter, we explore disruptive architecture proposals. In the previous chapter, we showed that it is possible to obtain very compact reconfigurable in-field computation cells. These cells require architectural adaptations. For example, the use of a fine-grain logic cell is not possible in the context of a conventional reconfigurable circuit scheme, since it would incur a large interconnection overhead.

We thus proposed an architecture for this compact logic, characterized by the association of a logic layer, to adapt the granularity and the use of fixed interconnection topologies to reduce the routing impact. The FPGA organizational scheme is then adapted to this new layer.

To compare, in an objective way, our approach with conventional FPGAs, it was necessary to develop a specific toolflow suited to our requirements, able to describe the designed architecture. Based on the VTR toolflow, the tool integrates fixed topology routing and the specific organization of the layered architecture.

Benchmarking simulations were performed. In a first approach, a local exploration of the proposed layer is done, in order to study the impact of the fixed interconnect topologies. We showed that the Modified Omega topology gives the best mapping rates on the structure with about 90% of mapping success for 6-node graphs. In a second approach, complete architectural benchmarking was conducted and we showed that the proposed architecture leads to an improvement, in area saving, of 46% in average, with respect to CMOS. We also discovered that the routing delay is less distributed and tends to be more controllable than in the traditional approach.

The proposed architecture shows good promise for future FPGA systems. It is also worth noticing that the toolflow that we developed opens the way towards more exploration of disruptive architectures. The potential of ultra-fine grain logic is thus very broad and particularly promising for emerging reconfigurable circuits.

In the previous chapter, we introduced compact logic cells. With emerging technologies, an elementary logic function can be efficiently implemented by hardware. The efficiency is based on several different aspects of the technology. The augmented functionality of the device was used to reduce the number of transistors required per function. Higher device integration density enables many devices to be embedded in a small area, which obviously compacts the size of the implemented logic. However, this aspect does not at first appear to be a particular advantage for reconfigurable logic. In section 2.1.3, we stated that only a low proportion of the overall circuit is used by elementary combinational and sequential logic (14% of total area), as compared to routing and memories (86% of total area).

Nevertheless, with the emergence of a novel ultra-compact logic cell, we face new architectural requirements, which could lead to a change in balance between logic resources and routing/memory resources. In this chapter, we will describe a new architecture for reconfigurable ultra-fine grain computing. To address the specificity of this new architecture, we will adapt the benchmarking methodology and present a specific packing tool. Finally, we will size the new architectural model and compare it to a standard MOS counterpart.

## 6.1 Introduction

In chapter 2, we already presented the *Field Programmable Gate Array* (FPGAs) architecture. This is currently the most common architectural scheme for reconfigurable logic. The architecture is designed for current CMOS technologies, and its regularity is compatible with most scaled technological nodes. The logic is performed through the use of *Configurable Logic Blocks* (CLBs), able to perform coarse-grain combinational and sequential logic functions. These CLBs are arranged into islands and interconnected through a complete (but resource-limited) programmable interconnect. The CLBs are formed by several *Basic Logic Elements* (BLEs) interconnected by a fully programmable interconnect infrastructure. BLEs are composed of *Look-Up Table* (LUT) and *Flip-Flop* (FF) elements. In modern FPGAs, several other blocks can be found, in order to perform specific optimized tasks, such as hardware multipliers or memory banks. This heterogeneous nature of the structures helps to optimize the performance for specific application classes. Nevertheless, since we focus on generic reconfigurable logic structures, we will consider in this chapter only homogeneous structures. Obviously, in conventional CMOS FPGAs, all these building blocks are sized to ensure a correct granularity of the structure and to maximize the performance of the assembled system.

In this work, we looked at disruptive technologies and saw that they could lead to very compact logic elements. These elements will be proposed as the foundation stone for all reconfigurable arrangements. Nevertheless, the organization of such blocks induces new problems to solve. As an illustration, if we consider the in-field reconfigurable cell proposed in 5.2.4, the compactness of the cell means that it is not possible to transpose directly the traditional FPGA architecture to our ultra-fine grain cell, since complete interconnectivity would be required between the elements. Again, we may consider that about 43% of silicon area is used for interconnects, while only 14% is used for logic [38]. However, in the case of ultra fine-grain logic, this ratio would further worsen, since the routing structures do not scale in the same way as logic. Thus, there will be a significant area imbalance between logic and interconnect, consequently leading to a severe loss of efficiency in the structure. Another area imbalance question might be raised by memories. Indeed, even if the size of the logic element has been decreased, the amount of memory required remains the same. In chapter 3, proposals have been made to reduce the complexity of memory in an FPGA. Nevertheless, their area impact is not zero. It is thus fundamental to bear in mind the extra circuitry required by these configuration memories during performance evaluation. Finally, another constraint is present - that of reliability. The use of billions of unreliable emerging devices can be handled by the use of specific fault-tolerance mechanisms. Even if these questions are not addressed in this thesis, it is important to keep this aspect in mind.

To obtain a credible disruptive block and architecture based on emerging technologies, all these considerations must be addressed. In this chapter, we propose an architectural scheme that addresses these issues. We will specifically address the questions regarding the use of ultra-fine grain logic cells in a reconfigurable circuit. It is worth noticing that the resulting architecture will be in aligned with the architectural template of the study. Nevertheless, due to the specificity of the ultra-fine grain approach, the template will be enhanced. Then, in order to evaluate in an objective way the potential of the structures and technologies with respect to their CMOS counterparts, we will also propose specific tools. This tool flow will be compatible with standard approaches, and will only be enhanced to address the specific architecture. Finally, we will evaluate the proposal and compare it to a standard equivalent circuit.

## 6.2 Architectural Proposals

### 6.2.1 Introduction

As already mentioned, FPGAs possess a layered structure, as depicted in figure 101-b. Four layers are traditionally considered: elementary logic with *Look-Up Tables* (LUT layer); fine grain combinational and sequential logic with *Basic Logic Elements* (BLE layer); coarse grain logic with *Configurable Logic Blocks* (CLB layer); and finally an island style arrangement (FPGA layer). This organization follows the conventional hierarchical template, shown in figure 101-a. All blocks were detailed in chapter 2.

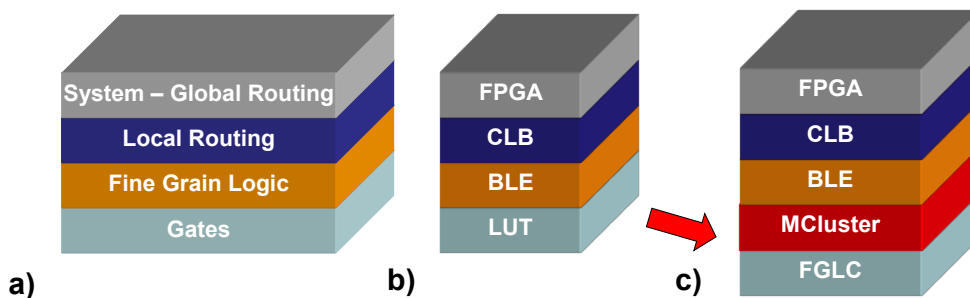


Figure 101. Conventional hierarchical template (a), FPGA model (b) and modified levels (c)

In our approach, we replace the traditional LUTs by a *Fine Grain Logic Cell* (FGLC). However, if we were to project this onto the conventional approach, each FGLC would be directly connected to a full connectivity unit, such as a set of switchboxes. Considering the in-field reconfigurable gate as an example, only seven transistors are used in the cell, while a similar number of transistors (at least six) are used for a 1-bit switchbox. This projection would therefore result in a large overhead in terms of connections and, as mentioned previously, would worsen the already significant imbalance between routing and logic resources at the FPGA level. The overall structure would consequently be left with an intrinsic granularity imbalance.

To solve this issue, we propose to modify the FPGA scheme, with the layered organization shown in figure 101-c. To correct the granularity issues, we propose to pack the cells into an intermediate structure which is compatible in terms of size with the traditional levels of FPGAs: *MClusters* (for Matrix Clusters). Thus, we obtain fine-grain logic blocks with enhanced functionality with respect to conventional LUTs. Within this organization, the connectivity issue between logic cells must be addressed, considering that full connectivity is not an option. We thus propose a solution based on fixed interconnect topologies. Finally, the convergence with the conventional FPGA scheme will occur at the BLE level. In particular, the MClusters will be used as the replacement block for LUTs.

### 6.2.2 MCluster Organization

#### 6.2.2.1 MCluster: Adding a "Logic Layer"

At the MCluster layer, elementary ultra-fine grain logic cells are used, where the term "ultra-fine grain" is twofold. First, it covers the size of the cell, which is highly compact, as opposed



to its CMOS equivalent. Second, it also covers the granularity of the logic functions covered. In an FPGA, the smallest block is in general the LUT, and it has been shown that an LUT with 4-inputs is the best solution for routability efficiency [175]. With a 4-input LUT, it is possible to realize 65536 ( $2^{2^4}$ ) functions. The targeted logic block (the in-field reconfigurable logic cell) uses a more restricted set of functions. This means that we also introduce an imbalance between fine-grain logic and ultra-fine grain logic.

To increase the logic coverage of the structure, we propose a 2-D matrix assembly of ultra-fine grain logic. The logic cells are arranged in a layered structure, with connections only existing between adjacent layers in order to avoid long connections and to maximize local connectivity. This organization, called *MCluster*, is depicted in figure 102 and will handle fine-grain logic operations in the proposed organization. In the following, MClusters will be defined by the w (width) and d (depth) parameters and written as: MCluster\_d\_w.

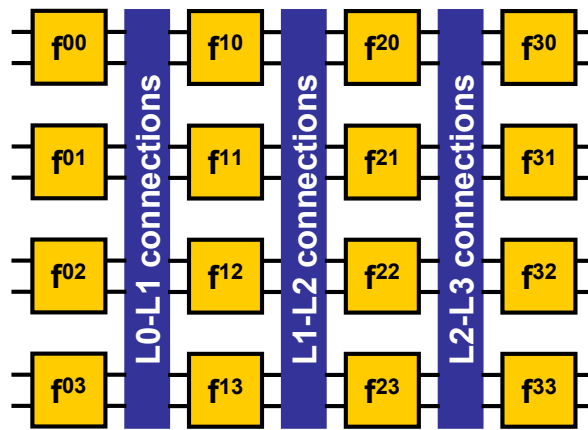


Figure 102. MCluster approach for reconfigurable architectures (MCluster\_4\_4)

#### 6.2.2.2 *MCluster: Simplifying the Interconnect Overhead*

For intra-matrix interconnect, any total interconnectivity topology is prohibited, because of the associated wiring complexity and the extra area requirements. Instead, and through analogy to computer networks, our approach is to adapt incomplete interconnection sets to the matrix architecture. In fact, *Multistage Interconnection Networks* (MINs) are designed to interconnect computing layers in an efficient way and can be applied in this context. In computer science, MINs are used to interconnect layers of switchboxes in order to route information packets only. This concept has been reported several times in interconnect strategies for VLSI, for example in [176, 177], and used in FPGAs to reduce the complexity of wiring through logic blocks [178, 179], by proposing an architectural organization for the routing circuits. It is worth noticing that the main difference here, with respect to the network context, is that switchboxes have been replaced by logic cells, thus introducing computing directly inside the network. Furthermore, the use of very local MIN-style interconnect has a drastic impact on the size and the wire length is reduced accordingly.

There are many MIN topologies and combinations, but in this work, we focus on four typical permutations [180, 181]: Banyan (Figure 103-a), Baseline (Figure 103-b), Flip (Figure 103-c) and Modified Omega (Figure 103-d), where the modifications to standard Omega maximize the shuffling in this topology. Since the interconnect topology is fixed and static, the choice of topology must be made by the designer during architectural development. We will discuss the intrinsic performance metrics of each topology in the following sections of this chapter.

While fixed interconnect topologies are useful for area considerations, we must also consider the cost in terms of logic mapping performance. In this context, we could explore trade-offs between fixed interconnection patterns and full connectivity. In particular, since we have seen that vertical nanowire field-effect transistor technologies are useful to build smart vias, we could envisage a topology, which can be reconfigured between two or more sets of

interconnect, with the number of interconnect sets quite low to avoid a large number of configuration signals.

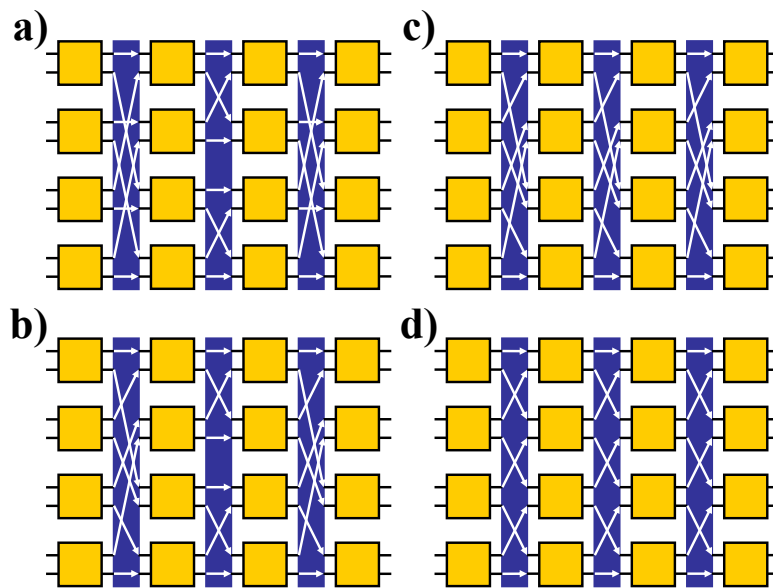


Figure 103. Matrix of 16 reconfigurable gates with fixed interconnect topology (a)Banyan, b)Baseline, c)Flip, d)Modified Omega)

### 6.2.3 BLE and CLB organization

We described above a matrix arrangement for ultra-fine grain logic cells. This arrangement is able to perform combinational logic functions and can thus be considered to be a replacement element at the LUT level. We now consider how the BLE/CLB organization should be adapted using the MClusters.

Figure 104 shows the organization of the logic blocks at the CLB level. The BLE are composed of  $MClusters_{d \times w}$  used to perform the combinational logic. In a traditional approach, all outputs of the LUTs may be latched and multiplexed to the CLB outputs. In our approach, a partial depopulation of the output latches is implemented as shown in figure 104, where one of the two outputs are not latched. The subsequent CLB organization remains similar to the LUT approach, where all outputs of the BLEs are connected to the global reconfigurable interconnect, and are also fed back to the inputs of the CLB, thus achieving a full connectivity pattern. The  $d \times N$  feedback signals and the  $I$  inputs of the CLB may be routed to any of the MClusters by using the reconfigurable multiplexers. Finally, a single clock signal is used to control the sequential elements.

### 6.2.4 FPGA Final Organization

At the top level of organization, the granularity (as well as the configuration and routing) assessments are the same as for conventional FPGAs. Thus, we will consider the FPGA layer, as depicted in figure 105. Hence, the logic blocks are interconnected by a complete interconnect set, but limited in terms of resources (i.e. all the connections can not be realized simultaneously). As in conventional FPGAs, this architecture uses switchboxes and connection boxes to ensure the connectivity between the CLBs. Nevertheless, we can see from figure 105 that this "sea of logic" is not directly connected to I/O blocks. While the standard FPGA approach could be used, the reconfigurable circuit is connected in this approach to a logic manager, which can then communicate data to global interconnect through a Network-On-Chip interface. In this way, the architecture can be integrated for example as a reconfigurable hardware accelerator in a Multi-Processor System-on-Chip or in a heterogeneous circuit. Finally, specific blocks are also included, such as a fault-tolerant logic router and a performance improvement decision controller. This kind of IP can operate as a service block to help in re-defining the structural placement and routing, allowing defective CLBs to be avoided, or power consumption to be optimized in real-time throughout the circuit.

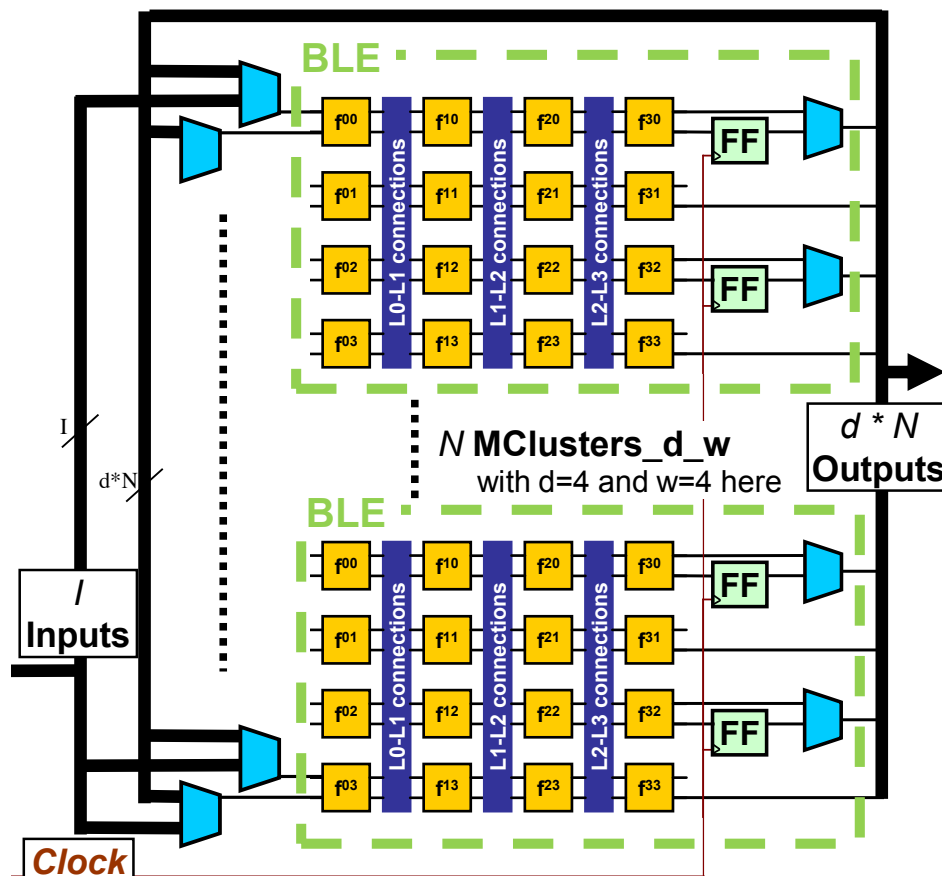


Figure 104. MCluster-based CLB proposal

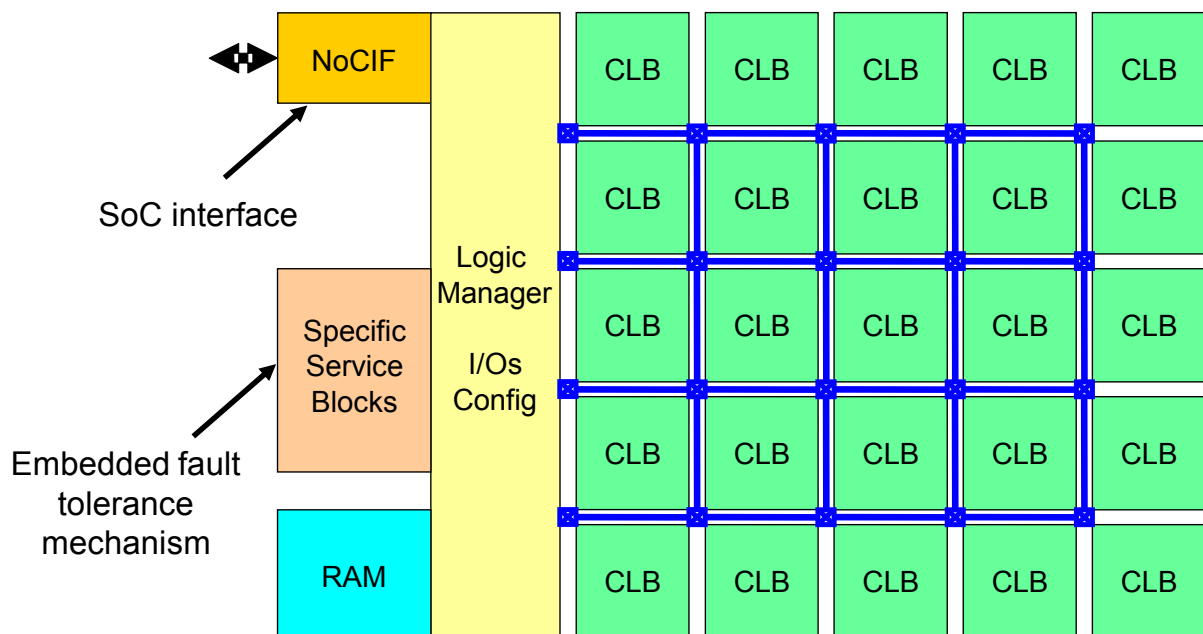


Figure 105. Final FPGA layer organization

### 6.3 Benchmarking Tool

In the previous section, a specific organization for ultra-fine grain logic cells was described. To analyze performance gain from an application point of view and compare it objectively to existing approaches, we must evaluate the performance of its implementation of well-known circuits (benchmarks). Thus, we propose a complete benchmarking tool flow suited to our proposal, as well as to existing architectures.

### 6.3.1 General Overview of the Flow

In 4.2, the traditional FPGA benchmarking flow was described. Based on the *VPR 5.0* tool suite, designed to handle LUT-based logic, the architectural description is basically limited to a homogeneous description of logic blocks. The toolflow must therefore be adapted in order to handle the complexity of our structures. This toolflow is now based on the *Verilog-To-Routing* (VTR) [182] project from Toronto University. This toolflow is able to describe the complexity present in modern commercial architectures. The VTR 6.0 flow uses a logic block description language that can express far more complex logic blocks than is currently possible with any other publicly available toolsets. It can describe complex logic blocks with arbitrary internal routing structures, can allow arbitrary levels of hierarchy within the logic block and can assign different functional modes to blocks. The latter property is necessary since, modern commercial FPGAs have different modes in their LUT usage or memory blocks. The architecture is described in XML, and the toolflow is supported by *AA-Pack* [187] for the packing of logic blocks and *VPR 6.0* [182, 185] for the place and route part. Furthermore, the set of benchmarks that can be used with the toolflow are not restricted to well-known (but old) BLIF-format circuits, which are somewhat outdated and no longer completely representative of current and future applications. Thus, a logic synthesizer (*ODIN* [184]) has been connected to the flow, in order to input directly high-level VERILOG benchmarks of large applications.

This tool suite has been developed for FPGA research. Hence, the use of elementary logic blocks other than LUTs is not supported natively. While it is not in our interest to build a new toolflow from scratch, we will complement it to meet our requirements. This allows us to ensure good compatibility with the FPGA architecture and to guarantee its efficiency. The use of fixed interconnect topologies is not supported by VTR tools. We must thus add a specific tool to the flow: *M-Pack*. This tool handles the packing of logic cells into an *M-Cluster* set. As in the previous flow, the *ABC* tool will be used to perform the technology mapping [183]. While it is always possible to map the circuits onto a LUT target, specific libraries will be used for this mapping. This is of course required while considering the use of reconfigurable logic gates with a lower set of achievable functions. It also allows benchmarks to be mapped with specific logic gates. For example, the tools developed by EPFL [188] can be connected to the flow in order to build a library of logic gates based on a very specific set of logic functions realized with ambipolar transistors.

The obtained flow is very versatile, since it is possible to evaluate several architectural scenarios. A set of addressable scenarios is depicted in figure 107. Figure 107-a shows the traditional LUT-based FPGA scenario. Obviously, the original flow is still compatible with FPGAs, since it is based on the FPGA-dedicated VTR flow. Figure 107-b corresponds to the architectural scheme which has been proposed in this work. The flow addresses the matrix clustering and packing by the specific *M-Pack* tool. Then, since the architectural organization is close to that of FPGAs, it is possible to pursue the packing and routing with the VTR toolflow. Figure 107-c illustrates that, after matrix packing, it is possible to explore new kinds of organization for the high level logic layers of the structure. In particular, some opportunities for these levels will be discussed at the end of this chapter. Figure 107-d shows that the toolflow is compatible with very regular approaches. These architectures have been proposed for computation architectures in environments with a high level of defects. Thus, it is possible to envisage the description of the *NanoFabric* organization [47] or the *Cell Matrix* approach [189], where large matrices of cells interconnected with their neighbourhood are used to create the logic plane. Finally, figure 107-e depicts the ability of the flow to further consider reconfigurable architectures with specific custom logic instead of programmable elementary components. Hence, the flow is compatible with approaches which propose a logic cell as the elementary combinational logic block [190, 191].

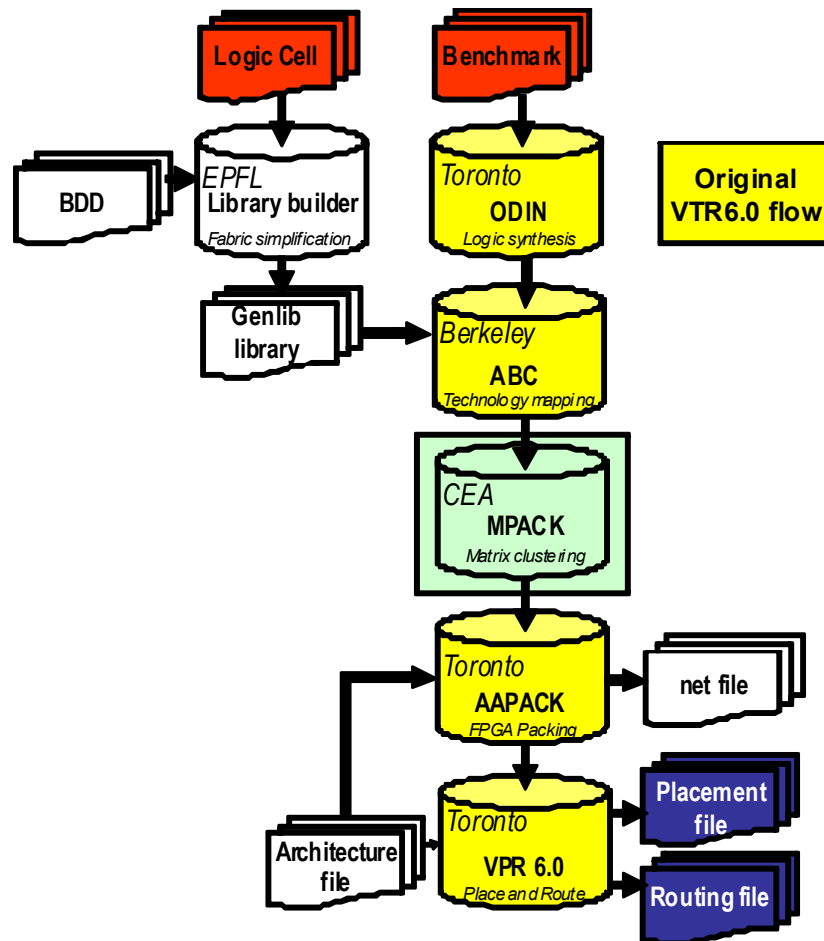


Figure 106. Disruptive technology compatible benchmarking flow diagram

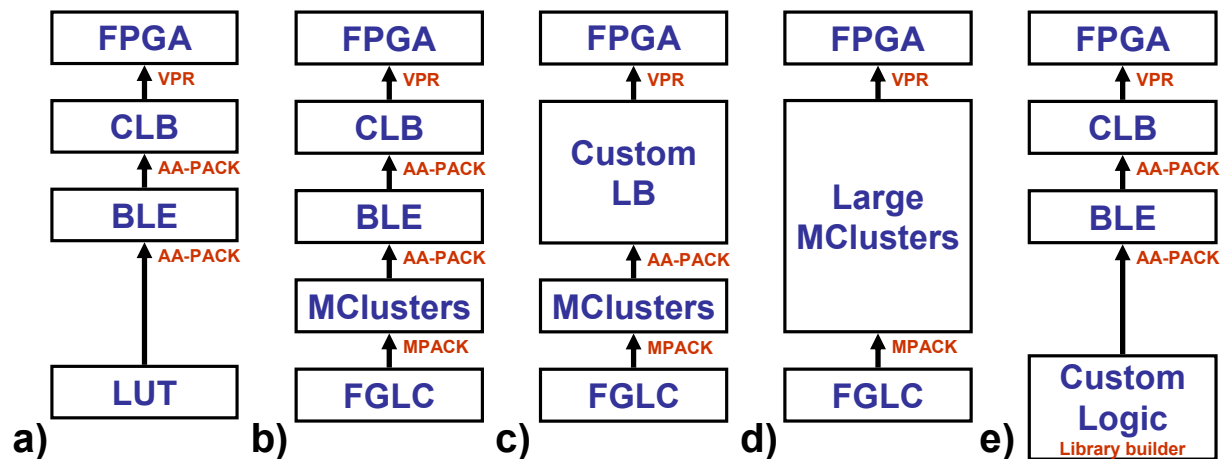


Figure 107. Illustration of possible scenarios compatible with the proposed toolflow

However, at the time of writing, it should be noted that *VPR 6.0* is still under development regarding timing analysis. More particularly, the tool cannot model inter-CLB timings. This makes it impossible to carry out exact comparisons to other architectures regarding delay. This feature will be enabled later.

### 6.3.2 MPACK: Matrix PACKer

The use of the MCluster arrangement of cells introduced several novelties to the architecture. In particular, the layered structure of the interconnect topology makes the use of traditional packing tools impossible. The proposed tool is able to manage specific interconnect topologies, as well as the associated buffering generation required by the layered structure.

6.3.2.1 Concept

The Matrix Packer has been designed to pack netlists of logic cells into a set of MClusters. Figure 108 describes the internal organization of the packing tool MPack. The tool is composed of two principal algorithms: the mapper and the clusterer.

The mapping algorithm aims to fit a netlist of logic cells onto an MCluster architecture. This module is split into two parts: the architectural optimization and the brute force mapper.

The clustering algorithm aims to group the logic cells into arrangements of MClusters. These arrangements are subsequently mapped onto the architecture in order to validate the legitimacy of the structure.

To ensure the integration into the whole flow, netlists are read and exported using the BLIF file format [192]. The Architecture Generator block generates the MCluster template. Finally, in order to help with tool debugging, it is possible to export every manipulated graph graphically, using the DOT language [193].

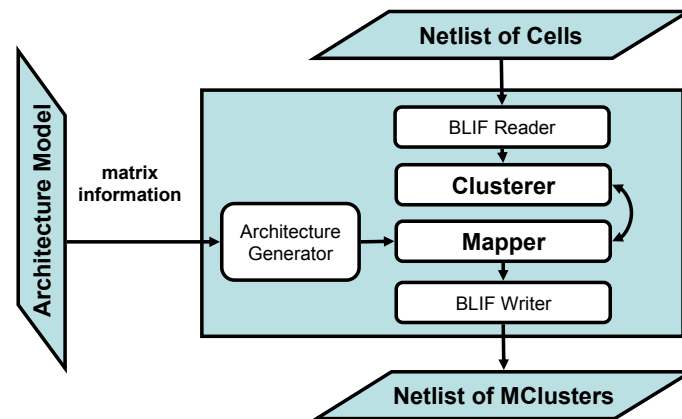


Figure 108. MPack model flow

The targeted MCluster shape is parameterized in terms of the size and topology scheme. An example of a targeted Banyan topology with a  $d=w=4$  is shown in figure 109-a. In the tool internal representation, the connectivity between nodes is represented by the adjacency matrix of the graph. In such matrices,  $(i, j)$  refers to the intersection of row  $i$  and column  $j$ . A 1 at the position  $(i, j)$  means that the point  $i$  is connected to the point  $j$ . As an example, the individual cross-connectivity matrices  $X_{nm}$  ( $X_{01}$ ,  $X_{12}$  and  $X_{23}$ ) between logic cell stages of depth  $n$  to  $m$  are shown in figure 109-b.

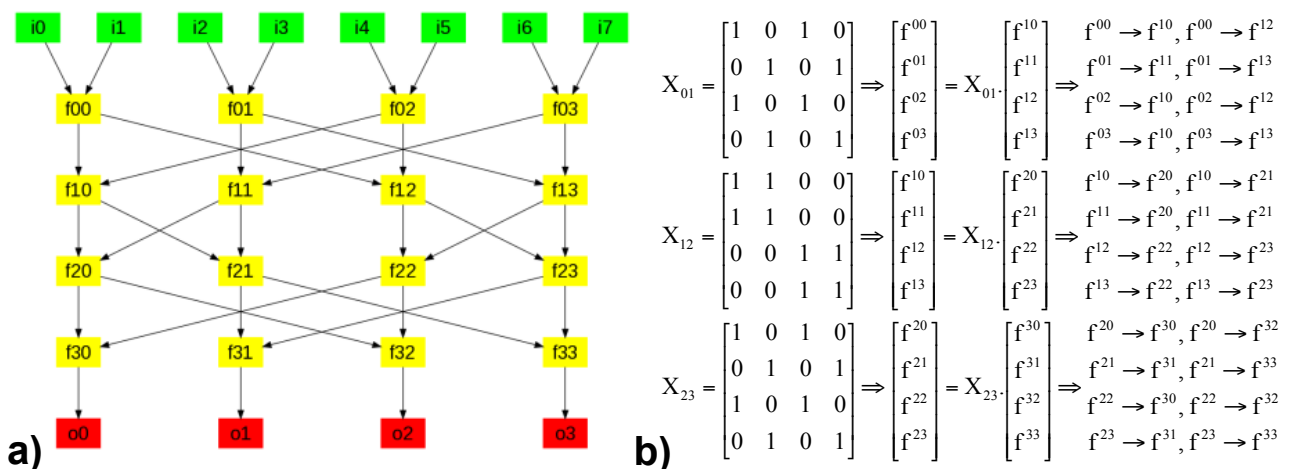


Figure 109. a) Banyan\_4\_4 topological arrangement (logic cells are in yellow, virtual input nodes are in green, virtual output nodes are in red) and b) associated cross-connectivity matrices

6.3.2.2 Mapping Algorithm : Architectural Optimization

In the mapping process, an architectural optimization is first performed in order to adapt the netlist to the physical architectural scheme. In fact, due to the layered structure, the logic can

be seen as pipelined. The input netlist of cell graphs has to be processed by adding necessary synchronization elements to extend the input and output data paths, as well as by adapting connections to conform to the physical topology. For instance, in an ABC output netlist, any logic gate output can be connected to several gate inputs. However, MClusters have a fixed fan-in and fan-out. The netlist is made MCluster-compatible by adding buffering nodes to limit the fan-out. The netlist is thus processed as follows:

- **Input Layer Correction:**

A check operation is performed to ensure that external data inputs only to the first layer. If an input node (i.e. a node with no input connections) is not positioned on the first layer, the hierarchical node position is updated. An illustration of this step is shown in figure 110.

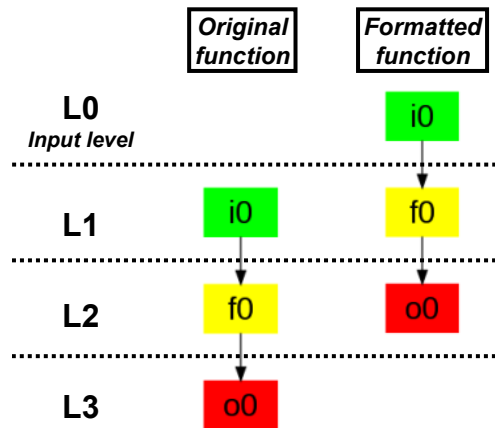


Figure 110. Input level correction illustration

- **Feedback Correction:**

In a real-life netlist, feedback connections are possible within a set of logic gates. Obviously, the layered organization of the structure does not allow a retroactive connection within the matrix. Any loop on a computation cell must thus be handled by the BLE external interconnect. Hence, it is necessary to modify any feedback connection into additional output and input terminals to ensure external connectivity. This case is illustrated in figure 111.

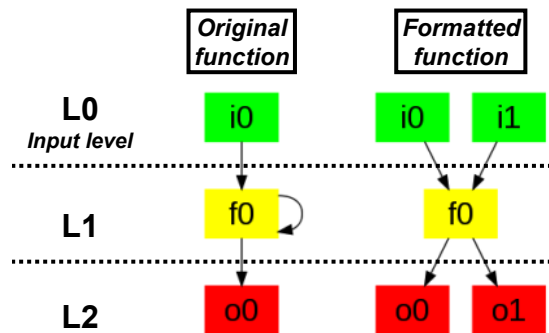


Figure 111. Feedback correction illustration

- **Jump Correction:**

Due to the layered structure, a connection which "jumps" at least one logic layer is not allowed by the topology. In fact, a connection between two layers must pass through a logic cell. Consequently, the creation of a path (by the addition of a buffer cell) instead of a jump is necessary to handle such a situation, as shown in figure 112.

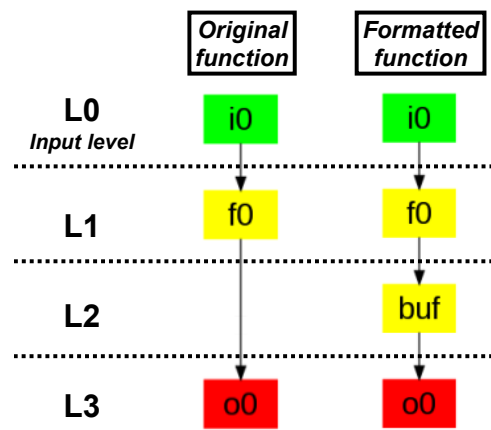


Figure 112. Jump correction illustration

• **Multiple Buffering/Inverting Path Correction:**

A logic cell which is connected to only one input signal is a buffer or an inverter. Hence, we can try to find multiple buffering/inverting paths in the logic function. Figure 113-a depicts the situation of an incoming path. Two values  $f_0$  and  $f_1$  are propagated by two buffer paths (formed by  $f_2$  and  $f_3$  respectively). The signals are then used by the  $f_4$  cell to compute the  $o_0$  result. It is clear that buffering both inputs to propagate them through the structure is not optimal. It is of course preferable to perform the operation as soon as possible in the path and to buffer only the result, thus reducing the number of buffers by 2x. Figure 113-b depicts the optimization of an outgoing path, which is analogous to the previously described case but is in terms of the output signals.

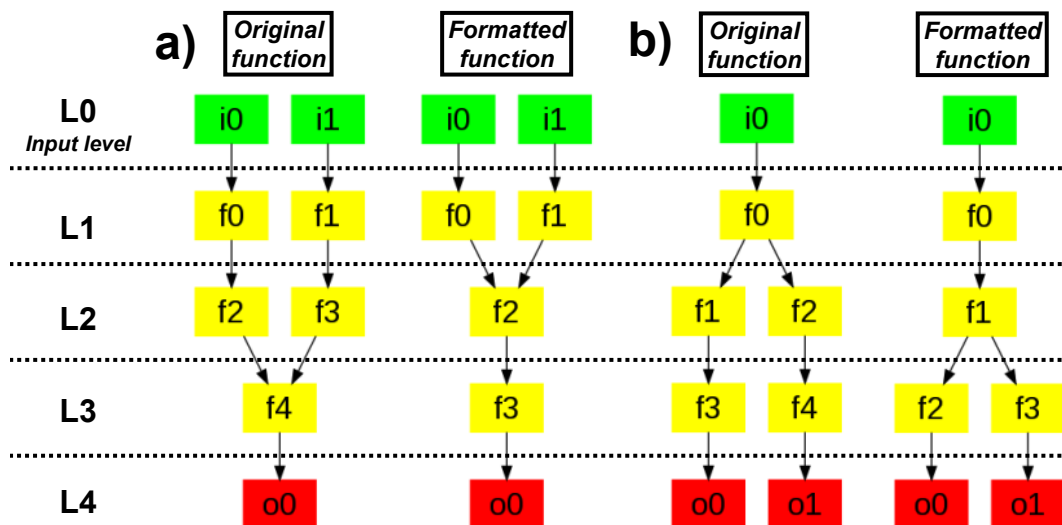


Figure 113. Multiple buffering/inverting path simplification illustration – a) incoming path – b) outgoing path

• **Multiple Output Correction:**

In the MCluster organization, the logic cells are connected to a limited number of other cells, i.e. fan-in/fan-out limitation. Hence, it is necessary to add buffers, in order to duplicate the signals in an architecturally compatible strategy. Figure 114 depicts an example case, where node  $f_0$  drives three different outputs. In this example, it is not possible to obtain more than two outputs per logic node. Hence, a buffer is added to reach the constraint.



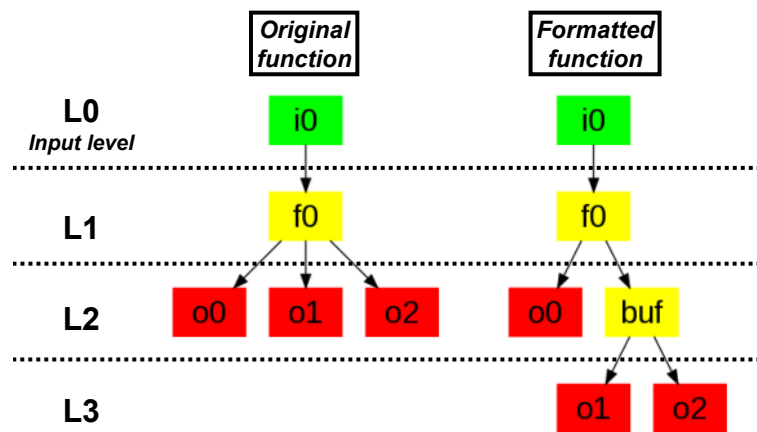


Figure 114. Illustration of multiple output correction

- **Output Layer Correction:**

In this final optimization step, all the outputs of the logic functions are placed on the physical output layer. In figure 115, we consider a logic function where the outputs are on layers L2 and L3. To meet the constraints of the architecture (here a matrix of depth 4), buffers are added to place the outputs on the layer L5.

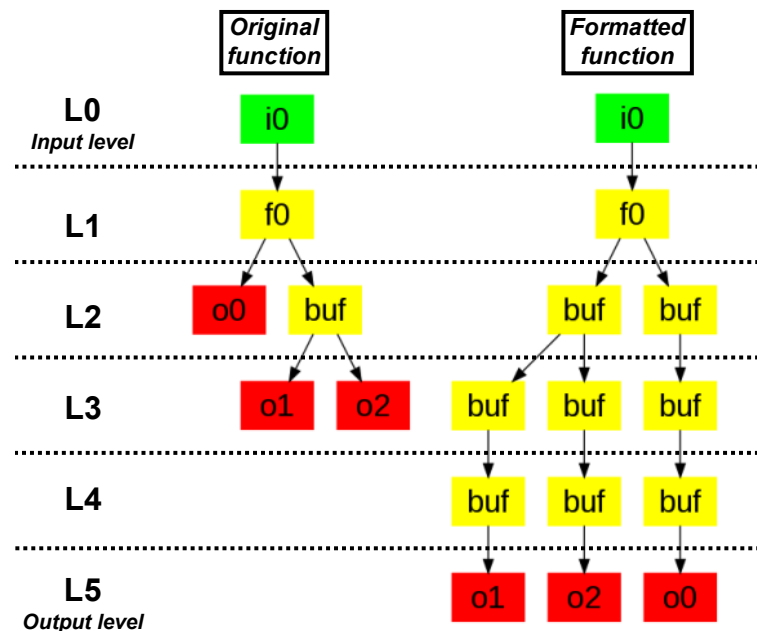


Figure 115. Output level correction illustration

### 6.3.2.3 Mapping Algorithm : Brute-Force Mapper

After the logic functions have been processed to meet the specific architectural requirements, the logic function netlist is brute-force mapped to the architecture netlist. The adapted netlist is analyzed in depth, meaning that for each node in the structure, child branches are identified and recursively explored. This depth exploration is performed without any consideration of the edge orientation in the graph. This allows identification of the connection between the nodes and initiation of the mapping sequence. Logic nodes are then assigned to cells. This is done according to the physical interconnections. Each layer's connections are compared to the relevant inter-layer architectural connectivity – allowing (or not) the assignment of functions to cells. Branching (i.e. the exploration of the immediately preceding alternative) is used when the arbitrary choice leads to a dead-end; the process is repeated until all functions are assigned to cells. The recursive algorithm of this step is shown in figure 116.

```

FUNCTION Mapping_procedure(current node)
  IF (all children/parents of current node already placed) THEN
    FOR (empty cell connected to children or parents) THEN
      Store the potential positions
    END FOR
  ELSE Store all the positions on the considered level
  END IF

  IF (No positions have been found) THEN
    MAPPING FAILURE
  END IF

  FOR (Each potential position)
    Place the current node at the position
    Ret=Mapping_procedure(Next node)
    IF (Ret is MAPPING FAILURE) THEN
      Unplace the current node from the position
    END IF
    IF (Ret is MAPPING SUCCESS AND No more nodes) THEN
      MAPPING SUCCESS
    END IF
  END FOR

  IF (No potential positions have been correct) THEN
    MAPPING FAILURE
  END IF
END FUNCTION

```

Figure 116. MPack mapping algorithm (pseudo-code)

As an example, we again consider a matrix which is four cells deep and four cells wide using a Banyan interconnect topology (Figure 109-a).

The example function to map is represented by a graph (shown in figure 117), generated by a random graph generator for test purposes. Its adjacency matrix is also shown in figure 117.

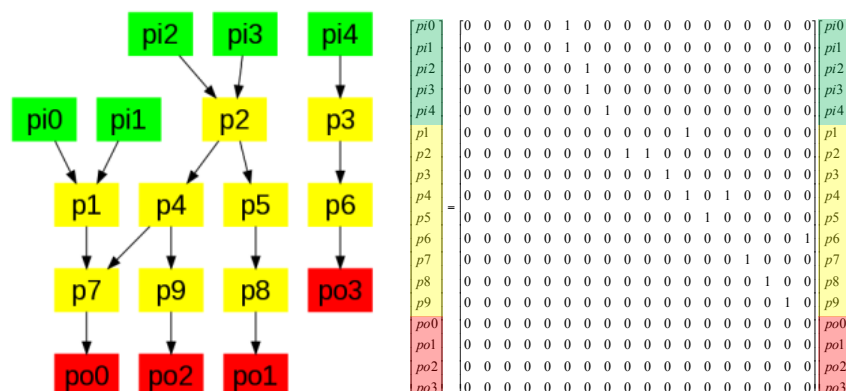


Figure 117. Function graph to map onto the Banyan\_4\_4 interconnect topology and its associated adjacency matrix

Following the previously described adaptation methodology, the original function graph is corrected in order to maximize the matching with the targeted MCluster. Hence, the position of inputs/outputs and inter-layer jumps will be corrected. The corrected graph is depicted in figure 118.

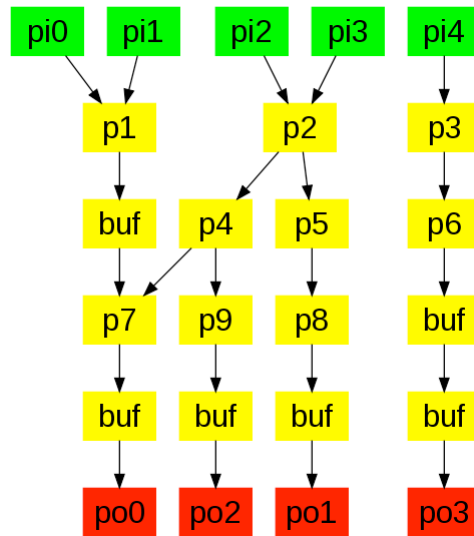


Figure 118. Function graph after correction step

The graph exploration is then launched and we obtain the following sequence (obviously, *buf* represents synchronization nodes):

*pi0-p1-pi1-buf-p7-p4-p2-pi2-pi3-p5-p8-buf-po1-p9-buf-po2-buf-po0-pi4-p3-p6-buf-buf-po3*

Finally, the graph assignment is performed. In the example, the first point  $p1$  is assigned to the cell  $f^{00}$ . According to the path defined in the previous step, a buffer is the next point to assign to a cell in the matrix. Since  $f^{00}$  is physically connected to  $f^{10}$  and  $f^{12}$ , the cell with lower  $y$ -index (here  $f^{10}$ ) is arbitrarily chosen for the buffer assignment, and the other possibility is memorized. In our example, the final programmed matrix is shown in figure 119. In this figure, we can see the nodes of the logic function graph and the nodes added for synchronization purposes correctly placed on the cell matrix.

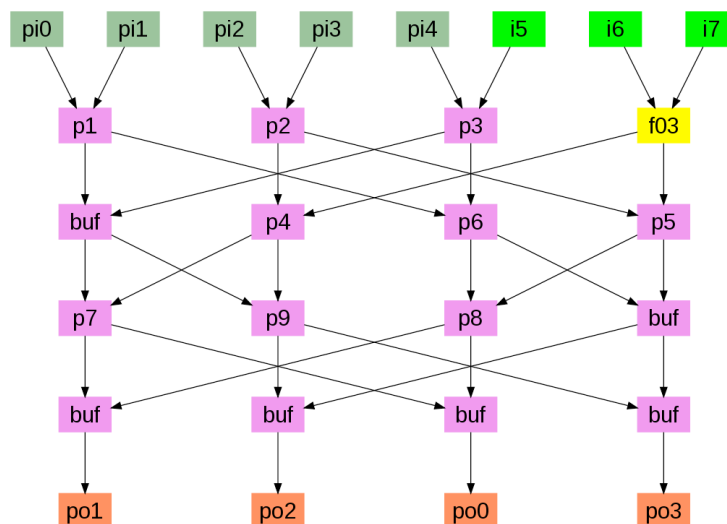


Figure 119. MCluster after function packing

#### 6.3.2.4 Clustering Algorithm

The clustering algorithm aims to group the logic cells into arrangements of MClusters. Such functionality is widely used in CAD tools.

There is large body of prior work on the packing problem for FPGAs. Most of it focuses on the optimization of area, delay, and/or power for the basic (LUT-based) soft logic complex

blocks. These algorithms include T-VPack [194], T-RPack [195], IRAC [196], HDPack [197], and others [198, 200, 199, 201].

In the MPack tool, we used algorithms derived from VPACK [17]. The tool constructs each MCluster sequentially, where the algorithm greedily packs the cells into MClusters. The pseudo-code for the algorithm is shown in figure 120.

It starts by choosing a seed cell for the current MCluster. As described by [17], the best way to choose the seed is to select the unclustered cell with the most used inputs. This approach prioritizes early placement of cells using the most cluster inputs, which are a scarce resource. Next, the tool selects the cell with the highest "attraction" to the current MCluster, which can legally be added to the current cluster. This evaluation of legality is performed by the previously described mapping algorithm. Essentially, we check if the cell added to the current MCluster leads to a valid and mappable MCluster. If the cluster is mappable, then the cell is definitively added to the current cluster. The attraction between a logic cell  $LC$  and the current MCluster  $MC$  is, as described in [17], the number of inputs and outputs they have in common:

$$Attraction(LC) = |Nets(LC) \cap Nets(MC)|$$

This procedure of greedily selecting the cells to add to the current MCluster continues until either (i) the MCluster is full or (ii) adding any unclustered cell would make the current cluster illegal. If the cluster is full, a new seed is selected and the packing starts for another MCluster.

```

Remaining_Cells: set of unclustered logic cells
C: set of cells packed in the current MCluster
MClusters: set of MClusters

MClusters = NULL
WHILE (Remaining_Cells != NULL) // More Cells to cluster
  C = Seed (Remaining_Cells) // Most Used Inputs and legal Cell
  WHILE (|C| < (MCluster_w.MCluster_d)) // Cluster is not full
    Best_Cell = MaxAttractionLegalCell(C, Remaining_Cells)

    IF (Best_Cell == NULL)
      BREAK
    END IF

    Remaining_Cells = Remaining_Cells - Best_Cell
    C = C U Best_Cell
  END WHILE
MClusters = MClusters U C
END WHILE

```

Figure 120. MPack' clustering algorithm (pseudo-code)

## 6.4 Evaluation of Fixed Interconnect Topologies

This section aims to evaluate the impact of the choice of interconnect topology in the MCluster architecture. The topology is evaluated considering performance metrics in the context of the packing method. It is worth pointing out that the evaluations proposed in this section are based only on the packing method presented above. The results are then

considered valid at the cluster level, and are used to validate the methodology. Large benchmarks, based on the complete tool flow, will be evaluated in the next section.

### 6.4.1 Methodology

Our analyses were carried out on an MCluster\_4\_4 using the previously mentioned intra-matrix interconnection topologies. We evaluated various metrics: the success rate of packing function graphs, the fault tolerance and the average interconnect length. We have carried out detailed analyses to compare the efficiency of the different intra-matrix interconnect topologies. We use a random graph generator to generate static sets of function graphs containing 6-16 points, in order to obtain fixed comparison criteria between topologies. No graphs contain isolated nodes, as here we focus on fixed interconnect layers, which are severely penalized by isolated nodes. Each set, which corresponds to a given number of points in the function graph, contains 1000 samples. Using the previously described packing method, each function is programmed onto the MCluster\_4\_4 using the various intra-matrix interconnect topologies (ideal or faulty) and metrics are calculated. MCluster\_4\_4 was an arbitrary choice because there is a good tradeoff between complexity and simulation time. Figure 121 summarizes the evaluation methodology and the associated parameters.

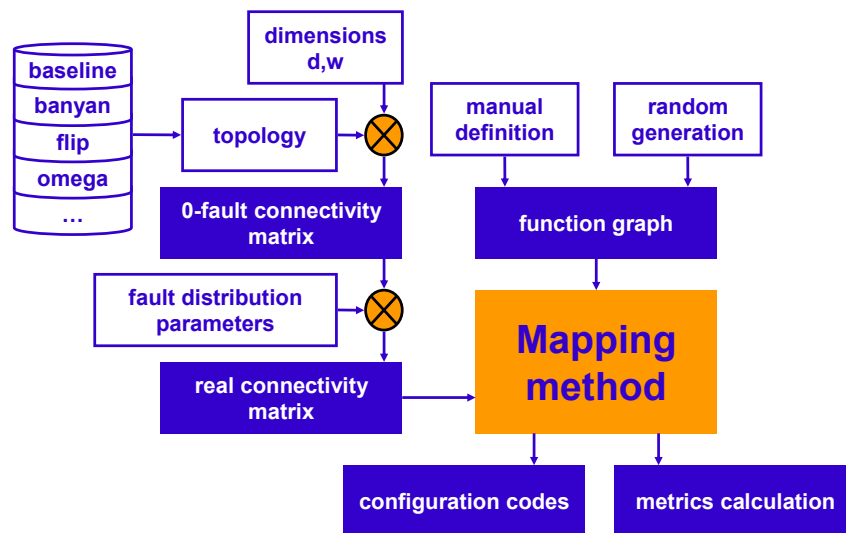


Figure 121. Performance evaluation method for fixed interconnection topologies

### 6.4.2 Packing Success Rate

Applying static sets to ideal interconnect topologies, we can test the ability of the matrix-topology ensemble to have complex functions packed onto it. Considering the percentage of function graphs successfully packed onto the matrix with respect to the number of samples in the set, we obtain the *success rate*. Figure 122 shows the comparison of success rates for Banyan, Modified Omega, Flip and Baseline topologies. For Banyan, Flip and Baseline interconnect topologies, the success rate is about 80% when the function graphs have 6 points. At 12 points, the success rate is about 25%. The difference between these three topologies is thus relatively small. However for the Modified Omega interconnect topology, the success rate is about 90% for 6-point function graphs and about 40% for 12-point graphs. This clearly shows that the Modified Omega interconnect topology is more suitable for this type of matrix.

This is because this topology has lower redundancy than the other topologies and spreads calculations over cells occupying less width, which seems to correspond better to typical function graphs. In fact in the matrix, there are pairs of cells which have the same inputs. For two cells which have the same inputs, the sum of the number of functions they can achieve is 14. For two cells which do not have the same inputs, the sum of the number of functions they can implement is  $14+14 = 28$ . In the Banyan topology for example, there are 6 pairs of cells which have the same inputs, while in the Modified Omega topology, every inputs are

different. This is the main reason why the Modified Omega topology has the potential to realize more functions than other topologies.

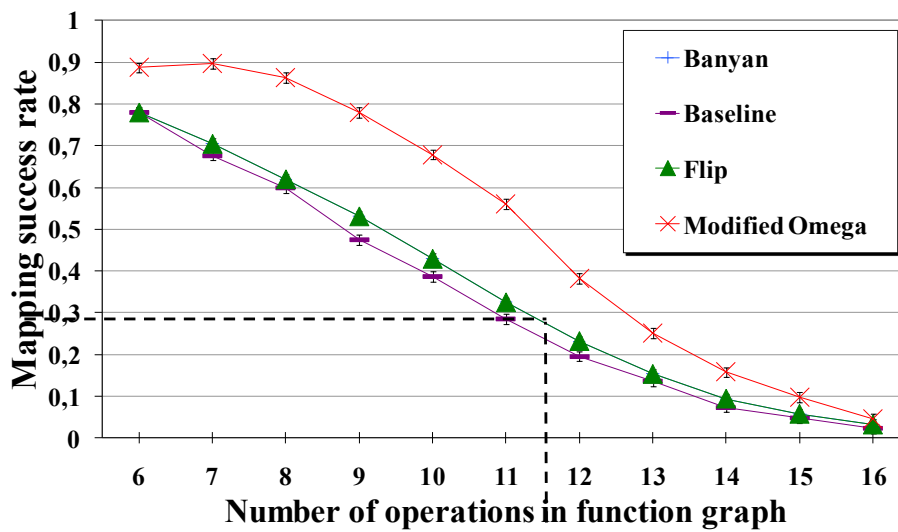


Figure 122. Programmability success rates for Banyan, Modified Omega, Flip and Baseline interconnect topologies within 4-deep 4-wide matrices

It is worth noticing that a traditional LUT approach would lead to a 100% mapping success rate, inasmuch as all possible combinational functions could be achieved by a LUT with a sufficient number of inputs. Then, the use of a MIN reduces the number of packed logic functions. Such a problem is managed at a higher hierarchical level. For example, if a function graph cannot be packed onto a single matrix, we can split it and map the sub graphs onto different matrices.

### 6.4.3 Fault Tolerance

#### 6.4.3.1 Physical Origin of Defects at Nanoscale

A significant concern in the use of nano-scale devices for computing architectures is the reliability of these individual devices and that of the resulting system. For instance, the chemical processes for building devices in a bottom-up approach will have significantly lower yields than those obtained via current fabrication practices, resulting in aggregates with high defect rates. For example, when considering transistors based on CNT, defects introduced by the CNT synthesis process could impact the behavior of the CNFET [202]. It is also necessary to consider that traditional processes have to shrink to build such devices. Due to variability, the overall system performance is thus likely to be dominated by unreliability as concerns individual device characteristics. Thus, it will be utopian to consider that interconnection layers will be defect-free. Most probably, at least one link could be destroyed or stuck in a Boolean state [203].

#### 6.4.3.2 Analysis Protocol and Results

Typically, unreliable systems use defect-avoidance techniques [53] to increase their reliability. Such techniques are based on defect detection and are packing-aware. To study fault tolerances in matrices, we must introduce defects in the interconnection layers, which will force the packing method to work around them. This will simulate the behavior we could achieve with a defect avoidance technique. With these considerations, physical interconnection defects were simulated by randomly deleting links in layers. Cell non-functionality is simulated by deleting all related input/output links. By so-doing, we assume that the faults are only static and due to fabrication processes issues, i.e. they do not change dynamically during runtime. It is also worth noticing that, by using defect-avoidance techniques, we assume that the faults are testable.

Using the Banyan topology as a reference, figure 123 shows the comparison of success rates for Modified Omega and Banyan topologies in the cases of one or two faulty links on the first

layer and one faulty cell, which are the most representative cases. For all topologies, when the function graphs have 6 points or 16 points, faults have no influence on the success rate, because the interconnect topologies are not a decisive factor at these limits. For 12 points, the success rate of the Modified Omega topology falls faster than for the Banyan topology. The Banyan topology can thus be considered to be more robust than the Modified Omega topology in terms of sensitivity to faults, but in all situations, the absolute success rate for the Modified Omega topology remains higher than that of the Banyan topology.

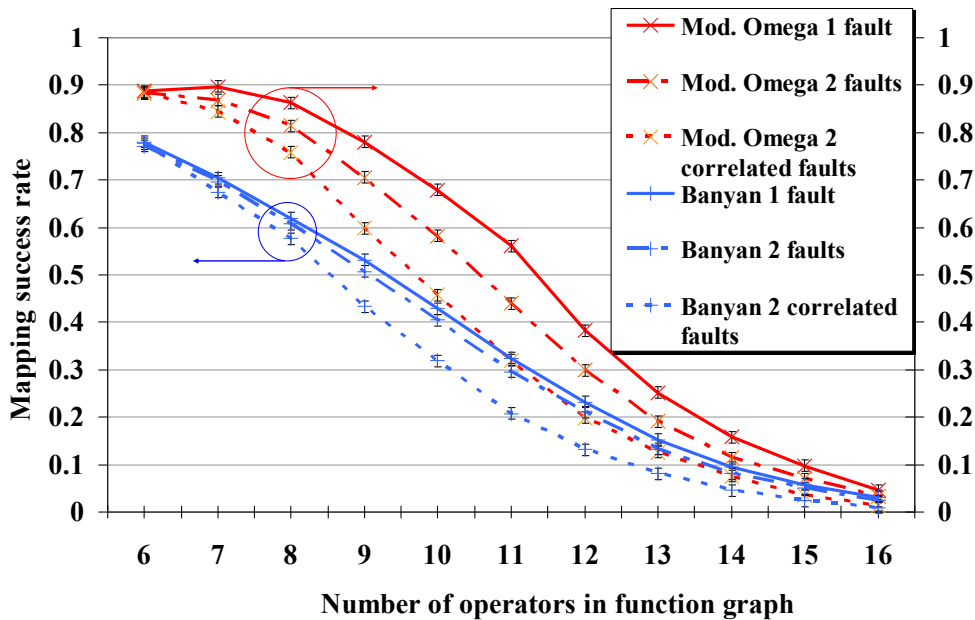


Figure 123. Comparison of programmability success rates for Banyan and Modified Omega interconnect topologies within 4-deep 4-wide matrices in the case of faulty links and faulty cells

### 6.4.4 Average Interconnect Length

Arbitrary coefficients are assigned to interconnect based on how long they are. Given an interconnect between cells at coordinates  $(x, y_1)$  and  $(x+1, y_2)$  respectively, the coefficient value is given by  $|y_2 - y_1| + 1$  (assuming a Manhattan-based interconnect geometry).

Directly after the packing process, it is therefore possible to sum the coefficients of all used interconnects. This method gives information about the interconnect length used in the system. Considering layout and design rules, it is possible to extract parasitic capacitances, and therefore to deduce the maximum working speed of one pipeline stage.

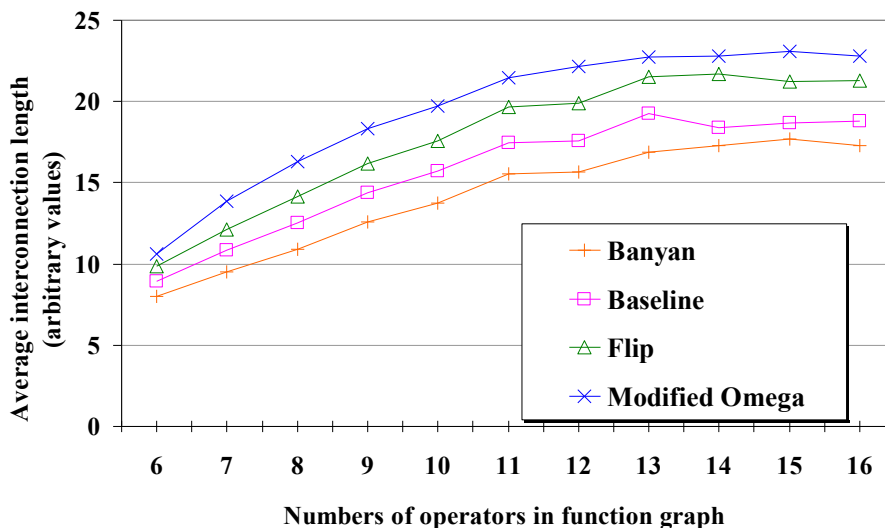


Figure 124. Comparison of average interconnect length for Banyan, Modified Omega, Flip and Baseline interconnect topologies within 4-deep 4-wide matrices

Figure 124 shows the average interconnect length comparison. The Banyan topology displays the lowest interconnection length due to its straight-ahead paths. However, the Modified Omega topology, with its long shuffle links, is the worst case with 25% greater overall interconnects than Banyan. This shows that the average system speed is likely to be higher with the Banyan topology than with the others.

### 6.4.5 Discussion

The results, obtained above, are summarized in Table XXIV. The metric “Performance” is calculated as the inverse of the average interconnection length and all the values are normalized considering the use of the Banyan topology for a 12-node graph as reference, such that

$$Metric_{Normalized}(X) = \frac{Metric_{at12nodes}(X)}{Metric_{at12nodes}(Banyan)}$$

Table XXIV. Analysis summary

	Packing efficiency		Performance (1/int. length)
	no faults	with faults	
<i>Banyan</i>	1	1	<b>1</b>
<i>Baseline</i>	0.84	0.86	0.89
<i>Flip</i>	1	1	0.78
<i>Modified Omega</i>	<b>1.65</b>	<b>1.42</b>	0.71

The Modified Omega topology gives the best results in terms of success rate and fault tolerance, whereas it appears to be the worst for performance. However, considerations about performance will require deeper analyses considering the layout, and real values of extracted parasitic capacitances. The three other topologies analyzed are quite similar in results. As mentioned above, these results might be explained by the redundancy of the interconnect topologies. In fact, highly redundant topologies will give good results in terms of fault tolerance, while their packing success rate is less convincing.

## 6.5 Global Architecture Evaluations

### 6.5.1 Methodology

In this section, we evaluate the performance of the proposed architectural scheme. We will use the previously described toolflow and map the MCNC and ISCAS’89 benchmarks [204] to the target architecture. Different scenarios will be used to evaluate the impact of the granularity and the latch depopulation, and to compare the global performance with the CMOS counterpart. The CMOS reference architecture is formed by CLBs organized with 10 BLEs of 4-input LUTs. 22 inputs will then come from the external routing structures into the block. The physical parameters of the different structures are extracted using 65-nm Back-End data. The used metrics are the area and the critical routing delay. These metrics are computed during the Place and Route iteration of the flow. The area corresponds to the sum of the logic area, i.e. the area of used CLBs, and the routing area, i.e. the area of the used routing resources. These values are expressed using the  $\lambda$  parameter, which is defined as half of the minimum gate length. All areas are normalized with respect to the largest circuit implementation. The critical routing delay corresponds to the most constrained delay through the routing structures, i.e. the longest path external to the CLBs.

### 6.5.2 Evaluation Results

#### 6.5.2.1 Impact of the Granularity

Figure 125 depicts the area estimation of an FPGA, following the proposed architectural scheme. While the MClusters are totally populated in terms of latches, we study the impact of their granularity by varying the size of the clusters. The best results are obtained with MClusters\_3\_3, where an improvement of 12% on average, compared to a 2 by 2 granularity,



can be observed. It is worth noticing that for a smaller, as well as for a higher granularity, the used area increases drastically. For large clusters, this can be explained by the large size of MClusters with respect to their poor mapping ability. In such cases, the incomplete interconnectivity becomes a major hurdle, because of the large number of wasted cells used only for buffering purposes. On the other hand, when a 1 by 1 granularity is used, the amount of external circuitry required for the routing heavily dominates the area, as compared to the computation cell. This leads to a large area overhead, and confirms the initial hypothesis of this work, where the direct transposition of an FPGA architecture to ultra-fine grain would lead to a large interconnection overhead.

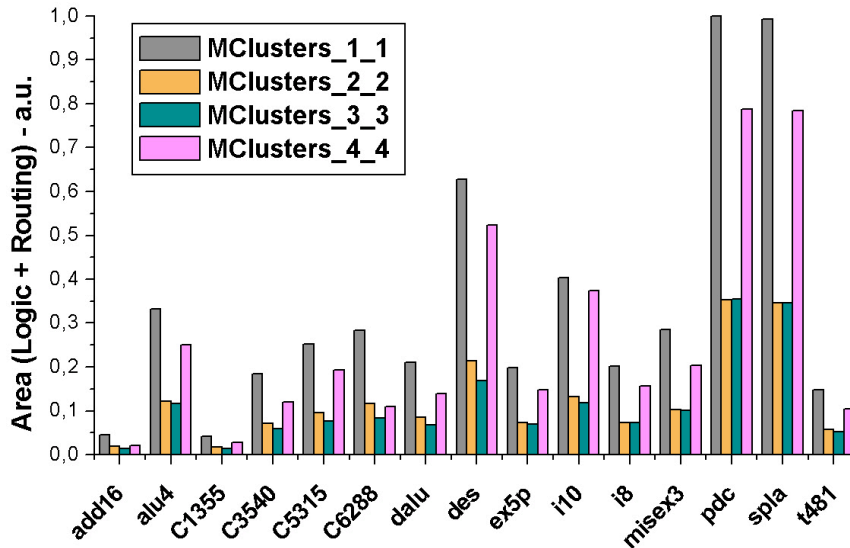


Figure 125. Area estimation for MCluster-based FPGAs with various granularities

### 6.5.2.2 Impact of Latch Depopulation

Figure 126 depicts the area estimation of the MCluster-FPGA with a 3 by 3 granularity and latch depopulation.

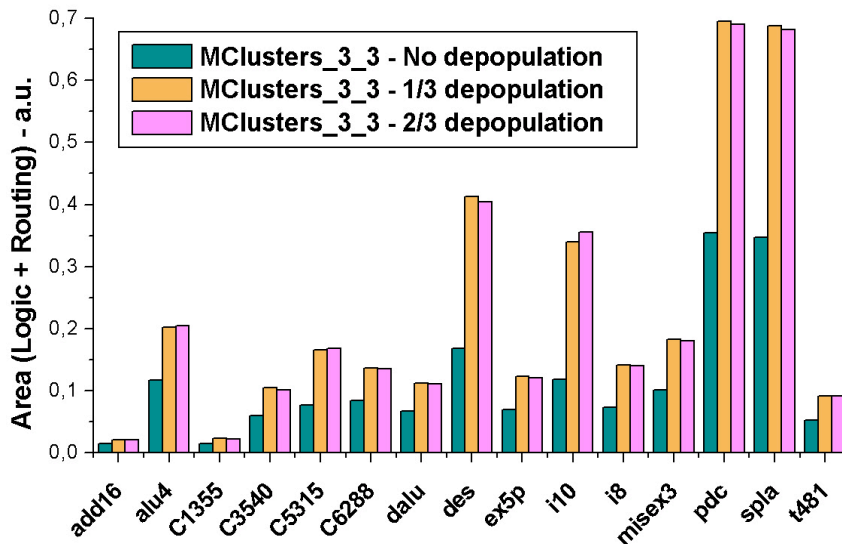


Figure 126. Area estimation for MCluster\_3\_3-based FPGAs with different latch depopulation scenarios

The depopulation ratio corresponds to the number of latches that have been removed in the structure, with respect to the maximum number of possible latches. We can observe that the use of latch depopulation is not a good strategy for area performance, because of a nearly 2x increase in used area. While the idea of removing some latches is interesting to handle the size of the cluster (the size of a latch is much greater than that of an elementary cell), it appears detrimental to the routing efficiency. In fact, the position of the latch is not taken into account by the MCluster packing tool. This leads to a less efficient placement iteration, and consequently to the inefficiency of the depopulation.

### 6.5.2.3 Performance Comparison with CMOS

Figure 127 depicts the area estimation for a 3 by 3 granularity MCluster-based FPGA and compares it to its CMOS counterpart. The benchmarks show an area reduction ranging from 38% to 62%, with an average of 46%. This is clearly due to the low area impact of ultra-fine grain logic cells, compared to the rather larger area required by a CMOS LUT. In fact, a 3 by 3 cluster costs an area of  $63433\lambda^2$  (with  $\lambda$  the half of the minimum gate length) while a 4-input LUT costs  $54292\lambda^2$ . It is worth noticing that, while their size remains in the same order of magnitude, the functionality of an MCluster is much higher than that of an equivalent LUT. With the same size, MClusters can input 50% more data and can output 3x more results. When correlated with the efficiency of the packing tool for matrix clustering, this demonstrates a clear advantage of the complete proposal as compared to the CMOS approach.

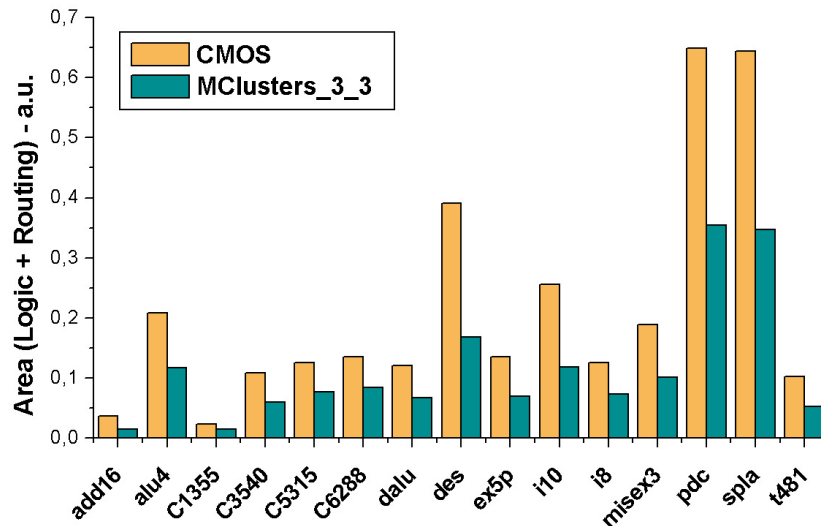


Figure 127. Area estimation for 3 by 3 MCluster-based and CMOS-based FPGAs

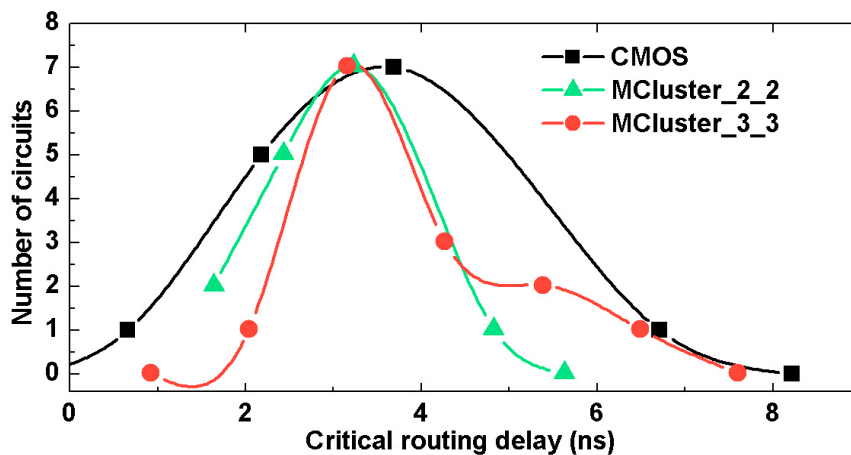


Figure 128. Critical routing delay repartition for 2 by 2, 3 by 3 MCluster-based and CMOS-based FPGAs

Figure 128 depicts the distribution of the critical routing delay through all the benchmarks using three different scenarios. While the average delay is not changed significantly (delay reduction up to 10%), it is worth noticing that the use of MCluster structures allows a reduction in the standard deviation of the delay distribution. In fact, we can observe that, while the CMOS distribution is quite large, the use of a 2 by 2 MCluster implies a lowering of the extremes, and tends to globally improve the performance of the mapped circuits. This behavior can be explained through the global impact of ultra-fine granularity on the benchmarking toolflow. In fact, ultra-fine granularity induces a predominance of local inter-CLB interconnect instead of long wire connections, thus leading to the reduction of long critical paths.

### 6.5.3 Discussion

In this chapter, we studied how logic blocks with compact size and ultra-fine granularity might be used in reconfigurable systems.

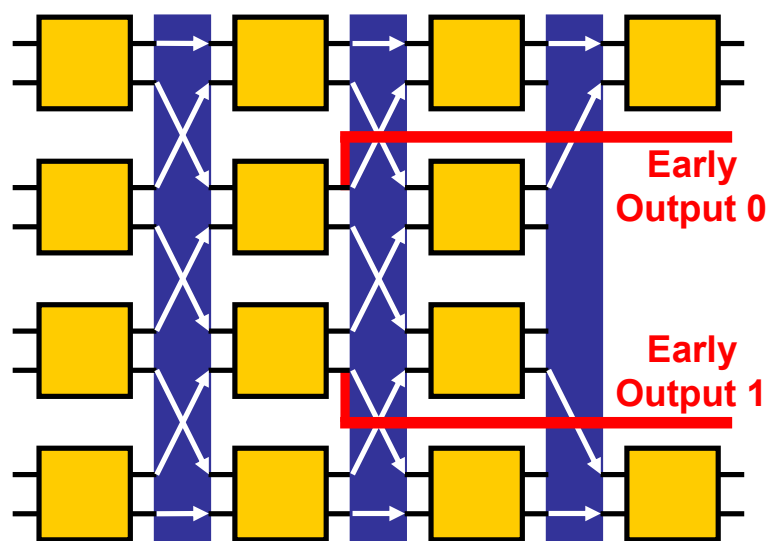
To manage the logic abstraction and to adapt the granularity, we proposed a layered structure, derived from conventional FPGAs. Nevertheless, it was necessary to add a new layer in order to build a fine-grain combinational block from a population of ultra-fine grain logic elements.

In the envisaged architecture, ultra-fine grain logic cells are arranged into matrices, called MClusters. To prevent inherent interconnect overheads, the cells are interconnected through fixed and incomplete topologies. Several topologies could be envisaged. The MClusters are further organized into BLE/CLBs, with complete interconnect and the CLBs are arranged in islands to form the FPGA.

To allow the performance evaluation of this new architectural FPGA scheme, we proposed a complete benchmarking toolflow. While this flow is based on existing and traditional tools, coming from the VTR flow, we add specific functionalities to meet the specificities of the architectures. In particular, the Matrix Packer (MPack) tool has been proposed to pack netlists of logic cells into netlists of MClusters. This tool is thus able to model the MCluster architectures, in terms of several parameters, such as the topology and the dimensions.

This tool has been used to study, in a first approach, the impact of interconnect topology on the structure, and we showed that the mapping success rate is about 90% for 6-point function graphs and about 40% for 12-point graphs when using the Modified Omega interconnect topology in an MCluster<sub>4\_4</sub>. Moreover, we have shown that the Modified Omega interconnect topology is still more robust than other topologies in the presence of defects. In contrast, the Modified Omega topology is not a good choice for wire size, displaying 25% greater overall interconnect length within the cell than the Banyan topology.

It is worth noticing that several other topologies can be investigated. In other work, a cross-cap pattern has been proposed for the matrix arrangements [205, 206] to overcome logic depth and data directivity concerns. Another solution to constant logic depth, where the layered matrix structure requires a large number of buffers to correctly output results, is to modify the structure to output some of the internal nodes of the matrix and delete some cells often used for only buffering, as shown in figure 129.

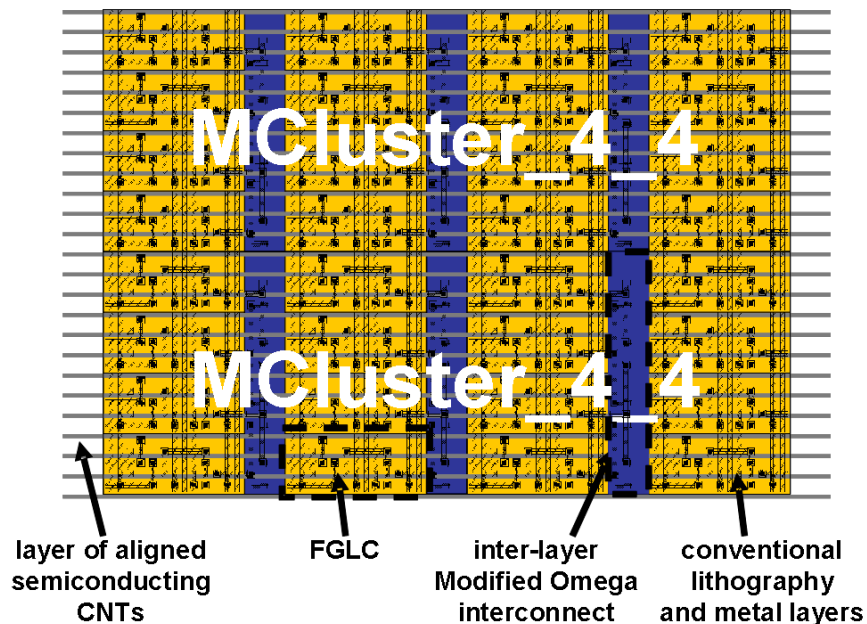


**Figure 129: Modified-Omega based MCluster with early internal outputs and cell depopulation**

From an integration point of view, it is worth noticing that a matrix organization is particularly well-suited to the structural regularity required by CNFET fabrication processes. Matrix patterns fundamentally exhibit structural regularity, which is easily transposable to highly regular fabrication processes. For example, in the case of in-field reconfigurable cell,

even if there are many fundamental CNFET integration issues which are still unsolved, the architecture is compatible with recent proposals and fabrication results [207, 55], since it is possible for long bundles of nanotubes to span several cells in the same column (Figure 130).

In a second approach, we analyzed the global architecture through large benchmarks using the developed methodology. We evaluated the potential of the architecture in terms of MCluster granularity and latch depopulation, and compared the structure to a CMOS FPGA. We showed that the best granularity for used area is a 3 by 3 cluster, while the latch remains entirely populated. We finally demonstrated an improvement with respect to CMOS with an average area saving of 46%, and we discovered that the routing delay is less distributed and tends to be more controllable than in the traditional approach.



**Figure 130.** Illustration of two MClusters\_4\_4 physical implementation on parallel carbon nanotubes layer

While the proposed solution appears particularly interesting for future FPGAs, it will also be of high interest to study other architectures, which go beyond conventional CLB organization. In fact, the approach presented in this chapter is quite close to the traditional FPGA scheme. In figure 131 and figure 132, we present a more speculative proposition for the Configurable Nano Logic Block. Figure 131 shows an MCluster\_4\_4 connected to four sequential elements, dedicated to latch the outputs of the cluster. Peripheral circuitries are taken into account for configuration. Since we also envisage non-volatile memories, it is possible to use local heavy-duty circuitry to allow serial loading of configuration data, thus reducing the number of interconnect metals lines that are used for configuration. Figure 132 shows the organization of the CNLB. We propose the use of a logic pipeline. In this arrangement, the logic structure presented in figure 131 is duplicated many times and cascaded into many pipelined paths. A configurable interconnection is possible between the pipelines. This organization can be considered as another layered arrangement, similar to that of the MCluster. Furthermore, management circuitry will be added to this circuit, to ensure service tasks, such as scheduling or synchronization between tasks and paths through the pipeline, as well as ensuring the correct behaviour of the whole. For example, a defect mapping operation could be in charge of a dedicated state-machine-based test circuitry to return the information to upper-level blocks in charge of defect avoidance.

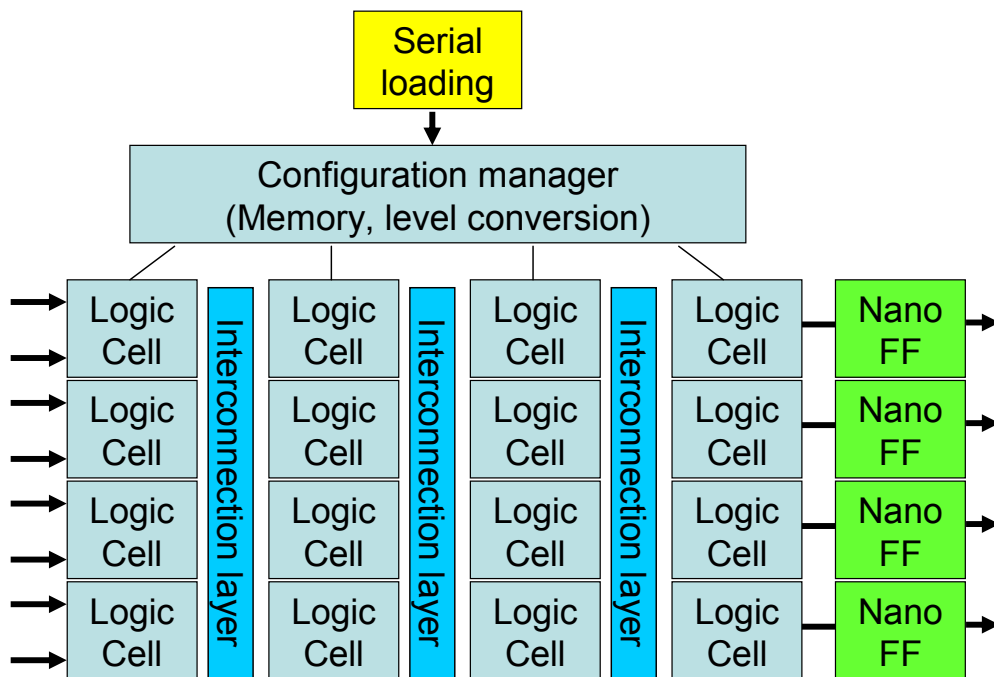


Figure 131. MCluster organization for speculative CNLB

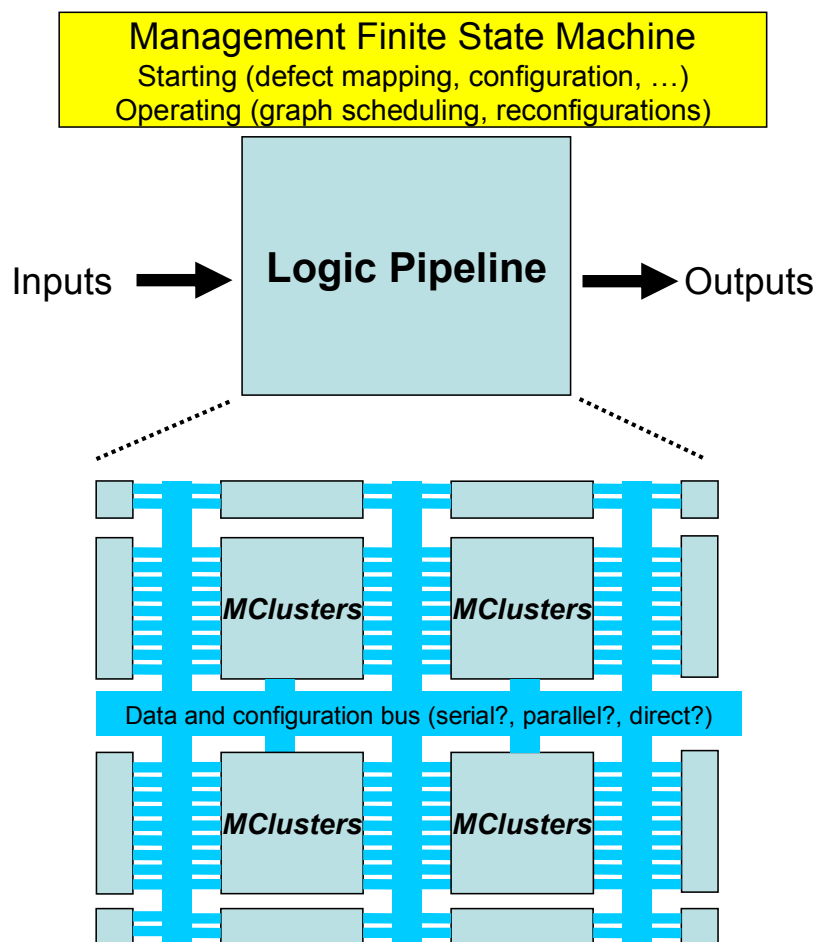
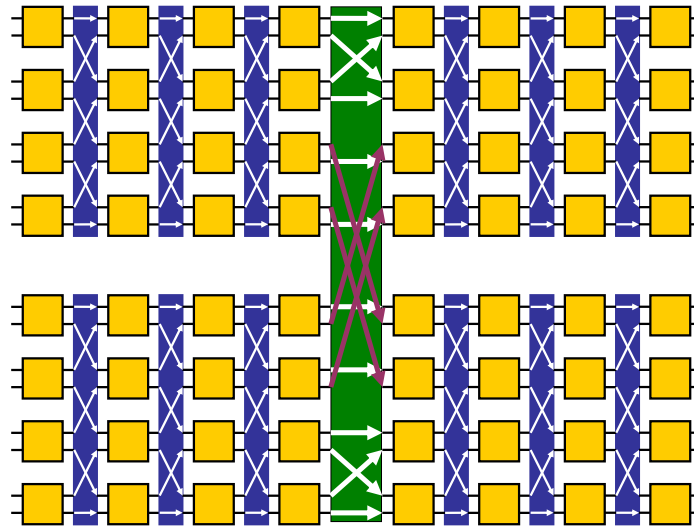


Figure 132. Speculative pipelined organization of CNLB

Another final possibility is also to pursue the straightforward MCluster organization. For this case, figure 133 shows a double-stage MCluster organization with four MClusters<sub>4\_4</sub> interconnected using a Modified Omega topology (Blue interconnect). The four internal MClusters are then packed into another MCluster organization of two blocks in depth and two in width. These blocks are then interconnected by a Modified Banyan topology (Green interconnect), which aims to interconnect the various blocks from side to side (purple arrows),

while maximizing direct interconnections. This final proposal helps to show that the design exploration space, reachable by the proposed benchmarking flow for ultra-fine grain, is both broad and promising.



**Figure 133. Double-stage MCluster organization with two-level hierarchic interconnection strategy**

Finally, several routing arrangements can also be envisaged instead of the traditional mesh pattern. In [178, 179], the authors have proposed a convenient hierarchical organization for intra-CLB routing. These approaches are based on the use of MINs and fixed interconnection topologies that could efficiently replace the complete CLB interconnect. Thus, an interesting opportunity could be found in mixing the hierarchical approaches for routing (i.e. hierarchical routing) and logic (i.e. MCluster organization).

## 6.6 Conclusion

In this chapter, we dealt with the problem of architectural organization for emerging logic circuits. Reconfigurable logic uses the largest part of its area for “peripheral circuitry”. It is then of high interest to explore the implication of a drastic reduction of the logic part. While it is not obviously of first order, this exploration is motivated by the fact that reducing the logic part may lead to changes in interconnect strategies and consequently to new architectural organizations which improve the global performance of the reconfigurable device. Specifically, emerging technologies lead to the emergence of highly compact logic cells, such that using conventional reconfigurable computing schemes with ultra-fine granularity would lead to a large interconnect overhead. To prevent this, we propose the use of matrices of logic cells, interconnected by layers and using a specific fixed and incomplete interconnection pattern. This new combinational computation block is then packed with FFs to realize a BLE scheme. It is worth noticing that not all outputs will be latched, in order to reduce the size of the BLE. The final island-style remains the same.

In order to benchmark the performance of the proposed circuit, we enhanced the conventional toolflow. In particular, we added a specific packing tool, called MPack. This tool is intended to perform a packed netlist of logic cells into netlists of MClusters. The tool is based on two principal algorithms: greedy clustering and brute force mapping. The clustering chooses which blocks must be packed together, while the cluster is legally checked by brute-force mapping onto the target architecture. With this toolflow, it is possible to benchmark the proposed structure.

This method has thus been used to analyze intra-matrix topologies with respect to various metrics, and demonstrated that the mapping success rate is about 90% for 6-point function graphs and about 40% for 12-point graphs when using the Modified Omega interconnect topology in a 4x4 matrix. Moreover, we have shown that the Modified Omega interconnect topology is still more robust than other topologies in the presence of defects. In contrast, the

Modified Omega topology is not a good choice for speed, displaying 25% greater overall interconnect length within the cell than the Banyan topology.

We subsequently evaluated the potential of the architecture in terms of MCluster granularity and latch depopulation, and compared the structure to a CMOS FPGA. We showed that the best granularity for used area is a 3 by 3 cluster, while the latching remains entirely populated. We finally showed an improvement with respect to CMOS with an average area saving of 46%, and we discovered that the routing delay is less distributed and tends to be more controllable than in the traditional approach.

To summarize, ultra-fine grain architectures appear to be highly relevant for future computation systems. Furthermore, many logic cell arrangements and interconnect topologies are possible at several levels. With the help of complete architectural exploration tools, the field could reveal interesting opportunities for future systems.

# CHAPTER 7 *Conclusions and Future Work*

---

The global objective of this work was to provide a framework, suited for the comparison of different emerging technologies. Within a reconfigurable architecture framework, several technologies have been benchmarked and compared. In this concluding chapter, we provide a brief final comparison between the different architectures. This global comparison extracts some guidelines and also provides feed-back to the technologists. We also summarize the overall contributions and the results of the work, and discuss future opportunities to pursue the work.

## 7.1 Global Comparison

Table XXV summarizes the metrics for the different architectures. The first three proposals correspond to the conventional FPGA architecture, which was improved in terms of memories and routing circuits (chapters 3 and 4). The last architecture is based on ultra-fine grain logic (chapter 5) in a specific MCluster-based organization (chapter 6). The depicted gain figures are expressed with regards to the baseline CMOS 65-nm FPGA architecture. The gain in area comes from the contribution to the logic and routing parts. The critical path delay corresponds to the slowest signal in the design. This metric is thus directly correlated to the maximal achievable performance of the system. It should be noted that the critical path delay gain is not given for ultra-fine grain architectures, due to the unavailability of the timing evaluation in the latest benchmarking tool (6.3.1).

**Table XXV. Summary of global gain figures for the proposed architectures with respect to a baseline CMOS 65-nm FPGA architectures**

	<b>Total Area Gain</b>	<b>Critical Path Delay Gain</b>
<i>PCM</i>	13%	44%
<i>Monolithic 3-D</i>	21%	22%
<i>Vertical NWFET</i>	46%	42%
<i>DG-CNFET Ultra-Fine Grain</i>	46%	-

Regarding area, it is worth noting that a gain of about 50% is achievable in two different contexts. The first means to improve area is to reduce the impact of routing and memory resources. With 3-D techniques (e.g. with PCM, Monolithic 3-D and vertical FETs), it is possible to reduce the size of these peripheral circuitries. With highly integrated 3-D circuits such as Vertical NWFET, all these circuits can mostly be stacked above the data path logic, leading to an area reduction of peripheral circuits of up to 46%. The second means investigated deals with the architectural organization of FPGAs. Emerging technologies have been envisaged to realize very compact logic gates. This leads to new computation based on ultra-fine grain logic blocks. Again, we demonstrated a gain of 46% in terms of area. This improvement is due to the repercussions of the new architectural logic block organization. In



fact, the use of ultra-fine grain logic cells and the MCluster arrangements yields a compactness of logic requirements compared to equivalent size LUTs, such that the number of required CLBs is significantly reduced, and the number of routing resources accordingly. It is worth pointing out that the improvement of area in this situation is of a different nature than that of the previous examples. Indeed, in the latter architecture, we work on the logic block size and their arrangement. The organization of routing resources is however identical to that of CMOS. This must be considered with regards to the other technologies, which work to improve routing resources. Hence, we could say that working towards a severe reduction in routing resources leads to similar area improvements as skillful logic reorganization. This is particularly remarkable if we consider that peripheral circuitries occupy the largest part of the structure, while computing resources are relatively insignificant. Finally, it will be of high interest to merge the different technologies to improve all the limiting factors of the different blocks.

Regarding critical path delay, we should observe that numbers have not been extracted for ultra-fine grain architectures. With the current benchmarking tool, complete timing evaluation is not yet available. More precisely, it is not currently possible to evaluate the inter-logic timings. In fact, only studies on routing delay are currently possible. It is thus unfair to directly compare the ultra-fine grain technology to the others. Technologies that are working on routing resource improvement are seen as highly positive for global circuit performance. In the first three proposals, the performance of data path logic has been improved. In particular, the technologies allow the on-resistance of the pass switches to be reduced. This leads to a delay reduction of up to 44%, for the lowest on-state resistive device. With respect to the ultra-fine structure, it is worth noticing that the granularity improves the routing delay globally. In fact, the architecture and its associated mapping methodology favor the local interconnect. This leads to a promising global improvement of the performance levels. Finally, it is again possible to expect a symbiosis between the different approaches, in order to improve the specific routing circuits in the ultra-fine grain architecture. This will certainly provide a global improvement of the delay figures.

## **7.2 Conclusions and Contributions**

### **7.2.1 Global Conclusions**

In this thesis, we intended to evaluate the potential of emerging technologies for future computing application. This topic is of course so broad that significant effort will be invested in the costly development of emerging technologies even if the system perspective cannot be foreseen. In order to reduce the lead-time and the costs, we explored a fast evaluation methodology. This methodology is based on a generic architectural template (chapter 1). This was chosen to be a generalized FPGA architecture, since its high flexibility makes the circuit well-suited to several applications. After building a global overview of current and emerging reconfigurable architectures, we formalized the architectural template that will be the baseline of our study (chapter 2).

The standard reconfigurable architecture drawbacks are the routing and memories that are used to configure the whole architecture. We thus studied how emerging technologies can improve the configuration memory nodes, as well as the routing resources (chapter 3). Several memory and configurable routing nodes have been proposed using 3-D based technologies. These circuits are expected to improve the global FPGA architecture. These considerations have been addressed by benchmarking a structurally enhanced FPGA (chapter 4). The benchmarking is performed by using the traditional FPGA flow and by tuning the architectural description.

While improvements of memories and routing resources appeared as the most convenient way to optimize the FPGA structure, we also assessed the logic block circuits (chapter 5). In fact, computation blocks in standard FPGA systems occupy less than 20% of the total area. We thus investigated new computation cell paradigms that to enable the exploration of new

architectural organizations in which the predominance of the logic blocks is increased. Several technologies were investigated to compact the elementary logic block, leading to ultra-fine grain logic. The standard reconfigurable template was enhanced to take into account the novel ultra-fine grain paradigm. We provided a specific organization for ultra-fine grain logic FPGAs, based on matrices of cells interconnected by fixed and incomplete topologies (chapter 6). In order to evaluate this new scheme, a specific packing tool was developed and linked to the standard benchmarking flow (chapter 6). The ultra-fine grain architecture was therefore studied, in terms of topology selection and sizing, and was compared to its equivalent CMOS circuit (chapter 6).

### **7.2.2 Contributions to Methodology and Tools**

To allow fast and low-cost assessment of emerging technologies, we developed a new evaluation methodology. In fact, we defined a standard and generic architectural template that was intended to be efficiently realized by disruptive technologies. Thus, it became possible to provide a full set of evaluation tools that allow the evaluation of a wide range of technologies for a broad range of application contexts (chapter 1). To do so, we proposed a complete tool flow, able to instantiate our generic template customized with several technology parameters (chapters 4 and 6). Based on the conventional flow established for the exploration of FPGA architectures, several tools were connected together to allow the evaluation of (i) standard FPGAs, (ii) reconfigurable circuits with enhanced routing, (iii) reconfigurable circuits with gate-based computation blocks (instead of LUTs) and (iv) reconfigurable circuits with a specific architectural organization. In this final context, a specific packing tool was developed in order to allow regular arrangements of logic cells with special interconnection topologies (Chapter 6). This tool (called MPack) is fully compatible with the standard tools and has demonstrated performance in terms of computation time and mapping results compatible with its design exploration purpose.

### **7.2.3 Contributions to Memories and Routing Resources for FPGA**

In FPGAs, routing resources and memories occupy more than 80% of the chip area. Memories are distributed all over the circuit and are used to program the computational logic blocks and the active routing elements. Memories are traditionally implemented using SRAMs, which are costly in terms of area. Routing resources are used to create paths between logic blocks. The introduction of many active elements through the logic paths has two obstacles. Firstly, they are highly area consuming since they must be built in front-end silicon. This obviously also has an impact on interconnect size, since the lines must constantly connect back-end to front-end active layers and so on. Secondly, the introduction of active devices into the logic data path tends to globally decrease the performance. In order to assess these questions, we explored the interest of three 3-D based technologies (chapter 3): resistive phase-change memories, monolithic 3-D integration process and vertical NWFET.

Resistive memories are non-volatile passive devices that can be programmed between two stable resistive states. These devices are compatible with the back-end and can be embedded within the interconnection layers. A configuration memory node was proposed, based on resistive memories arranged in a voltage divider, in order to store intrinsically a logic voltage level. The node improves its equivalent non-volatile flash counterpart by 1.5x in area and 16.6x in write time, while it improves on an emerging magnetic memory based equivalent by 3.8x in area. Resistive memories are also remarkable for their low on-resistance (up to  $50\Omega$ , with respect to  $9.1\text{k}\Omega$  for minimum-sized CMOS 65-nm transistors). Thus, it is of high interest to introduce the device directly into the logic data paths, to yield high-performance switches. It is thus possible to create routing elements that combine pass-gate and configuration memory into a unique two-terminal resistive node. We demonstrated that this node leads to an improvement in area of 3.4x compared to flash and 14x compared to magnetic memories.

The monolithic 3-D integration process aims to stack several layers of active silicon with high via density. This appears promising with respect to the previous technology that stacks only

passive devices. We assessed the performance of FPGA elementary building blocks realized in a full 3-D scheme. Monolithic 3-D technology exhibits high alignment accuracy. This makes it possible to envisage splitting between the configuration memory and the data path logic. In standard 3-D techniques, this is not possible due to the large number of interconnections that are required between the two levels. We implemented a simple 2-bit LUT element and various cross point architectures. We showed that the LUT2 can be improved by 2.03x in area compared to an equivalent 2-D circuit. It has also been shown that the intrinsic delay is improved by 1.62x, the load factor by 6.11x and the average power by 2.02x. This tendency is also observed with the routing cross points where gains of 2.93x in area, 3.1x in intrinsic delay, 1.07x in load factor and 2.1x in average power were demonstrated.

More prospective vertical NWFET technology enables transistors to be entirely implemented between two metal layers. This opens the way towards large high-performance transistors with small footprint impact, due to the vertical orientation of the active volume. This technology also enables the realization of complete logic functions within the metal layers. In particular, we proposed a configurable via that switches on under a Boolean (single- or multiple-input) condition. It is also possible to realize complex logic gates, and so to place signal buffers, logic functions and configuration memories above the initial circuit. To evaluate this disruptive technology, we proposed a methodology based on TCAD simulations. After modeling the devices physically, transient simulations were carried out to extract the performance metrics of simple circuits. We demonstrated that such a technology yields an improvement of 31.2x in area, 2.5x in intrinsic delay and 14.5x in leakage power with respect to the equivalent CMOS gate.

#### 7.2.4 Contributions to Logic Blocks for FPGA

The largest part of the FPGA architecture is occupied by “peripheral circuitries”. In an incremental design approach, it is reasonable to focus on this part. Nevertheless, new computational block architectures should lead to the emergence of new system architectures, and potentially to a reduction of the imbalance between peripheral and logic circuitry resources. In this context, we drastically reduced the size of elementary computational blocks (chapter 5), and defined the obtained logic as ultra-fine grain logic. Two principal technological paradigms have been explored: the use of enhanced-functionality devices (i.e. devices with innovative functionalities with the same dimensions) and the use of density-increased devices (i.e. more devices in the same space). Regarding the technological obstacles for each technology, we proposed credible technological assumptions for all the expected processes.

For the functionality-enhanced device, we propose to use DG-CNFET technology as case study. A DG-CNFET can be configured to one of three different states (*n*, *p* and *off*) by changing its back-gate voltage. This device enabled the design, in prior work, of a compact dynamic logic cell able to perform 14 Boolean functions with only 7 transistors. We showed that such a cell improves area by 3.1x with respect to the equivalent CMOS counterpart due to the increase in functionality, and average power consumption by 2x due to the good electrical properties of carbon electronics. It is worth noticing that this evaluation was carried out with comparisons to a prospective CMOS technology (objectively extrapolated with a methodology based on the ITRS approach and figures). We then provide a generalization of the use of control polarity-gate voltages on dynamic logic cells. In particular, the back-gate is used to select the polarity of the transistor and to control the evaluation of the different stages. We demonstrated a potential gain up to 50% in area compared to standard dynamic MOS logic.

Two density-increased devices process were subsequently explored. Emerging technologies have opened the way towards the use of active devices in a high-density crossbar structure. We proposed the use of a sublithographic crossbar of nanowires to realize an ultra-fine grain logic cell. While such an organization has already been envisaged to build complete and

complex systems, we used it to implement small circuits only. The scheme was simplified, in order to manage the connection requirements between the micro and nano scales. In particular, we added internal inverters directly within the structure, instead of propagating inverted signals through the chip. We demonstrated that such a scheme would lead to an improvement of 4.1x in area and of 4.6x in intrinsic delay. This is due to the very small dimensions achieved by the device integration well below lithographic dimensions. Nevertheless, we should note that such a technology is highly controversial regarding the fabrication process. In fact, sublithographic dimensions often require complex alignments of bottom-up fabricated devices. Thus, we proposed a crossbar process flow derived from an industrial and lithographic FDSOI technology, with higher feasibility than the previous approach. We devised the technology guidelines and the layout structure for the logic cells using the crossbar arrangement. Considering all the parasitics, we then performed technological optimization in order to find the best process conditions. We finally demonstrated that the technology yields an area improvement of 6x and a power consumption improvement of 1.48x compared to the CMOS counterpart.

### **7.2.5 Contributions to Architectures**

As previously presented, the architecture could be improved in two ways: an improvement of routing resources (chapter 4) and an improvement of the logic blocks that lead to a new architectural organization (chapter 6).

Improvement of the architecture in terms of memory and routing resources has been done using the various technologies introduced. For each technology explored, a different FPGA organization has been proposed. Thanks to the traditional FPGA benchmarking flow and the tuning of architectural parameters, we showed that it is possible to improve the area, compared to the conventional FPGA scheme, by 46%. This value is obtained for the vertical NWFET technology. This technology is the best in terms of routing circuit area improvement. This obviously impacts the architecture accordingly. Others technologies led to a gain of 21% in area for monolithic 3-D and 13% for phase-change memories. Regarding critical path delay, a maximum gain (up to 44%) was observed for phase-change memories, which is obviously due to the low on-resistance of the technology (i.e. good switch performances). Vertical NWFET technology also led to good results with a gain of 42%. Monolithic 3-D improved the delay by 22%, due to the good electrical properties of the FDSOI technology.

In terms of architectural organization, we described a new scheme suited to computation at ultra-fine grain (chapter 6). In particular, we explored the use of compact logic cells (chapter 5). While it is not possible to use a standard FPGA organization due to a large interconnection overload, we proposed the use of matrices of logic cells. These matrices are expected to perform combinational computation. To prevent interconnection complexity, we further proposed the use of a layered interconnect, based on fixed and incomplete topologies. These arrangements were named MClusters, and were introduced in the standard FPGA scheme instead of LUTs. Using the specific benchmarking tools that we developed, we studied the topology impact on the performance of MClusters. We found that the best topology is the Modified Omega, thanks to its properties of interconnect shuffling. We then performed an evaluation of the complete FPGA architecture, and we demonstrated that the best MCluster sizing was obtained for matrices of 3 by 3 logic cells. Such a scheme yielded an improvement of 46% in the total area for the FPGA and a global improvement of the critical path delay repartition.

## **7.3 Future Work**

### **7.3.1 Design Generalization and Technology Heterogeneity**

In this thesis, several emerging technologies were studied within a reconfigurable architectural template. Reconfigurable architectures represent a good trade-off between flexibility and performance which makes them promising for future computing systems.

Within this framework, a generic tool flow has been proposed in order to evaluate the impact of an emerging technology.

In the present work, we used emerging technologies in addition to the CMOS technology, and only a single emerging technology was envisaged for each proposal. In fact, we expect that only a single technology will first be mixed with CMOS in a hybrid approach. This makes sense regarding the cost issues that must be overcome for full process development. Nevertheless, it is detrimental not to mix the technologies that can improve the specific functions of a complete circuit. For example, an ultra-fine grain FPGA architecture with resistive memory-based routing resources would be of clear interest. Indeed, we demonstrated that an ultra-fine grain organization yields an area improvement of about 50%, while the resistive memories improve the critical path delay by also almost 50%. With such alliance combination, it is possible to drastically improve the performance of the structure regarding all the different metrics. However, careful precautions must be taken regarding the heterogeneity of the structure. In fact, low-cost technological heterogeneity can be achieved if and only if the various technology steps required for each approach do not overlap. In the above example, the fabrication of ultra-fine cells is a front-end process, while resistive memory integration is done at the back-end. This means that the solution will be efficient regarding the cost of production.

In terms of circuits and tools, it could be of interest to broaden the study to standard ASIC systems. The micro-electronics industry is still attached to specific standard cell-based circuit design and FPGAs are usually considered only for prototyping or low-volume applications. Thus, providing a methodology and a tool suite that would enable the fast evaluation of an emerging technology within an ASIC flow would be relevant, and would share a good deal of functionality with the proposed reconfigurable template-based tool flow. Furthermore, several industrial ASIC tools can be derived, as in [209] where a 2-D placement methodology was derived to manage monolithic 3-D technology. This implies that the development of such a tool should be aware of with several available background contexts.

### **7.3.2 Fault Tolerance**

For the purposes of simplification in this thesis, we have not directly addressed the question of reliability in future devices. Nonetheless, the questions regarding robustness must be addressed in priority for the viability of future computing systems.

Dealing with emerging technologies at the nano scale is inherently synonymous with examining reliability. It is obvious that while the CMOS technology dimensions decrease, system robustness decreases. In addition, the co-integration of advanced CMOS and non-mature technologies will necessarily lead to unreliable systems. The lack of maturity of fabrication processes means that the device properties are heavily impacted by the technological variability. This implies that architectural design must be intimately connected to knowledge of envisaged manufacturing techniques. As an illustration, we consider designs with carbon nanotubes. The process of manufacturing circuits requires the growth of semiconducting CNTs, which are grown on or transferred to a substrate. Two major challenges remain standing before the large-scale integration of CNFETs can become feasible. First, CNT synthesis techniques are not able to grow exclusively semi-conducting CNTs, which results in the creation of short-circuits in the transistors. Second, CNT manufacturing cannot accurately control the position of CNTs, which results in incorrect logic behavior of fabricated logic structures.

The efficient use of the projected tens of billions of elementary, unreliable, nanometric devices will in all probability lead to the emergence of reconfigurable platforms as the principal computing fabric. The reconfigurable approach demonstrates broad defect tolerance capabilities. Conventional robustness techniques could be used, such as defect avoidance [53], error correction coding [54], or redundancy. However, a robust design approach is preferred to increase the reliability of a circuit early in the design. In [55], a special layout technique was proposed in order to improve the alignment accuracy of the CNFET technology. In [56],

a method was proposed to ensure good immunity to the conduction of metallic nano-tubes. To summarize, we should remark that the reliable computation of future systems must be tackled by considerations at each architectural level from the layout level to complex architecture level.

Robust techniques are inevitably costly in terms of performance and circuit area. Designers must keep in mind the broad nature of effect and quantify the impact of the circuit robustness with respect to the original element. For example, we can consider triple modular redundancy that copies the same elements three times and adds a majority voting element to eject a potential faulty decision. Such a solution is more than 3x bigger than the original circuit and obviously impacts the global circuit performance. Hence, the evaluation of robustness should be integrated into the evaluation tool. Based on the presented benchmarking tool, it is possible to integrate the defect tolerance as a metric. Novel metrics can thus be added such as a modified Mean-Time Before Failure or expected performance impact. This is a novel opportunity for the exploration of the design space, since the evaluation of the strengths and weaknesses of the robustness techniques from the architectural point of view will allow a fair comparison of the expected gain of novel over existing technologies.

# References

---

- [1] J. Bardeen and W. H. Brattain, “Three-electrode circuit element utilizing semiconductor materials,” US Patent No. 2524035, 1948.
- [2] D. Kahng, “Electric field controlled semiconductor device,” US Patent No. 2524035, 1960.
- [3] G. E. Moore, “Cramming more components onto integrated circuits,” *Electronics*, vol.38, no. 8, 19 April 1965.
- [4] Transistor Count and Moore's Law, Wikipedia, 2011.
- [5] H. Iwai, “Roadmap for 22nm and beyond (Invited Paper)”, *Microelectronic Engineering*, vol. 86(7-9), pp. 1520-1528, 2009.
- [6] Executive Summary, Updated Edition, International Technology Roadmap for Semiconductors, 2010.  
<http://www.itrs.net/Links/2010ITRS/Home2010.htm>
- [7] Emerging Research Devices and Materials Chapters, Updated Editions, International Technology Roadmap for Semiconductors, 2010.  
<http://www.itrs.net/Links/2010ITRS/Home2010.htm>
- [8] T. Hoffmann, G. Doorribos, I. Ferain, N. Collaert, P. Zimmerman, M. Goodwin, R. Rooyackers, A. Kottantharayil, Y. Yim, A. Dixit, K. De Meyer, M. Jurczak, and S. Biesemans, “GIDL (gate-induced drain leakage) and parasitic schottky barrier leakage elimination in aggressively scaled HfO<sub>2</sub>/TiN FinFET devices,” pp. 725–728, Dec. 2005.
- [9] T. Hori, “Drain-structure design for reduced band-to-band and band-to-defect tunneling leakage,” pp. 69–70, June 1990.
- [10] R. H. Dennard, F. H. Gaensslen, H.-N. Yu, V. L. Rideout, E. Bassous, and A. R. Leblanc, “Design of ion-implanted MOSFET’s with very small physical dimensions,” *IEEE Journal of Solid-State Circuits*, vol. 9, no. 5, pp. 256–268, 1974.
- [11] System Drivers Chapter, Updated Edition, International Technology Roadmap for Semiconductors, 2010.  
<http://www.itrs.net/Links/2010ITRS/Home2010.htm>
- [12] Design Chapter, Updated Edition, International Technology Roadmap for Semiconductors, 2010.  
<http://www.itrs.net/Links/2010ITRS/Home2010.htm>
- [13] D. Marple and L. Cooke, “An MPGA Compatible FPGA Architecture,” *IEEE Custom Integrated Circuits Conference*, pp.4.2.1-4.2.4, 3-6 May 1992.

- [14] P. Dillien, "Adaptive hardware becomes a reality using electrically reconfigurable arrays (ERAs)," IEE Colloquium on User-Configurable Logic - Technology and Applications, pp.2/1-2/10, 1 Mar 1991.
- [15] A. El Gamal, J. Greene, J. Reyneri, E. Rogoyski, K.A. El-Ayat, A. Mohsen, "An architecture for electrically configurable gate arrays," IEEE Journal of Solid-State Circuits, vol.24, no.2, pp.394-398, Apr 1989.
- [16] B. Small, "The flexibility of the QuickLogic FPGA architecture," WESCON, pp.688-691, 27-29 Sep 1994.
- [17] V. Betz, J Rose and A. Marquart, "Architecture and CAD for Deep-Submicron FPGAs", Kluwer Academic Publishers, New York, 1999.
- [18] J. Rose, S. Brown, "Flexibility of interconnection structures for field-programmable gate arrays," IEEE Journal of Solid-State Circuits, vol.26, no.3, pp.277-282, Mar 1991.
- [19] E. Ahmed, "The effect of logic block granularity on deep-submicron FPGA performance and density", Master thesis, University of Toronto, 2001.
- [20] Xilinx Virtex-6 FPGA Family Overview, 24 March 2011.  
[http://www.xilinx.com/support/documentation/data\\_sheets/ds150.pdf](http://www.xilinx.com/support/documentation/data_sheets/ds150.pdf)
- [21] M. Hutton, V. Chan, P. Kazarian, V. Maruri, T. Ngai, J. Park, R. Patel, B. Pedersen, J. Schleicher and S. Shumarayev, "Interconnect enhancements for a high-speed PLD architecture", ACM/SIGDA 10th International Symposium on. FPGA, pp.3, 2002.
- [22] D. Lewis, E. Ahmed, G. Baeckler, V. Betz, M. Bourgeault, D. Cashman, D. Galloway, M. Hutton, C. Lane, A. Lee, P. Leventis, S. Marquardt, C. McClintock, K. Padalia, B. Pedersen, G. Powell, B. Ratchev, S. Reddy, J. Schleicher, K. Stevens, R. Yuan, R. Cliff and J. Rose, "The Stratix II logic and routing architecture," ACM/SIGDA 13th International Symposium on. FPGA, pp.14, 2005
- [23] E. Kusse and J. Rabaey, "Low-energy embedded FPGA structures," International Symposium on Low Power Electronics and Design, pp.155, 1998.
- [24] L. Shang, A. S. Kaviani and K. Bathala, "Dynamic power consumption in Virtex-II FPGA family," ACM/SIGDA 10th International Symposium on. FPGA, pp.157, 2002.
- [25] V. Degalahal and T. Tuan, "Methodology for high level estimation of FPGA power consumption," Design Automation Conference, pp.657, 2005.
- [26] I. Kuon and J. Rose, "Measuring the gap between FPGAs and ASICs," ACM/SIGDA 14th International Symposium on. FPGA, pp.21, 2006.
- [27] J. Birkner, "Reduce random-logic complexity", Electronic Design, vol.26, no.17, pp.98-105, January 1978.
- [28] Altera Quartus-II: Altera FPGAs Design IDE.  
<http://www.altera.com/products/software/sfw-index.jsp>
- [29] Xilinx ISE: Xilinx FPGAs Design IDE.  
<http://www.xilinx.com/products/design-tools/ise-design-suite/index.htm>
- [30] W. Carter, K. Duong, R. H. Freeman, H. Hsieh, J. Y. Ja, J. E. Mahoney, L. T. Ngo, and S. L. Sze, "A user programmable reconfigurable gate array," Proceedings of the Custom Integrated Circuits Conference, pp.233-235, May 1986.
- [31] Xilinx Spartan-3A FPGA Product Overview, 4 April 2011.  
[http://www.xilinx.com/support/documentation/data\\_sheets/ds557.pdf](http://www.xilinx.com/support/documentation/data_sheets/ds557.pdf)



- [32] J. McCollum, H.-S. Chen, and F. Hawley, “Non-volatile programmable memory cell for programmable logic array,” US Patent No. 0064484, 2007.
- [33] J. McCollum, G. Bakker, and J. Greene, “Non-volatile look-up table for an FPGA,” US Patent No. 0007293, 2008.
- [34] J. Lipp, D. Freeman, U. Broze, M. Caywood, and G. Nolan, “A general purpose, non-volatile reprogrammable switch,” WO Patent No. 01499, 1996.
- [35] K. J. Han, N. Chan, S. Kim B. Leung. V. Hecht, and B. Cronquist, “A Novel Flash-based FPGA Technology with Deep Trench Isolation,” IEEE Non-Volatile Semiconductor Memory Workshop, pp.32-33, 26-30 Aug. 2007.
- [36] H. S. Stone, “Parallel Processing with the Perfect Shuffle,” IEEE Transactions on Computers, vol.C-20, no.2, pp. 153- 161, Feb. 1971.
- [37] N. Bruchon, L. Torres, G. Sassatelli, G. Cambon, “New nonvolatile FPGA concept using magnetic tunneling junction,” IEEE Computer Society Annual Symposium on Emerging VLSI Technologies and Architectures, pp.6, 2-3 March 2006.
- [38] M. Lin, A. El Gamal, Y.-C. Lu and S. Wong, “Performance Benefits of Monolithically Stacked 3-D FPGA,” IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, vol.26, no.2, pp.216-229, Feb. 2007.
- [39] J.J. Hopfield, “Neural networks and physical systems with emergent collective computational abilities,” Proceedings of the National Academy of Sciences of the United States of America, vol.79, pp.2254-2258; 1982.
- [40] R.C. Eberhardt, and R.W. Dobbins, “Neural Networks PC Tools – A practical guide,” Academic Press Inc., October 1990.
- [41] J. Hoekstra, and E. Rouw, “Modeling of dendritic computation: The single dendrite,” The American Institute of Physics, vol.517, pp.308-322; 2000.
- [42] C. Mead, “Analog VLSI and Neural Systems,” Addison Wesley, 1989.
- [43] D.B. Strukov, K.K. Likharev, “Reconfigurable Hybrid CMOS/Nanodevice Circuits for Image Processing,” IEEE Transactions on Nanotechnology, vol.6, no.6, pp.696-710, Nov. 2007.
- [44] W.Wu, G.-Y. Jung, D. L. Olynick, J. Straznicky, Z. Li, X. Li, D. A. A. Ohlberg, Y. Chen, S.-Y. Wang, J. A. Liddle, W. M. Tong, and R. S. Williams, “One-kilobit cross-bar molecular memory circuits at 30-nm half-pitch fabricated by nanoimprint lithography,” Applied Physics A: Materials Science and Processing, vol. 80, no. 6, pp. 1173–1178, 2005.
- [45] Y. Luo, C. P. Collier, J. O. Jeppesen, K. A. Nielsen, E. DeIonno, G. Ho, J. Perkins, H.-R. Tseng, T. Yamamoto, J. F. Stoddart, and J. R. Heath, “Two-dimensional molecular electronics circuits,” Journal of Chemical Physics and Physical Chemistry, vol. 3, pp. 519–525, 2002.
- [46] J. E. Green, J. Wook Choi, A. Boukai, Y. Bunimovich, E. Johnston-Halperin, E. Deionno, Y. Luo, B. A. Sheriff, K. Xu, Y. Shik Shin, H.-R. Tseng, J. F. Stoddart, and J. R. Heath, “A 160-kilobit moleculelectronic memory patterned at 1011 bits per square centimetre,” Nature, vol. 445, pp. 414–417, 2007.
- [47] S. C. Goldstein, and M. Budiu, “NanoFabrics: spatial computing using molecular electronics”, 28th Annual International Symposium on Computer Architecture, pp.178-189, 2001.
- [48] S. Goldstein and D. Rosewater, “Digital logic using molecular electronics,” IEEE International Solid-State Circuits Conference, vol. 1, pp. 204–459, 2002.

- [49] A. DeHon, P. Lincoln, and J. Savage, "Stochastic Assembly of Sublithographic Nanoscale Interfaces," *IEEE Transactions on Nanotechnology*, vol. 2, no. 3, pp. 165–174, 2003.
- [50] H. Yan, H. S. Choe, S.W. Nam, Y. Hu, S. Das, J. F. Klemic, J. C. Ellenbogen, and C. M. Lieber, "Programmable nanowire circuits for nanoprocessors," *Nature Letter*, vol.470, pp.240-244, 10 February 2011.
- [51] M.H. Ben Jamaa, D. Atienza, Y. Leblebici, and G. De Micheli, "Programmable logic circuits based on ambipolar CNFET," 45th ACM/IEEE Design Automation Conference, pp.339-340, 8-13 June 2008.
- [52] M.H. Ben Jamaa, K. Mohanram, and G. De Micheli, "Novel library of logic gates with ambipolar CNTFETs: Opportunities for multi-level logic synthesis," *Design, Automation & Test in Europe Conference & Exhibition*, pp.622-627, 20-24 April 2009.
- [53] J.R. Heath, P. J. Kuekes, G. S. Snider and R. S. Williams, "A Defect-Tolerant Computer Architecture: Opportunities for Nanotechnology," *Science*, vol.280,no.5370, pp.1716-1721, 12 June 1998.
- [54] W. Huffman, V. Pless, "Fundamentals of error-correcting codes," Cambridge University Press, 2003.
- [55] N. Patil, J. Deng, A. Lin, H.S.-P. Wong and S. Mitra, "Design Methods for Misaligned and Mis-positioned Carbon-Nanotube-Immune Circuits," *IEEE Transactions on Computer-Aided Design*, 2008.
- [56] N. Patil, A. Lin, J. Zhang, H. Wei, K. Anderson, H.-S.P. Wong and S. Mitra, "Scalable Carbon Nanotube Computational and Storage Circuits Immune to Metallic and Mis-positioned Carbon Nanotubes," *IEEE Transactions on Nanotechnology*, 2010
- [57] H.J. Hovel and J.J. Urgell, "Switching and Memory Characteristics of ZnSe - Ge Heterojunctions," *Journal of Applied Physics*, vol. 42, pp.5076, 1971.
- [58] I.G. Baek, M.S. Lee, S. Seo, M.J. Lee, D.H. Seo, D.-S. Suh, J.C. Park, S.O. Park, H.S. Kim, I.K. Yoo, U.-I. Chung, and J.T. Moon, "Highly scalable nonvolatile resistive memory using simple binary oxide driven by asymmetric unipolar voltage pulses," *International Electron Devices Meeting*, pp. 587-590, 13-15 December 2004.
- [59] K. Tsunoda, K. Kinoshita, H. Noshiro, Y. Yamazaki, T. Iizuka, Y. Ito, A. Takahashi, A. Okano, Y. Sato, T. Fukano, M. Aoki, and Y. Sugiyama, "Low Power and High Speed Switching of Ti-doped NiO ReRAM under the Unipolar Voltage Source of less than 3 V," *International Electron Devices Meeting*, pp.767-770, 10-12 December 2007.
- [60] C. Nauenheim, C. Kugeler, S. Trelenkamp, A. Rudiger, and R. Waser, "Phenomenological considerations of resistively switching TiO<sub>2</sub> in nano crossbar arrays," 10<sup>th</sup> International Conference on ULIS, pp.135-138, 18-20 March 2009.
- [61] H.Y. Lee, Y.S. Chen, P. S. Chen, T. Y. Wu, F. Chen, C.C. Wang, P.J. Tzeng, M.-J. Tsai, and C. Lien, "Low-Power and Nanosecond Switching in Robust Hafnium Oxide Resistive Memory With a Thin Ti Cap," *IEEE Electron Device Letters*, vol.31, no.1, pp. 44-46, January 2010.
- [62] W.C. Chien, Y.C. Chen, K.P. Chang, E.K. Lai, Y.D. Yao, P. Lin, J. Gong, S.C. Tsai, S.H. Hsieh, C.F. Chen, K.Y. Hsieh, R. Liu, and C.-Y. Lu, "Multi-Level Operation of Fully CMOS Compatible WOX Resistive Random Access Memory (RRAM)," *International Memory Workshop*, pp.1-2, 10-14 May 2009.
- [63] P. Zhou, H.J. Wan, Y.L. Song, M. Yin, H.B. Lu, Y.Y. Lin, S. Song, R. Huang, J.G. Wu, and M.H. Chi, "A Systematic Investigation of TiN/Cu<sub>x</sub>O/Cu RRAM with Long

- Retention and Excellent Thermal Stability,” International Memory Workshop, pp.1-2, 10-14 May 2009.
- [64] C. A. Moritz, T. Wang, “Latching on the Wire and Pipelining in Nanoscale Designs,” 3rd Workshop on Non-Silicon Computation (NSC-3), June 2004.
- [65] T. Wang, P. Narayanan, and C. A. Moritz, “Combining 2-level Logic Families in Grid-based Nanoscale Fabrics,” IEEE/ACM International Symposium on Nanoscale Architectures (NANOARCH), October 2007.
- [66] P. Vijayakumar, P. Narayanan, I. Koren, C. M. Krishna, and C. A. Moritz, “Impact of Nanomanufacturing Flow on Systematic Yield Losses in Nanoscale Fabrics,” IEEE/ACM International Symposium on Nanoscale Architectures (NANOARCH), June 2011.
- [67] Z. Wei, Y. Kanzawa, K. Arita, Y. Katoh, K. Kawai, S. Muraoka, S. Mitani, S. Fujii, K.; Katayama, M. Iijima, T. Mikawa, T. Ninomiya, R. Miyanaga, Y. Kawashima, K. Tsuji, A. Himeno, T. Okada, R. Azuma, K. Shimakawa, H. Sugaya, T. Takagi, R. Yasuhara, K. Horiba, H. Kumigashira, M. Oshima, “Highly reliable TaOx ReRAM and direct evidence of redox reaction mechanism,” International Electron Devices Meeting, pp.1-4, 15-17 Dec. 2008.
- [68] A. Pirovano, F. Pellizzer, I. Tortorelli, A. Rigano, R. Harrigan, M. Magistretti, P. Petruzza, E. Varesi, A. Redaelli, D. Erbetta, T. Marangon, F. Bedeschi, R. Fackenthal, G. Atwood, and R. Bez, “Phase-change memory technology with self-aligned  $\mu$ Trench cell architecture for 90 nm node and beyond,” Solid-State Electronics, vol. 52, no. 9, pp.1467-1472, September 2008.
- [69] K. SangBum, Z. Yuan, J.P. McVittie, H. Jagannathan, Y. Nishi, and H.-S.P. Wong, “Integrating Phase-Change Memory Cell With Ge Nanowire Diode for Crosspoint Memory—Experimental Demonstration and Analysis,” IEEE Transactions on Electron Devices, vol.55, no.9, pp.2307-2313, September 2008.
- [70] S. Lai, “Current status of the phase change memory and its future”, IEDM Technical Digest, pp.225-228, Dec. 2003.
- [71] S. Raoux, G. W. Burr, M. J. Breitwisch, C. T. Rettner, Y.-C. Chen, R. M. Shelby, M. Salinga, D. Krebs, S.-H. Chen, H.-L. Lung, C. H. Lam, “Phase-change random access memory: A scalable technology”, IBM Journal of Research and Development, vol. 52, no.4-5, pp.465-479, 2008.
- [72] G. Bruns, P. Merkelbach, C. Schlockermann, M. Salinga, M. Wuttig, T. D. Happ, J. B. Philipp and M. Kund, “Nanosecond switching in GeTe phase change memory cells”, Applied Physics Letters, vol.95, no.4, 2009.
- [73] G.Servalli, “A 45nm Generation Phase Change Memory Technology”, IEDM Technical Digest, pp. 113-116, 2009.
- [74] J.H. Oh, J. H. Park, Y. S. Lim, H. S. Lim, Y. T. Oh, J. S. Kim, J. M. Shin, Y. J. Song, K. C. Ryoo, D. W. Lim, S. S. Park, J. I. Kim, J. H. Kim, J. Yu, F. Yeung, C. W. Jeong, J. H. Kong, D. H. Kang, G. H. Koh, G. T. Jeong, H. S. Jeong and K. Kinam, “Full Integration of Highly Manufacturable 512Mb PRAM based on 90nm Technology”, IEDM Tech.Dig., p.49-52, 2006.
- [75] A.L. Lacaita and D.J. Wouters, “Phase-change memories”, physica status solidi (a), vol. 205, no.10, pp.2281-2297, October 2008.
- [76] J.H. Kyung, N. Chan, K. Sungraen, L. Ben, V. Hecht and B. Cronquist,, “A novel flash-based FPGA technology with deep trench isolation”, IEEE Non-Volatile Semiconductor Memory Workshop, pp.32-33, 2007.

- [77] D. Ielmini and M. Boniardi, "Common signature of many-body thermal excitation in structural relaxation and crystallization of chalcogenide glasses", *Applied Physics Letters*, vol.94, no.09, pp.091906, 2009.
- [78] D. Ielmini and Y. Zhang, "Analytical model for subthreshold conduction and threshold switching in chalcogenide-based memory devices", *Journal of Applied Physics*, vol.102, no.5, pp.054517, 2007.
- [79] G. Betti Beneventi, A. Calderoni, P. Fantini, L. Larcher and P. Pavan, "Analytical model for low-frequency noise in amorphous chalcogenide-based phase-change memory devices", *Journal of Applied Physics*, vol.106, no.5, pp.054506, 2009.
- [80] T. Morikawa, K. Kurotsuchi, M. Kinoshita, N. Matsuzaki, Y. Matsui, Y. Fuiisaki, S. Hanzawa, A. Kotabe, M. Terao, H. Moriya, T. Iwasaki, M. Matsuoka, F. Nitta, M. Moniwa, T. Koga, N. Takaura, "Doped In-Ge-Te Phase Change Memory Featuring Stable Operation and Good Data Retention," *IEEE International Electron Devices Meeting*, pp.307-310, 10-12 Dec. 2007.
- [81] B. Gleixner, F. Pellizzer, R. Bez, "Reliability characterization of Phase Change Memory," *10th Annual Non-Volatile Memory Technology Symposium*, pp.7-11, 25-28 Oct. 2009.
- [82] A. Fantini, L. Perniola, M. Armand, J.-F. Nodin, V. Sousa, A. Persico, J. Cluzel, C. Jahan, S. Maitrejean, S. Lhostis, A. Roule, C. Dressler, G. Reibold, B. De Salvo, P. Mazoyer, D. Bensahel and F. Boulanger, "Comparative Assessment of GST and GeTe Materials for Application to Embedded Phase-Change Memory Devices," *IEEE International Memory Workshop*, pp.1-2, 10-14 May 2009.
- [83] L. Perniola, V. Sousa, A. Fantini, E. Arbaoui, A. Bastard, M. Armand, A. Fargeix, C. Jahan, J.-F. Nodin, A. Persico, D. Blachier, A. Toffoli, S. Loubriat, E. Gourvest, G. Betti Beneventi, H. Feldis, S. Maitrejean, S. Lhostis, A. Roule, O. Cueto, G. Reibold, L. Poupinet, T. Billon, B. De Salvo, D. Bensahel, P. Mazoyer, R. Annunziata, P. Zuliani and F. Boulanger, "Electrical Behavior of Phase-Change Memory Cells Based on GeTe," *IEEE Electron Device Letters*, vol.31, no.5, pp.488-490, May 2010.
- [84] T. H. Jeong, M. R. Kim, H. Seo, J. W. Park and C. Yeon, "Crystal Structure and Microstructure of Nitrogen-Doped  $\text{Ge}_2\text{Sb}_2\text{Te}_5$  Thin Film", *Japan Journal of Applied Physics*, vol.39, pp.2775-2779, 2009.
- [85] Y. Lai, B. Qiao, J. Feng, Y. Ling, L. Lai, Y. Lin, T. Tang, B. Cai and B. Chen, "Nitrogen-doped  $\text{Ge}_2\text{Sb}_2\text{Te}_5$  films for nonvolatile memory," *Journal of Electronic Materials*, vol.34, no.2, pp.176-181, 2005.
- [86] R. Waser, "Electrochemical and thermochemical memories," *IEEE International Electron Devices Meeting*, pp.1-4, 15-17 Dec. 2008.
- [87] M. Kund, G. Beitel, C.-U. Pinnow, T. Rohr, J. Schumann, R. Symanczyk, K.-D. Ufert, and G. Muller, "Conductive bridging RAM (CBRAM): an emerging non-volatile memory technology scalable to sub 20nm," *IEEE International Electron Devices Meeting*, pp.754-757, 5 Dec. 2005.
- [88] M.-J. Lee, Y. Park, B.-S. Kang, S.-E. Ahn, C. Lee, K. Kim, W. Xianyu, G. Stefanovich, J.-H. Lee, S.-J. Chung, Y.-H. Kim, C.-S. Lee, J.-B. Park, I.-K. Yoo, "2-stack 1D-1R Cross-point Structure with Oxide Diodes as Switch Elements for High Density Resistance RAM Applications," *IEEE International Electron Devices Meeting*, pp.771-774, 10-12 Dec. 2007.
- [89] S. Brown, R. Francis, J. Rose, and Z. Vranesic; "Field-Programmable Gate Arrays with Embedded Memories," *Kluwer Academic Publisher*, 1992.

- [90] M. Reyboz, O. Rozeau, L. Perniola and G. Betti Beneventi, "Compact Modeling of a PCRAM cell", MOS AK workshop, April 2010, Italy.
- [91] A. Pirovano, A. L. Lacaita, F. Pellizzer, S. A. Kostylev, A. Benvenuti and R. Bez, "Low-field amorphous state resistance and threshold voltage drift in chalcogenide materials", IEEE Transactions on Electron Devices, vol.51, no.5, pp.714-719, May 2004.
- [92] J. Z. Sun and D. C. Ralph, "Magnetoresistance and spin-transfer torque in magnetic tunnel junctions", Journal of magnetism and magnetic materials, vol. 320, no.7, pp. 1227-1237, 2008.
- [93] Y. Guillemenet, L. Torres and G. Sassatelli, "Non-volatile run-time field-programmable gate arrays structures using thermally assisted switching magnetic random access memories," IET Computers & Digital Techniques, vol.4, no.3, pp.211-226, May 2010.
- [94] Y. Akasaka and T. Nishimura, "Concept and basic technologies for 3-D IC structure," International Electron Devices Meeting, vol.32, pp. 488- 491, 1986.
- [95] A. Rahman, J. Trezza, B. New and S. Trimberger, "Die Stacking Technology for Terabit Chip-to-Chip Communications," IEEE Custom Integrated Circuits Conference, CICC '06., pp.587-590, 10-13 Sept. 2006.
- [96] D. Henry, S. Cheramy, J. Charbonnier, P. Chausse, M. Neyret, C. Brunet-Manquat, S. Verrun, N. Sillon, L. Bonnot, X. Gagnard and E. Saugier, "3D integration technology for set-top box application," IEEE International Conference on 3D System Integration, 3DIC 2009, pp.1-7, 28-30 Sept. 2009.
- [97] P. Batude, M. Vinet, A. Pouydebasque, L. Clavelier, C. LeRoyer, C. Tabone, B. Previtali, L. Sanchez, L. Baud, A. Roman, V. Carron, F. Nemouchi, S. Pocas, C. Comboroure, V. Mazzocchi, H. Grampeix, F. Aussenac and S. Deleonibus, "Enabling 3D monolithic Integration," ECS journal, vol.6, pp.47, 2008.
- [98] Y.-H. Son, J.-W. Lee; P. Kang, M.-G. Kang, J. B. Kim, S. H. Lee, Y.-P. Kim, I. S. Jung, B. C. Lee; S. Y. Choi; U.-I. Chung, J. T. Moon, B.-I. Ryu, "Laser-induced Epitaxial Growth (LEG) Technology for High Density 3-D Stacked Memory with High Productivity," IEEE Symposium on VLSI Technology, pp.80-81, 12-14 June 2007.
- [99] S. E. Steen, D. LaTulipe, A. W. Topol, D. J. Frank, K. Belote, D. Posillico, "Overlay as the key to drive wafer scale 3D integration", Microelectronic Engineering, vol.84, no.5-8, pp.1412-1415, May-August 2007.
- [100] P. Batude, M. Vinet, A. Pouydebasque, C. Le Royer, B. Previtali, C. Tabone, J.-M. Hartmann, L. Sanchez, L. Baud, V. Carron, A. Toffoli, F. Allain, V. Mazzocchi, D. Lafond, N. Bouzaida, O. Thomas, O. Cueto, A. Amaral, S. Deleonibus and O. Faynot, "Advances in 3D CMOS Sequential Integration," IEEE International device Meeting, 2009.
- [101] P. Batude, M. Vinet, A. Pouydebasque, C. Le Royer, B. Previtali, C. Tabone, L. Clavelier, S. Michaud, A. Valentian, O. Thomas, O. Rozeau, P. Coudrain, C. Leyris, K. Romanjek, X. Garros, L. Sanchez, L. Baud, A. Roman, V. Carron, H. Grampeix, E. Augendre, A. Toffoli, F. Allain, P. Grosgeorges, V. Mazzochi, L. Tosti, F. Andrieu, J.-M. Hartmann, D. Lafond, S. Deleonibus and O. Faynot, "GeOI and SOI 3D monolithic cell integrations for high density applications," 2009 Symposium on VLSI Technology, pp.166-167, 16-18 June 2009.
- [102] S.-M. Jung, Y. Rah, T. Ha, H. Park, C. Chang, S. Lee, J. Yun, W. Cho, H. Lim, J. Park, J. Jeong, B. Son, J. Jang, B. Choi, H. Cho, K. Kim, "Highly cost effective and high performance 65nm S<sub>3</sub> (stacked single-crystal Si) SRAM technology with 25F<sub>2</sub>,

- 0.16 $\mu\text{m}^2$  cell and doubly stacked SSTFT cell transistors for ultra high density and high speed applications”, Digest of Technical Papers of Symposium on VLSI Technology, pp.220, 2005.
- [103] S.-M. Jung, H. Lim, C. Yeo, K. Kwak, B. Son, H. Park, J. Na, J.-J. Shim, C. Hon, K. Kim, “High Speed and Highly Cost effective 72M bit density  $\text{S}_3$  SRAM Technology with Doubly Stacked Si Layers, Peripheral only CoSix layers and Tungsten Shunt W/L Scheme for Standalone and Embedded Memory”, Digest of Technical Papers of Symposium on VLSI Technology, pp.82, 2007.
- [104] J. Feng, Y. Liu, P.-B. Griffin, J.-D. Plummer, “Integration of Germanium-on-Insulator and Silicon MOSFETs on a Silicon Substrate”; IEEE Electron Device Letters, vol.27, no.11, pp. 911, 2006.
- [105] D.-S. Yu, A. Chin, C.-C. Liao, C.-F. Lee, C.-F. Cheng, M.-F. Li, W.-J. Yoo, S.-P. McAlister, “Three-dimensional metal gate-high- $\kappa$ -GOI CMOSFETs on 1-poly-6-metal 0.18- $\mu\text{m}$  Si devices”; IEEE Electron Device Letters, vol.26, no.2, pp.118, 2005.
- [106] M. Jeong, B. Doris, J. Kedzierski, K. Rim and M. Yang, “Silicon Device Scaling to Sub-10-nm Regime”, Science, vol. 306, no. 5704, pp.2057-2060, 2004.
- [107] H. Wei, N. Patil, A. Lin, H.-S.P. Wong and S. Mitra, “Monolithic Three-Dimensional Integrated Circuits using Carbon Nanotube FETs and Interconnects,” IEEE International. Electron Devices Meeting, December 2009.
- [108] T. Ernst, E. Bernard, C. Dupre, A. Hubert, S. Becu, B. Guillaumot, O. Rozeau, O. Thomas, P. Coronel, J.-M. Hartmann, C. Vizioz, N. Vulliet, O. Faynot, T. Skotnicki and S. Deleonibus, “3D multichannels and stacked nanowires technologies for new design opportunities in nanoelectronics,” IEEE International Conference on Integrated Circuit Design and Technology and Tutorial, ICICDT 2008, pp.265-268, 2-4 June 2008.
- [109] J. Hahm and C.M. Lieber, “Direct Ultrasensitive Electrical Detection of DNA and DNA Sequence Variations Using Nanowire Nanosensors,” Nano Letter, vol.4, pp.51-54, 2004.
- [110] Y. Cui, Z. Zhong, D. Wang, W. U. Wang and C.M. Lieber, “High Performance Silicon Nanowire Field Effect Transistors,” Nano Letter, vol.3, pp.149-152, 2003.
- [111] A.-L. Bavencove, G. Tourbot, E. Pugeoise, J. Garcia, P. Gilet, F. Levy, B. André, G. Feuillet, B. Gayral, B. Daudin and Le S. Dang, “GaB-based nanowires: From nanometric-scale characterization to light emitting diodes,” Physica Status Solidi (a), vol.207, no.6, pp. 1425-1427, 2010.
- [112] J. Goldberger, A. I. Hochbaum, R. Fan, and P. Yang, “Silicon Vertically Integrated Nanowire Field Effect Transistors,” Nano Letter, vol.6, no.5, pp.973-977, 2006.
- [113] V. Schmidt, H. Riel, S. Senz, S. Karg, W. Riess, U. Gösele, “Realization of a silicon nanowire vertical surround-gate field-effect transistor,” Small, vol.2, no.1, pp.85-88, 2006.
- [114] ATLAS User’s Manual, SILVACO, 2008.
- [115] G.Betti Beneventi, L. Perniola, A. Fantini, D. Blachier, A. Toffoli, E. Gourvest, S. Maitrejean, V. Sousa, C. Jahan, J. F. Nodin, A. Persico, S. Loubriat, A. Roule, S. Lhostis, H. Feldis, G. Reimbold, T. Billon, B. De Salvo, L. Larcher, P. Pavan, D. Bensahel, P. Mazoyer, R. Annunziata and F. Boulanger, “Carbon-doped GeTe Phase-Change Memory featuring remarkable RESET current reduction”, Proceedings of the European Solid-State Device Research Conference (ESSDERC), pp.313-316, 14-16 September 2010.

- [116] G. Betti Beneventi, E. Gourvest, A. Fantini, L. Perniola, V. Sousa, S. Maitrejean, J. C. Bastien, A. Bastard, A. Fargeix, B. Hyot, C. Jahan, J. F. Nodin, A. Persico, D. Blachier, A. Toffoli, S. Loubriat, A. Roule, S. Lhostis, H. Feldis, G. Reimbold, T. Billon, B. De Salvo, L. Larcher, P. Pavan, D. Bensahel, P. Mazoyer, R. Annunziata and F. Boulanger, "On Carbon doping to improve GeTe-based Phase-Change Memory data retention at high temperature", IEEE International Memory Workshop (IMW), pp.1-4, 16-19 May 2010.
- [117] S. Onkaraiah, P.-E. Gaillardon, M. Reyboz, F. Clermidy, J.-M. Portal, M. Bocquet, C. Muller, "Using OxRRAM Memories for Improving Communications of Reconfigurable FPGA Architectures," IEEE/ACM International Symposium on Nanoscale Architectures (NanoArch), 08-09 June 2011, San Diego (CA), USA.
- [118] K. Siozios, K. Sotiriadis, V. F. Pavlidis and D. Soudris, "Exploring Alternative 3D FPGA Architectures: Design Methodology and CAD Tool Support," 17th International Conference on Field Programmable Logic and Applications (FPL), pp. 652-656, 2007.
- [119] J. Rubin, R. Sundararaman, M. K. Kim and S. Tiwari, "A single lithography vertical NEMS switch," IEEE 24th International Conference on Micro Electro Mechanical Systems (MEMS), pp.95-98, 23-27 January 2011.
- [120] F. Li, Y. Lin, L. He, D. Chen and J. Cong, "Power modeling and characteristics of field programmable gate arrays," IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems, vol.24, no.11, pp. 1712- 1724, November 2005.
- [121] J. Liu, "Architectures reconfigurables à base de CNTFET double grille," Thesis report, Ecole Centrale de Lyon, 2008.
- [122] I. O'Connor, J. Liu, F. Gaffiot, F. Pregaldiny, C. Lallement, C. Maneux, J. Goguet, S. Fregonese, T. Zimmer, L. Anghel, T.-T. Dang and R. Leveugle, "CNTFET Modeling and Reconfigurable Logic-Circuit Design," IEEE Transactions on Circuits and Systems I: Regular Papers, vol.54, no.11, pp.2365-2379, Nov. 2007.
- [123] S. Iijima and T. Ichihashi, "Single-shell carbon nanotubes of 1-nm diameter," Nature, vol. 363, pp. 603-605, 1993.
- [124] K. S. Novoselov, A. K. Geim, S. V. Morozov, D. Jiang, Y. Zhang, S. V. Dubonos, I. V. Grigorieva, and A. A. Firsov, "Electric Field Effect in Atomically Thin Carbon Films," Science, vol. 306, no. 5696, pp. 666-669, 2004.
- [125] A. H. Castro Neto, F. Guinea, N. M. R. Peres, K. S. Novoselov and A. K. Geim, "The electronic properties of grapheme," Review of Modern Physics, vol. 81, no. 1, pp. 109-163, January 2009.
- [126] I. Meric, M. Han, A. F. Young, B. Ozyilmaz, P. Kim, and K. L. Shepard, "Current saturation in zero-bandgap, top-gated grapheme field-effect transistors," Nature Nanotechnology, vol. 3, no. 11, pp. 654-659, November 2008.
- [127] Y.-M. Lin, K. A. Jenkins, A. Valdes-Garcia, J. P. Small, D. B. Farmer, and P. Avouris, "Operation of grapheme transistors at gigahertz frequencies," Nano Letters, vol. 9, no. 1, pp. 422-426, January 2009.
- [128] I. Meric, N. Baklitskaya, P. Kim, and K. L. Shepard, "RF performance of top-gated, zero-bandgap grapheme field-effect transistors," IEDM Technical Digest, 2008.
- [129] V. T. Renard, M. Jublot, P. Gergaud, P. Cherns, D. Rouchon, A. Chabli and V. Jousseume, "Catalyst preparation for CMOS-compatible silicon nanowire synthesis," Nature Nanotechnology, vol.4, pp.654-657, 2009.
- [130] J. S. Moon, D. Curtis, M. Hu, D. Wong, C. McGuire, P. M. Campbell, G. Jernigan, J. L. Tedesco, B. VanMil, R. Myers-Ward, C. Eddy Jr., and D. K. Gaskill, "Epitaxial-

- graphene RF field-effect transistors on Si-face 6H-SiC substrates,” *IEEE Electron Device Letter*, vol. 30, no. 6, pp. 650-652, June 2008.
- [131] C. Berger, Z. Song, T. Li, X. Li, A. Y. Ogbazghi, R. Feng, Z. Dai, A. N. Marchenkov, E. H. Conrad, P. N. First, and W. A. de Heer, “Ultrathin Epitaxial Graphite: 2D Electron Gas Properties and a Route toward Graphene-based Nanoelectronics,” *Journal of Physichal Chemistry B*, vol. 108, no. 52, pp. 19912-19916, December 2004.
- [132] X. Li, X. Wang, L. Zhang, S. Lee, and H. Dai, “Chemically derived, ultrasmooth graphene nanoribbon semiconductors,” *Science*, vol. 319, no. 5867, pp. 1229–1232, 2008.
- [133] K.-T. Lam and G. Liang, “An ab initio study on energy gap of bilayer graphene nanoribbons with armchair edges,” *Applied Physics Letters*, vol. 92, no. 22, p. 223106, 2008.
- [134] C. L. Lu, C. P. Chang, Y. C. Huang, J. M. Lu, C. C. Hwang, and M. F. Lin, “Low-energy electronic properties of the ab-stacked few-layer graphites,” *Journal of Physics: Condensed Matter*, vol. 18, no. 26, p. 5849, 2006.
- [135] H. Min, B. Sahu, S. K. Banerjee, and A. H. MacDonald, “Ab initio theory of gate induced gaps in graphene bilayers,” *Phys. Rev. B*, vol. 75, no. 15, p. 155115, April 2007.
- [136] Y. Zhang, T.-T. Tang, C. Girit, Z. Hao, M. C. Martin, A. Zettl, M. F. Crommie, Y. R. Shen, and F. Wang, “Direct observation of a widely tunable bandgap in bilayer graphene,” *Nature*, vol. 459, no. 7248, p. 820, 2009.
- [137] M. S. Dresselhaus, G. Dresselhaus, and P. Avouris, “Carbon Nanotubes: Synthesis, Structure, Properties, and Applications.,” Berlin, Germany:Springer-Verlag, 2001.
- [138] Z. Yao, C. L. Kane, and C. Dekker, “High-field electrical transport in single-wall carbon nanotubes,” *Physical Review Letters*, vol. 84, pp. 2941-2944, 2000.
- [139] M. S. Fuhrer, B. M. Kim, T. Duerkop, and T. Brintlinger, “High-mobility nanotube transistor memory,” *Nano Letters*, vol. 2, pp. 755-759, 2002.
- [140] M. S. Fuhrer, M. Forero, A. Zettl, and P. L. McEuen, “Ballistic transport in semiconducting carbon nanotubes,” *XV International Winterschool/Euroconference on ELECTRONIC PROPERTIES OF MOLECULAR NANOSTRUCTURES*, vol. 591, pp.401-404, 2001.
- [141] R. Sordan; K. Balasubramanian, M. Burghard, and K. Kern, “Exclusive-OR gate with a single carbon nanotube,” *Applied Physics Letters*, vol.88, no.5, pp.053119-053119-3, Jan 2006.
- [142] Y.-M. Lin; J. Appenzeller, J. Knoch, and P. Avouris, “High-performance carbon nanotube field-effect transistor with tunable polarities, *IEEE Transactions on Nanotechnology*, vol.4, no.5, pp. 481- 489, Sept. 2005.
- [143] J. Appenzeller, J. Knoch, V. Derycke, R. Martel, S. Wind, and P. Avouris, “Field-Modulated Carrier Transport in Carbon Nanotube Transistors,” *Physical Review Letters*, vol. 89, no. 12, pp. 126801-126804, 2002.
- [144] A. A. Kane, T. Sheps, E. T. Branigan, V. A. Apkarian, M. H. Cheng, J. C. Hemminger, S. R. Hunt, and P. G. Collins, “Graphitic Electrical Contacts to Metallic Single-Walled Carbon Nanotubes Using Pt Electrodes,” *Nano Letters*, vol. 9, no. 10, pp. 3586-3591, 2009.
- [145] N. Patil, A. Lin, E.R. Myers, K. Ryu; A. Badmaev, C. Zhou; H.-S.P. Wong and S. Mitra, “Wafer-Scale Growth and Transfer of Aligned Single-Walled Carbon



- Nanotubes,” *IEEE Transactions on Nanotechnology*, vol. 8, no. 4, pp.498-504, July 2009.
- [146] S. Frégonèse, C. Maneux and T. Zimmer, “A Compact Model for Dual-Gate One-Dimensional FET: Application to Carbon-Nanotube FETs,” *IEEE Transactions on Electron Devices*, vol.58, no.1, pp.206-215, Jan. 2011.
- [147] J. Knoch and J. Appenzeller, “Tunneling phenomena in carbon nanotube field-effect transistors,” *physica status solidi (a)*, vol. 205, no. 4, pp. 679-694, 2008.
- [148] S. Frégonèse, C. Maneux and T. Zimmer, “Implementation of Tunneling Phenomena in a CNTFET Compact Model,” *IEEE Transactions on Electron Devices*, vol. 56, no. 10, pp. 2224-2231, Oct. 2009.
- [149] J. W. G. Wildöer, L. C. Venema, A. G. Rinzler, R. E. Smalley and C. Dekker, “Electronic structure of atomically resolved carbon nanotubes,” *Nature*, vol. 391, no. 6662, pp. 59-62, 1998.
- [150] J. Liu, I. O’Connor, D. Navarro and F. Gaffiot, “Dynamically reconfigurable CNTFET logic cell matrix programming method”, *VLSI SoC*, 2008.
- [151] N.F. Goncalves, H. De Man, “NORA: a racefree dynamic CMOS technique for pipelined logic structures,” *IEEE Journal of Solid-State Circuits*, vol.18, no.3, pp. 261-266, Jun 1983.
- [152] K.H. Yeo, S. D. Suk; M. Li, Y.-y. Yeoh, K. H. Cho, K.-H. Hong, S.K. Yun; M.S. Lee, N. Cho, K. Lee, D. Hwang, B. Park, D.-W. Kim, D. Park; B.-I. Ryu, “Gate-All-Around (GAA) Twin Silicon Nanowire MOSFET (TSNWFET) with 15 nm Length Gate and 4 nm Radius Nanowires”, *International Electron Devices Meeting 2006*, pp.1-4, 11-13 Dec. 2006.
- [153] Y. Huang, X. Duan, Y. Cui, L. Lauhon, K. Kim, and C. M. Lieber “Logic Gates and Computation from Assembled Nanowire Building Blocks,” *Science* 294, 1313-1317 (2001).
- [154] A. DeHon and M. J. Wilson, “Nanowire-based sublithographic programmable logic arrays”, *Proceedings of the 2004 ACM/SIGDA 12<sup>th</sup> international Symposium on Field Programmable Gate Arrays*, 2004.
- [155] Y. Wu, J. Xiang, C. Yang, W. Lu and C.M. Lieber, “Single-crystal metallic nanowires and metal/semiconductor nanowire heterostructures”, *Nature* 430, 61-65 (2004).
- [156] A. DeHon, “Array-based architecture for FET-based, nanoscale electronics,” *IEEE Transactions on Nanotechnology*, vol.2, no.1, pp.23-32, March 2003
- [157] T. Wang, P. Narayanan, M. Leuchtenburg and C.A. Moritz, “NASICs: A nanoscale fabric for nanoscale microprocessors,” *IEEE International Nanoelectronics Conference*, 24-27 March 2008.
- [158] R. S. Wagner, W. C. Ellis, “Vapour-liquid-solid mechanism of single crystal growth,” *Applied Physics Letters*, 4, 1964.
- [159] Y. Cui, L. J. Lauhon, M. S. Gudixsen, J. Wang, and C. M. Lieber, “Diameter-controlled synthesis of single crystal silicon nanowires,” *Applied Physics Letters*, vol.78, no.15, pp.2214–2216, 2001.
- [160] Y. Wu, Y. Cui, L. Huynh, C. J. Barrelet, D. C. Bell, and C. M. Lieber, “Controlled growth and structures of molecular-scale silicon nanowires,” *Nanoletters*, vol. 4, no. 3, pp. 433–436, 2004.
- [161] C. J. Kim, D. Lee, H. S. Lee, G. Lee, G. S. Kim, and M. H. Jo, “Vertically aligned Si intranowire p-n diodes by large-area epitaxial growth,” *Applied Physics Letters*, vol. 94, 173105, 2009.

- [162] L.J. Lauhon, M.S. Gudiksen, D. Wang, and C.M. Lieber, "Epitaxial Core-Shell and Core-Multi-Shell Nanowire Heterostructures," *Nature* 420, 57-61, 2002.
- [163] A. Ulman, "An introduction to ultrathin organic films: from Langmuir-Blodgett to self assembly," Academic Press, New York, 1991.
- [164] R.K. Brayton, A.L. Sangiovanni-Vincentelli, C.T. McMullen and G.D. Hachtel, "Logic Minimization Algorithms for VLSI Synthesis," Kluwer Academic Publishers, 1984.
- [165] T. Wang, M. Ben-Naser, Y. Guo, C.A. Moritz, "Wire-Streaming Processors on 2-D Nanowire Fabrics," Nano Science and Technology Institute, California, May 2005.
- [166] D. Whang, S. Jin, Y. Wu and C. M. Lieber, "Large-Scale Hierarchical Organization of Nanowire Arrays for Integrated Nanosystems," *Nano Letters*, vol. 3, no.9, pp. 1255-1259, 2003.
- [167] B. Cousin, M. Reyboz, O. Rozeau, M. Jaud, T. Ernst and J. Jomaah, "A Continuous Compact Model of Short-Channel Effects for Undoped Cylindrical Gate-All-Around MOSFETs," 9th Workshop on Compact Modeling, pp. 793-796, 2010.
- [168] B. Cousin, M. Reyboz, O. Rozeau, M.-A. Jaud, T. Ernst and J. Jomaah, "A unified short-channel compact model for cylindrical surrounding-gate MOSFET", *Solid-State Electronics*, vol.56, no.1, pp.40-46, February 2011.
- [169] D. B. Strukov and K. K. Likharev, "CMOL FPGA: a reconfigurable architecture for hybrid digital circuits with two-terminal nanodevices", *Nanotechnology*, vol. 16, no. 6, 2005.
- [170] Y. Chen, G.-Y. Jung, D. A. Ohlberg, X. Li, D. R. Stewart, J. O. Jeppesen, K. A. Nielsen, J. F. Stoddart and R. S. Williams, "Nanoscale molecular-switch crossbar circuits", *Nanotechnology*, vol. 14, no. 4, 2003.
- [171] T. Ernst T. Ernst, L. Duraffourg, C. Dupré, E. Bernard, P. Andreucci, S. Bécu, E. Ollier, A. Hubert, C. Halté, J. Buckley, O. Thomas, G. Delapierre, S. Deleonibus, B. de Salvo, P. Robert, and O. Faynot, "Novel Si-based nanowire devices: Will they serve ultimate MOSFETs scaling or ultimate hybrid integration?", *IEEE International Electron Devices Meeting*, 2008.
- [172] C. Dupré, A. Hubert, S. Becu, M. Jublot, V. Maffini-Alvaro, C. Vizioz, F. Aussenac, C. Arvet, S. Barnola, J.-M. Hartmann, G. Garnier, F. Allain, J.-P. Colonna, M. Rivoire, L. Baud, S. Pauliac, V. Loup, P. Rivallin, B. Guillaumot, G. Ghibaud, O. Faynot, T. Ernst, and S. Deleonibus, "15nm-diameter 3D Stacked Nanowires with Independent Gates Operation:  $\Phi$ FET", *IEEE International Electron Devices Meeting*, 15-17 Dec. 2008.
- [173] D. Lenoble and A. Grouillet, "The fabrication of advanced transistors with plasma doping", *Surface and coatings technology*, vol. 156, pp. 262-266, 2002.
- [174] A. DeHon, "Architecture approaching the Atomic Scale", *Proceedings of the 33rd European Solid-State Circuits Conference, ESSCIRC 2007*.
- [175] E. Ahmed and J. Rose, "The effect of LUT and cluster size on deep-submicron FPGA performance and density", *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol.12, no.3, pp. 288- 298, March 2004.
- [176] T. Ye, L. Benini and G. De Micheli, "Packetization and routing analysis of on-chip multiprocessor networks", *Journal of Systems Architecture*, vol.50, no.2-3, pp.81-104, 2004.

- [177] D.L. Lewis, S. Yalamanchili and H.-H. S. Lee, “High Performance Non-blocking Switch Design in 3D Die-Stacking Technology”, Proceedings of the 2009 IEEE Computer Society Annual Symposium on VLSI, 25-30, 2009.
- [178] H. Mrabet, Z. Marrakchi, P. Souillot and H. Mehrez, “Performances Improvement of FPGA using Novel Multilevel Hierarchical Interconnection Structure”, International Conference on Computer-Aided Design (ICCAD'2006), , pp. 675-679, November 2006.
- [179] Z. Marrakchi, H. Mrabet, C. Masson and H. Mehrez, “Performances Comparison between Multilevel Hierarchical and Mesh FPGA Interconnects”, International Journal of Electronics, vol. 95, no. 3, pp. 275-289, January 2008.
- [180] C.-L. Wu and T.-Y. Feng, “On a Class of Multistage Interconnection Networks”, IEEE Transactions on Computers, vol.C-29, no.8, pp.694-702, 1980.
- [181] G.B. Adams, D.P. Agrawal, and H.J. Siegel, “A Survey and Comparison of Fault-Tolerant Multistage Interconnection Networks”, Computer, 20(6), 14-27, 1987.
- [182] Verilog-To-Routing (VTR) Project, <http://www.eecg.utoronto.ca/vtr/>
- [183] ABC: Berkeley logic synthesis tool, <http://www.eecs.berkeley.edu/~alanmi/abc/>
- [184] ODIN-II: Verilog synthesis tool, <https://code.google.com/p/odin-ii/>
- [185] J. Luu, I. Kuon, P.Jamieson, T. Campbell, A. Ye, M. Fang, and J. Rose, “VPR 5.0: FPGA CAD and Architecture Exploration Tools with Single-Driver Routing, Heterogeneity and Process Scaling”, ACM Symposium on FPGAs, FPGA '09, February 2009, pp. 133-142.
- [186] Versatile packing, placement and routing tool for FPGA, <http://www.eecg.utoronto.ca/vpr/>
- [187] J. Luu, J. Anderson and J. Rose, “Architecture description and packing for logic blocks with hierarchy, modes and complex interconnect”, ACM/SIGDA 19th International Symposium on Field Programmable Gate Arrays, Monterey, 2011.
- [188] M. De Marchi, M. H. Ben Jamaa and G. De Micheli, “Regular Fabric Design with Ambipolar CNTFETs for FPGA and Structured ASIC Applications”, International Symposium on Nanoscale Architectures (NANOARCH), June 2010.
- [189] N. J. Macias, L. J. K. Durbeck, “Adaptive methods for growing electronic circuits on an imperfect synthetic matrix”, Biosystems, vol.73, no.3, pp.173-204, March 2004.
- [190] M. H. Ben Jamaa, P.-E. Gaillardon, S. Frégonèse, M. De Marchi, G. De Micheli, T. Zimmer, I. O'Connor and F. Clermidy, “FPGA Design with Double-Gate Carbon Nanotube Transistors”, The Electro-Chemical Society Transactions, vol. 34, no. 1, pp. 495-501, 2011.
- [191] J. Birkner, A. Chan, H.T. Chua, A. Chao, K. Gordon, B. Kleinman, P. Kolze, R. Wong, “A very-high-speed field-programmable gate array using metal-to-metal antifuse programmable elements”, Microelectronics Journal, vol.23, no.7, November 1992.
- [192] E. M. Sentovich, K. J. Singh, L. Lavagno, C. Moon, R. Murgai, A. Saldanha, H. Savoj, P. R. Stephan, R. K. Brayton, and A. L. Sangiovanni-Vincentelli, “SIS: A System for Sequential Circuit Synthesis”, Technical Report UCB/ERL M92/41, Electronics Research Lab, University of California, Berkeley, CA 94720, May 1992.
- [193] DOT language – Graphviz tool: <http://www.graphviz.org/Documentation.php>
- [194] A.S. Marquardt, V. Betz and J. Rose, “Using cluster-based logic blocks and timing-driven packing to improve FPGA speed and density”, ACM/SIGDA 7<sup>th</sup> International symposium on Field programmable gate arrays, pp.37-46, 1999.

- [195] E. Bozorgzadeh, S. Memik, X. Yang, and M. Sarrafzadeh, "Routability-driven Packing: Metrics and Algorithms for Cluster-Based FPGAs" *Journal of Circuits Systems and Computers*, vol.13, pp.77–100, 2004.
- [196] A. Singh, G. Parthasarathy and M. Marek-Sadowksa, "Efficient Circuit Clustering for Area and Power Reduction in FPGAs", *ACM Transaction on Design Automation of Electronic Systems*, vol.7, no.4, pp.643-663, Nov 2002.
- [197] D. Chen, K. Vorwerk and A. Kennings, "Improving Timing-Driven FPGA Packing with Physical Information", *International Conference on Field Programmable Logic and Applications*, pp.117-123, 2007.
- [198] G. Lemieux and D. Lewis, "Design of Interconnection Networks for Programmable Logic", Kluwer Academic Publishers, 2004.
- [199] J. Lin, D. Chen and J. Cong, "Optimal Simultaneous Mapping and Clustering for FPGA Delay Optimization", *ACM/IEEE Design Automation Conference*, pp. 472–477, 2006.
- [200] A. Ling, J. Zhu and S. Brown, "Scalable Synthesis and Clustering Techniques Using Decision Diagrams", *IEEE Transactions on CAD*, vol.27, no.3, pp.423, 2008.
- [201] K. Wang, M. Yang, L. Wang, X. Zhou and J. Tong, "A novel packing algorithm for sparse crossbar FPGA architectures", *International Conference on Solid-State and Integrated-Circuit Technology*, pp. 2345-2348, 2008.
- [202] P. Lambin, A.A., Lucas and J.C. Charlier, "Electronic properties of carbon nanotubes containing defects", *Journal of Physics and Chemistry of Solids*, vol.58, no.11, pp.1833-1837, 1997.
- [203] E.K. Vida-Torku, W., Reohr, J.A., Monzel and P. Nigh, "Bipolar, CMOS and BiCMOS circuit technologies examined for testability", *34th Midwest Symposium on Circuits and Systems*, pp.1015-1020, 14-17 May 1991.
- [204] BLIF circuit benchmarks: <http://cadlab.cs.edu/~kirill/>
- [205] N. Yakymets, I. O'Connor. Patent FR0958957 "Matrice interconnectee de cellules logiques reconfigurables avec une topologie d'interconnexion croisee", 2010.
- [206] N. Yakymets, K. Jabeur, I. O'Connor, and S. Le Beux, "Interconnect Topology for Cell Matrices Based on Low-Power Nanoscale Devices", *Faible Tension Faible Consommation*, Marrakech, Morocco, May 30 – June 1, 2011.
- [207] A. Lin, N. Patil, H. Wai, S. Mitra and H.-S.P. Wong, "A metallic-CNT-tolerant carbon nanotube technology using Asymmetrically-Correlated CNTs (ACCNT)", *Symposium on VLSI Technology*, pp.182-183, 16-18 June 2009.
- [208] N. Patil, A. Lin, E. Myers, H.-S.P. Wong, and S. Mitra, "Integrated wafer-scale growth and transfer of directional Carbon Nanotubes and misaligned-Carbon-Nanotube-immune logic structures", *Symposium on VLSI Technology*, pp.205-206, 17-19 June 2008.
- [209] S. K. Bobba, A. Chakraborty, O. Thomas, P. Batude, T. Ernst, O. Faynot, D. Z. Pan, and G. De Micheli, "CELONCEL: Effective Design Technique for 3-D Monolithic Integration targeting High Performance Integrated Circuits," *16th Asia and South Pacific Design Automation Conference*, pp.337-343, 2011.

# *List of Publications*

---

## **Articles published in refereed publications (journals)**

- [1] M. H. Ben-Jamaa, **P.-E. Gaillardon**, F. Clermidy, I. O'Connor, D. Sacchetto, G. De Micheli, Y. Leblebici, "Silicon Nanowire Arrays and Crossbars: Top-Down Fabrication Techniques and Circuit Applications", *Science of Advanced Materials*, vol.3, no.3, pp.466-476, June 2011.
- [2] **P.-E. Gaillardon**, I. O'Connor, M. Amadou, J. Liu, F. Clermidy, G. Nicolescu, "Matrix Nanodevice-Based Logic Architectures and Associated Functional Mapping Method", *ACM Journal on Emerging Technologies in Computing Systems*, vol. 7, no. 1, pp. 3:1-3:23, January 2011.
- [3] I. O'Connor, J. Liu, D. Navarro, R. Daviot, N. Abouchi, **P.-E. Gaillardon**, F. Clermidy, "Molecular electronics and reconfigurable logic", *International Journal of Nanotechnology*, vol. 7, no. 4/5/6/7/8 pp. 367 - 382, 2010.

## **Articles published in refereed international conference proceedings**

- [1] **P.-E. Gaillardon**, M. H. Ben-Jamaa, F. Clermidy, I. O'Connor, "Ultra-Fine Grain FPGAs: A Granularity Study", *IEEE/ACM International Symposium on Nanoscale Architectures (NanoArch)*, 08-09 June 2011, San Diego (CA), USA.
- [2] S. Onkaraiah, **P.-E. Gaillardon**, M. Reyboz, F. Clermidy, J.-M. Portal, M. Bocquet, C. Muller, "Using OxRRAM Memories for Improving Communications of Reconfigurable FPGA Architectures", *IEEE/ACM International Symposium on Nanoscale Architectures (NanoArch)*, 08-09 June 2011, San Diego (CA), USA.
- [3] **P.-E. Gaillardon**, M. H. Ben-Jamaa, P.-H. Morel, J.-P. Noël, F. Clermidy, I. O'Connor, "Can We Go Towards True 3-D Architectures?", *WACI session, 48th Design Automation Conference (DAC)*, 5-10 June 2011, San Diego (CA), USA.
- [4] **P.-E. Gaillardon**, M. H. Ben-Jamaa, F. Clermidy, I. O'Connor, "Evaluation of a Crossbar Multiplexer in a Lithography-Based Nanowire Technology", *IEEE International Symposium on Circuits and Systems (ISCAS)*, 15-18 May 2011, Rio de Janeiro, Brazil.
- [5] M. H. Ben Jamaa, **P.-E. Gaillardon**, S. Frégonèse, M. De Marchi, G. De Micheli, T. Zimmer, I. O'Connor and F. Clermidy, "FPGA Design with Double-Gate Carbon Nanotube Transistors", *The Electro-Chemical Society Transactions*, vol. 34, no. 1, pp. 495-501, 2011.
- [6] **P.-E. Gaillardon**, M. H. Ben-Jamaa, G. Betti Beneventi, F. Clermidy, L. Perniola, "Emerging Memory Technologies for Reconfigurable Routing in FPGA Architecture", *IEEE International Conference on Electronics, Circuits and Systems (ICECS)*, 12-15 December 2010, Athens, Greece.

- [7] **P.-E. Gaillardon**, M. Haykel Ben-Jamaa, M. Reyboz, G. Betti Beneventi, F. Clermidy, L. Perniola, I. O'Connor, "Phase-Change-Memory-Based Storage Elements for Configurable Logic" , International Conference on Field-Programmable Technology (FPT), 8-10 December 2010, Beijing, China.
- [8] K. jabeur, **P.-E. Gaillardon**, D. Navarro, I. O'Connor, M. H. Ben-Jamaa, F. Clermidy, "Reducing transistor count in clocked standard cells with ambipolar double-gate FETs", IEEE/ACM International Symposium on Nanoscale Architectures (NanoArch), June 17-18 2010, Anaheim (CA), USA.
- [9] **P.-E. Gaillardon**, F. Clermidy, I. O'Connor, R. Daviot, "Reconfigurable Reconfigurable nanoscale logic cells : a comarison study", IEEE International Conference on Electronics, Circuits and Systems (ICECS), 13-16 December 2009, Hammamet, Tunisia.
- [10] **P.-E. Gaillardon**, I. O'Connor, J. Liu, F. Clermidy, R. Daviot, "Mapping method of reconfigurable cell matrices based on nanoscale devices using inter-stage fixed interconnection scheme," IEEE International Conference on Electronics, Circuits and Systems (ICECS), 13-16 December 2009, Hammamet, Tunisia.
- [11] **P.-E. Gaillardon**, I. O'Connor, J. Liu, F. Clermidy, "Interconnection scheme and associated mapping method of reconfigurable cell matrices based on nanoscale devices", IEEE/ACM International Symposium on Nanoscale Architectures (NanoArch), July 30-31 2009, San Francisco (CA), USA.

## Patents

- [1] **P.-E. Gaillardon**, G. Betti Beneventi, L. Perniola, "Cellule Memoire", FR application 11 52127, 15 March 2011.
- [2] **P.-E. Gaillardon**, F. Clermidy, I. O'Connor, P.-H. Morel, "Reconfigurable Boolean Cells having a Criss-crossed Nanowire Matrix," WO2011070164 PCT application, FR2954022 application, 11 December 2009.

# *List of Acronyms*

---

2-D	Two Dimensions
3-D	Three Dimensions
5T-SRAM	Five Transistors Static Random Access Memory
ALD	Atomic Layer Deposition
ANN	Artificial Neural Network
ASIC	Application Specific Integrated Circuit
BEC	Bottom EleCtrode
BEOL	Back-End-Of-Line
BJT	Bipolar Junction Transistor
BLE	Basic Logic Element
BLIF	Berkeley Logic Interchange Format
BTB	Band-To-Band
CAD	Computer Aided Design
CB	Connection Box
CB-NWFET	CrossBar of NanoWires Field Effect Transistors
CBRAM	Conductive-Bridging Random Access Memory
CLB	Configurable Logic Block
CMOL	CMos/nanOeLectronic
CMOS	Complementary Metal-Oxide-Semiconductor
CN	Carbon Nanotube
CNLB	Configurable “Nano” Logic Block
CNT	Carbon NanoTube
CNTFET	Carbon NanoTube Field Effect Transistor
CVD	Chemical-Vapor-Deposition
DG-CNFET	Double-Gate Carbon Nanotube Field Effect Transistor
DIBL	Drain-Induced Barrier Lowering
DRAM	Dynamic Random Access Memory
DSP	Digital Signal Processor/Processing
DUV	Deep UltraViolet
EDA	Electronic Design Automation

EPFL	Ecole Polytechnique Fédérale de Lausanne
ERD	Emerging Research Devices
ERM	Emerging Research Materials
EUV	Extreme UltraViolet
EV	EValuation
FDSOI	Fully Depleted Silicon-On-Insulator
FEOL	Front-End-Of-Line
FET	Field Effect Transistor
FF	Flip-Flop
FGLC	Fine Grain Logic Cell
FO4	FanOut-of-four
FPGA	Field Programmable Gate Array
GAA	Gate-All-Around
GIDL	Gate-Induced Drain Leakage
GNR	Graphene NanoRibbon
GST	$\text{Ge}_2\text{Sb}_2\text{Te}_5$
HVT	High Threshold Voltage
IC	Integrated Circuit
IMS	Laboratoire d'Intégration du Matériau au Système
INL	Institut des Nanotechnologies de Lyon
IP	Intellectual Property
ITRS	International Technology Roadmap for Semiconductor
LB	Logic Block
LETI	Laboratoire d'Electronique et de Technologie de l'Information
LP	Low Power
LUT	Look-Up Table
LVT	Low Threshold Voltage
MAC	Multiplier Accumulator
MCluster	Matrix Cluster
MCNC	Microelectronics Center of North Carolina
MIN	Multistage Interconnection Network
MOS	Metal-Oxide-Semiconductor
MOSFET	Metal-Oxide-Semiconductor Field Effect Transistor
MPPA	Massively-Parallel Processor Array
MPSOC	Multi-Processor System-On-Chip
MPW	Multi Project Wafer
MRAM	Magnetic Random Access Memory
MUX	MULTipleXer



NASIC	Nanoscale Application Specific Integrated Circuit
NEMS	Nano-Electro Mechanical System
NOC	Network-On-Chip
NVM	Non-Volatile Memory
NW	NanoWire
NWFET	NanoWire Field Effect Transistor
OxRAM	Oxide Random Access Memory
PAL	Programmable Array Logic
PC	Phase-Change
PC	PreCharge
PCM	Phase-Change Memory
PDP	Power-Delay Profile
PIDS	Process Integration, Devices and Structures
PLA	Programmable Logic Array
PLAD	PLAsma Doping
PSP	State University-Philips
RAM	Random Access Memory
RR	Routing Resource
RRAM	Resistive Random Access Memory
ReRAM	Resistive Random Access Memory
SB	Switch Box
SCE	Short-Channel Effects
SNOW	Silicon Nanoelectronics On 300-mm Wafer
SiNW	Silicon NanoWires
SOC	System-On-Chip
SOI	Silicon-On-Insulator
SRAM	Static Random Access Memory
T-VPACK	Timing-driven Versatile PACKer
TCAD	Technology Computer Aided Design
TEC	Top EleCtrode
TSV	Through-Silicon-Vias
UFF	Unbalanced Flip-Flop
ULSI	Ultra Large Scale Integration
VLSI	Very Large Scale Integration
VPACK	Versatile PACKer
VPR	Versatile Place and Route tool
VTR	Verilog-To-Routing suite
KB	Wentzel-Kramers-Brillouin



# **APPENDIX A** *Memento of Presented Technologies*

---

A.1	Taxonomy of Emerging Technologies .....	xxii
A.2	State Variable Change: Phase-Change Memories .....	xxii
A.2.1	Brief Literature Overview .....	xxii
A.2.2	Technological Assumptions .....	xxiii
A.3	Material Change: Ambipolar Carbon Nanotubes .....	xxiii
A.3.1	Brief Literature Overview .....	xxiii
A.3.2	Technological Assumptions .....	xxv
A.4	Device Change: Monolithic 3-D Integration .....	xxvi
A.4.1	Brief Literature Overview .....	xxvi
A.4.2	Technological Assumptions .....	xxvii
A.5	Device Change: 1-D Active Element Structures .....	xxvii
A.5.1	1-D Dense NanoWire Crossbar Integration .....	xxviii
A.5.1.1	Brief Literature Overview .....	xxviii
A.5.1.2	Technological Assumptions .....	xxviii
A.5.2	Lithographic Crossbar Integration .....	xxix
A.5.2.1	Brief Literature Overview .....	xxix
A.5.2.2	Technological Assumptions .....	xxix
A.5.3	Vertical NanoWire FET Integration .....	xxx
A.5.3.1	Brief Literature Overview .....	xxx
A.5.3.2	Technological Assumptions .....	xxxi
A.6	Overall Comparison .....	xxxii

This aim of this appendix is to collate and review the information related to the technologies described in this thesis. Thus, it is expected to serve as a memento during the reading of this manuscript. We will first describe the taxonomy of the emerging technologies, in order to sort them according to their various novel properties. Then, each envisaged technology will be presented regarding their common properties. A brief literature review will be presented for each technology and the specific assumptions made in this thesis work will be described.

## A.1 Taxonomy of Emerging Technologies

The ITRS ERD and ERM chapters also suggest the possibility of using novel properties at different levels of hierarchy. Figure a-a presents the taxonomy of emerging devices. In this appendix, we only deal with the technology levels. Hence, the novelties can be identified on the three lower levels: the state variable level, the material level and finally the device level. Conversely to the organization in this thesis, we will group the explored technologies according to this taxonomy. Firstly, we will discuss the technology of PCM that uses the state of the material as its retention variable. Then, we will deal with the change of material to carbon, in order to manage the electronic charge with ambivalence. Finally, by moving to the device level, we will consider three processes that use scaled 1-D CMOS structures to improve the integration density and 3-D integration process that move the transistors into unconventional layers.

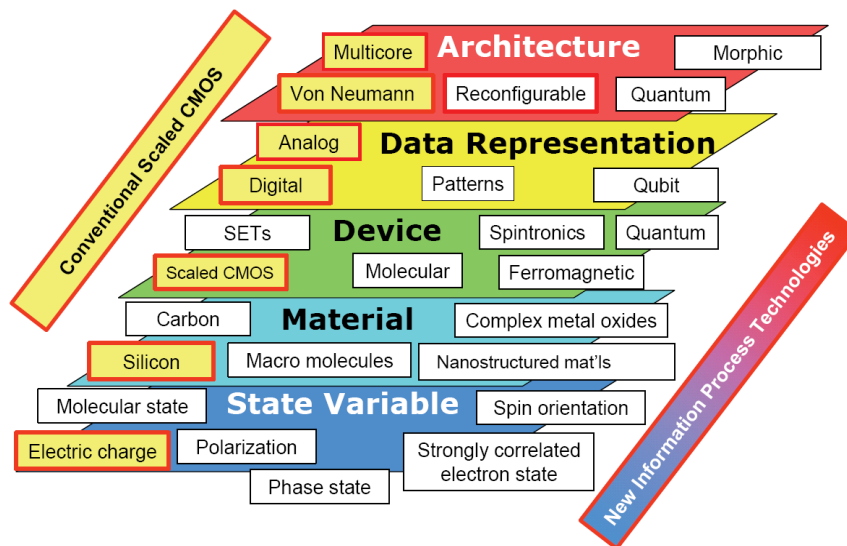


Figure a-a. Electronic systems hierarchical organization and opportunities [7]

## A.2 State Variable Change: Phase-Change Memories

### A.2.1 Brief Literature Overview

As introduced in figure a-b, several technologies could use efficiently novel state variables rather than common electronic charge. We especially consider here storage of the information through the resistance of a material. Phase-Change Memories are, as the name indicates, based on a material having two different stable physical phases leading to two different resistivity values. Several materials can be used, such as GeSbTe [82], GeTe [83], GeTeC [116] ... PCMs are considered today to be one of the most promising candidates for the next generation of non-volatile memory applications [70]. The interest in PCMs is due to various advantages, including: better scalability (down to a few nanometers) [71], faster programming time (of the order of few nanoseconds) [72] and improved endurance (up to  $10^9$  programming cycles) [73]. Some prototypes (such as a 60-nm 512-Mb [74] and a 45-nm 1-Gb [73] PCM technology) have been presented recently to showcase the viability of high-density standalone memories based on PCM technology from an industrial point of view.

## A.2.2 Technological Assumptions

PCM technology is CMOS-compatible. As in Flash-NOR arrays, each memory cell includes a storage phase-change node and a selector transistor in series (i.e. *1-resistor-1-transistor* configuration). The memory element may be fabricated either just after the Si contact forming step at the *Front-End-Of-Line* (FEOL) level or after the first steps of interconnections at the *Back-End-Of-Line* (BEOL) level, (e.g. on top of the Metal 0 or Metal 1 interconnect level) [75]. A schematic cross-section of the storage element architecture is shown in figure a-b. The PCM device, formed of a PC layer with Bottom (BEC) and Top (TEC) Electrode Contacts, is integrated between the M0 and M1 interconnection levels in the back-end-of-line. The MOSFET selector (bottom) is fabricated in the front-end-of-line. The pillar approach depicted in the figure is the simplest way to create a PCM device. First, a metallic heater is built. The heater is made by etching a via into the inter-layer dielectric and by filling it with a metal. The role of the heater is to help to channel the current in order to increase its density and thus maximize the heat control in the memory node. To improve the heater fabrication, several sublithographic techniques have been proposed [68, 69]. After the heater metal deposition, the via is filled by chalcogenide alloys with a room temperature deposition. The top electrode is obtained by a final metal deposition.

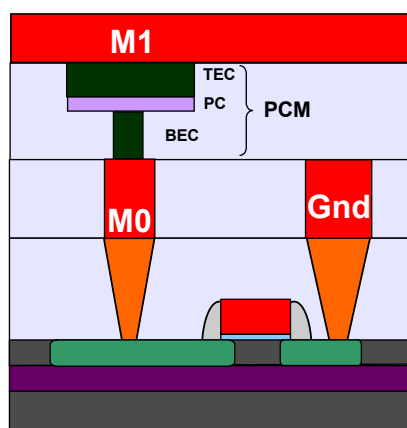


Figure a-b. Cross sectional schematic showing a PCM device integration.

## A.3 Material Change: Ambipolar Carbon Nanotubes

### A.3.1 Brief Literature Overview

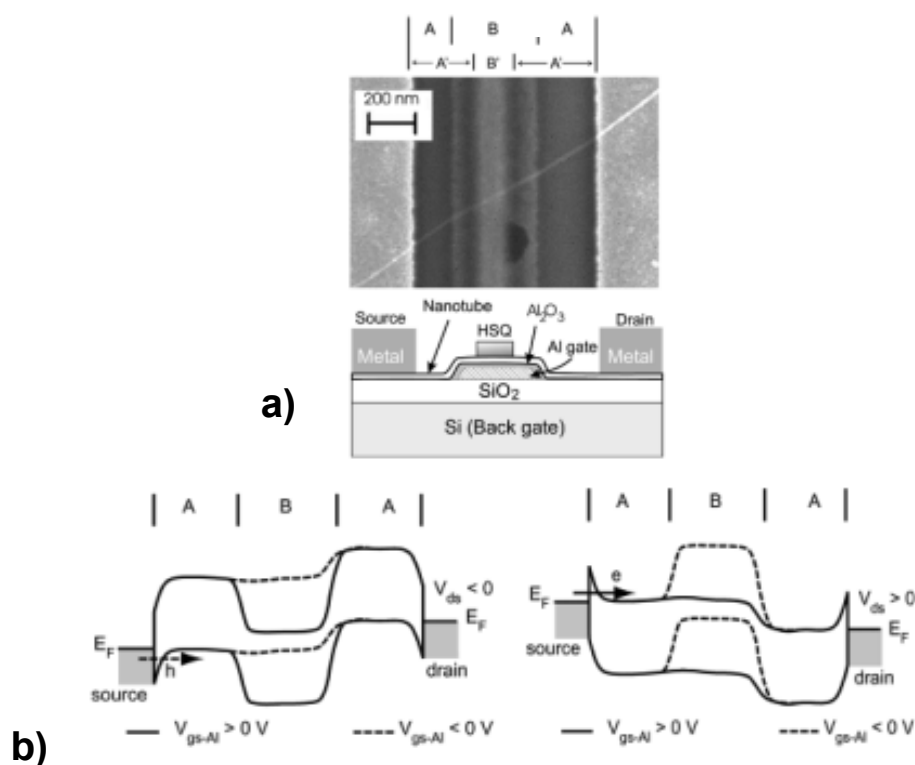
Since the discovery of *Carbon Nanotubes* (CNT) in the early 1990s [123] and the first isolation of a graphene sheet in 2004 [124], carbon electronics has witnessed a growing interest. In particular, the study of the intrinsic properties of carbon and its potential application to microelectronic devices has been of high interest.

Carbon nanotubes, which usually consist of a single atom thick sheet of carbon rolled up to form a seamless tube [137], possess exceptional electrical properties such as high current carrying capability ( $>10^9 \text{A/cm}^2$ ) [138] and excellent carrier mobility ( $9000 \text{cm}^2/\text{V.s}$ ) [139]. Due to the small diameter ( $\sim 1 \text{nm}$ ), nanotubes are ideal candidates to provide one-dimensional (1-D) electrical transport. Ballistic transport, even at room temperature, has been demonstrated over short distances ( $\lambda_{\text{MFP}} \approx 700 \text{nm}$ ) [139]. It should however be noted that only carbon nanotubes built from a single layer of carbon (single-wall carbon nanotubes) can behave as semi-conductors. Multiple-wall carbon nanotubes are composed of a number of coaxial single-wall nanotubes. High-performance electronics applications only use single-wall tubes, due to their small size and good semiconducting behavior.

Most *Carbon Nanotubes Field Effect Transistors* (CNFETs) studied so far have adopted a back-gate top-contact geometry. In this back-gate configuration, the nanotubes are dispersed or grown on a conducting substrate covered by an insulating layer. Two metal contacts are deposited on the nanotube to serve as source and drain electrodes, while the conducting

substrate is the gate electrode in this three-terminal device. The transistor exhibits an ambipolar transistor characteristic. Conversely to a conventional *Metal Oxide Semiconductor Field Effect Transistor* (MOSFET), switching of a CNFET is dominated by the modulation of Schottky barriers formed at the nanotube/metal contacts [143]. Since carbon nanotubes are intrinsically ambivalent, meaning that both electrons and holes can flow through the structure, the height of the Schottky barriers allows the selection of the carrier conduction. At sufficiently negative (resp. positive) gate voltage, the Schottky barrier is sufficiently thinned to enable hole (resp. electron) injection from the source (resp. drain) contact into the nanotube. In the behavior of a conventional FET, these barriers are not desirable because they reduce the controllability of the device. Several works have been proposed to improve the contact quality. In [144], a graphitic wet layer is used to improve the contact between metal and nanotube.

In [142], another approach is taken. Instead of eliminating the ambipolarity of the devices, a dual-gate approach is introduced. The dual-gate creates pure  $n$ - and/or  $p$ -type devices with excellent off-state performance, simply by using electrostatic control on the carriers. Figure a-c-a shows the scanning electron microscopy image, as well as the device cross section of a *Double-Gate CNFET* (DG-CNFET) structure. The DG-CNFET possesses an additional Al gate electrode placed underneath the nanotube between the source and drain contacts. The Al gate region is denoted as the inner region (B), and the area between the Al gate and the source/drain contacts are denoted as the outer regions (A). In the design, the Al gate is the primary gate that governs the electrostatics and the switching of the nanotube bulk channel in the inner region. The Schottky barriers at the nanotube/metal contacts are controlled by the Si back gate (substrate). This prevents the electrostatics in the outer regions from being influenced by the Al gate. In the device, electrostatic doping effects in CNFETs is utilized to eliminate ambipolar characteristics in a Schottky barrier CNFET and to obtain a bulk-switched transistor possessing a tunable polarity ( $n$  or  $p$ ), steep subthreshold swing and excellent *off*-state performance. This  $p$ -FET and  $n$ -FET behaviour of the dual-gate CNFET can be understood by the schematic band diagrams shown in figure a-c-b.



**Figure a-c.** a) SEM image (top) and schematic cross-sectional diagram of a DG-CNFET [142] b) Band diagrams of the structure (Left  $V_{gs-Si} < 0$ , and right  $V_{gs-Si} > 0$ ) [142]

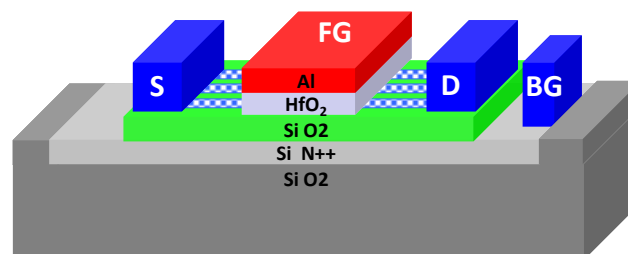
For a sufficiently negative (resp. positive) Si gate voltage, the Schottky barriers are thin enough to allow for hole (resp. electron) tunnelling from the metal contacts into the nanotube channel. Thus, regions A become electrostatically doped as  $p$ -type (resp.  $n$ -type), resulting in a  $p/i/p$  (resp.  $n/i/n$ ) band profile that allows only hole (resp. electron) transport in the nanotube channel. The dual-gate CNFET is switched *on* and *off* by varying the Al gate voltage that alters the barrier height for carrier transport across region B. In this configuration, regions A serve as extended source and drain, and the device operates similarly to a conventional MOSFET through bulk-switching in region B.

### A.3.2 Technological Assumptions

In [142], the authors used an integration process based on a generalized and common back-gate electrode. This process is simple for research and characterization purposes. Nevertheless, it requires a global and shared back-gate for all the devices. However for circuit design, the principal advantage of the device is its unique in-field reconfigurability, meaning that each back-gate needs an individual control.

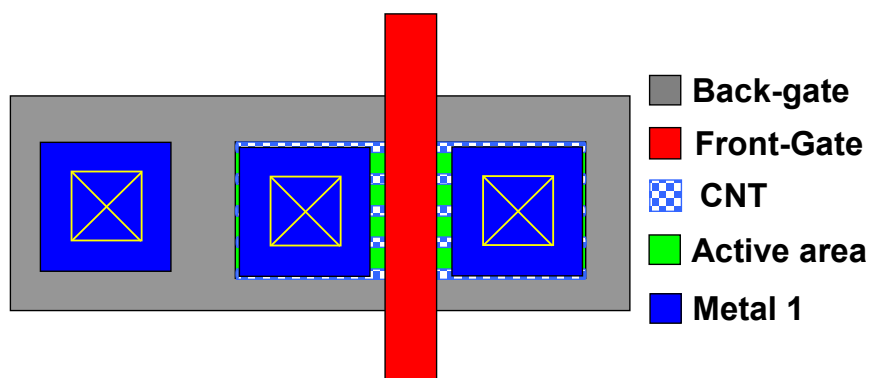
The proposed process flow is based on *Silicon-On-Insulator* (SOI) wafers. In this process, it is possible to build silicon mesas (i.e. islands of silicon surrounded by oxide) to realize the DG-CNFETs back-gate. The individual control requirement will thus be guaranteed. The final device is shown in figure a-d.

The expected technology process is based on a realistic and CMOS-compatible process flow. Starting with a SOI wafer, where the silicon-on-insulator layer will be used as back-gate,  $N^{++}$  doping (or NiSi salicidation) can be realized to ensure a good conductivity of the back electrodes.  $\text{SiO}_2$  is subsequently deposited as back-gate oxide, and intrinsic carbon nanotubes are transferred on top of this [145]. Then, the gate oxide ( $\text{HfO}_2$ ) and the metal (Al) of the top gate are deposited and patterned. Then, the active area and the back gate are defined by  $\text{SiO}_2$  and Si etching respectively. Subsequently, the metal is sputtered onto the contacts to drain, source and both gates.



**Figure a-d.** DG-CNFET device schematic using the proposed process-flow and showing the source (S), drain (D), front- (FG) and back-gate (BG) contacts

Using the previously defined process flow, we consider the layout requirements for our DG-CNTFET device. The expected layout is shown in figure a-e. This layout corresponds to parameterized-cell requirements and follows a standard industrial back-end rule set.



**Figure a-e.** DG-CNTFET device layout

## A.4 Device Change: Monolithic 3-D Integration

### A.4.1 Brief Literature Overview

IC integration in three dimensions appears to be a promising alternative path to scaling, and to some extent would avoid the huge investments required by scaling. While the concept is not entirely new [94], the development of the technology has witnessed significant growth over the last decade. In particular, the technology of *Through-Silicon-Vias* (TSVs) is currently the reference for 3-D technology processes. A TSV could be defined as a large via built through the substrate, in order to contact the active front-end to the reverse-side of the chip. Bumps are then used to contact the dies between them.

While such integration is challenging due to TSV density requirements (leading to a double specification of TSV aspect ratio and wafer thickness), the process is mature enough for industrial applications [95]. However, there are still some hurdles to face. Table a-a depicts the principal characteristics of TSV processes [96]. It is worth noticing that TSVs are quite area-hungry with diameter up to  $5\mu\text{m}$  and pitch up to  $10\mu\text{m}$ . This means that the connection density is quite poor (from 400 to 10000 TSVs/ $\text{mm}^2$ ) and this results in a loss of active size. Thus, designers are limited to high level interconnect, such as memory/core communication. However, this segregation is of great interest for performance, since processes will be tuned to optimize each layer to a given class of application (low power, general purpose ...)

Table a-a. Density survey of TSV technologies [96]

	Diameter ( $\mu\text{m}$ )	Pitch ( $\mu\text{m}$ )	Density (TSVs/ $\text{mm}^2$ )	CMOS 65-nm gate equivalent
<i>Low Density TSVs</i>	20	50	400	1250
<i>High Density TSVs</i>	5	10	10000	50

In a reconfigurable application, a large number of interconnections are required if separation between memory and logic is to be envisaged. This means that other integration processes should be used, to overcome the limitations of connection density. In an FPGA, connections between memory and logic are done at the gate level. We estimate the required density at about 500 000 3-D contacts/ $\text{mm}^2$ . Thus, instead of processing the layers separately and stacking them *a posteriori*, it is possible to use monolithic sequential integration. In monolithic integration, the circuit is processed from the bottom to the top. This means that the stacked layers are built by a set of technological steps above the already processed stack, as illustrated in figure a-f.

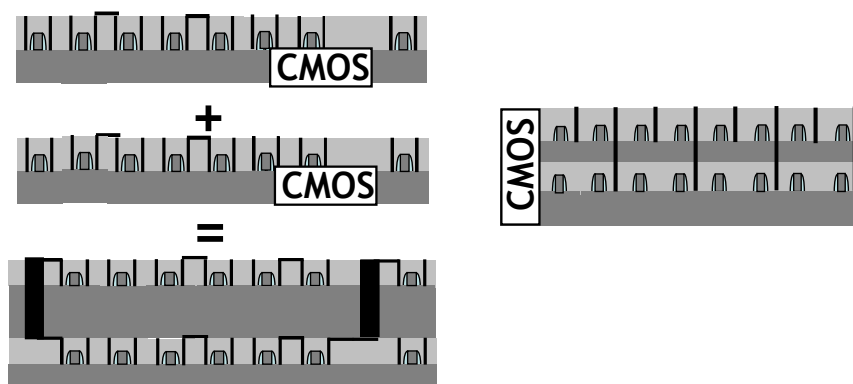


Figure a-f. Cross-sectional view of 3-D parallel integration (left) and 3-D monolithic sequential integration (right)

Such an integration scheme is promising for several reasons. Firstly, it enables a much higher via integration density. In fact, the contacts use a planar scheme step, as opposed to TSVs. Thus, it is possible to obtain alignment accuracy in the range of lithographical precision (around 10-nm), whereas TSV alignment accuracy is in the range of  $1\mu\text{m}$ . Various techniques



can be found in the literature, from which some particularly relevant publications in the field are [100, 101, 102, 103, 104, 105, 107].

#### A.4.2 Technological Assumptions

As briefly introduced above, the monolithic 3-D integration process integrates at least two layers of active silicon. These layers are processed sequentially, as shown in figure a-g. The principal process steps are (a) the realization of the bottom transistor, (b) the deposition of the top film silicon and (c) the realization of the top transistor. Several constraints appear at each step.

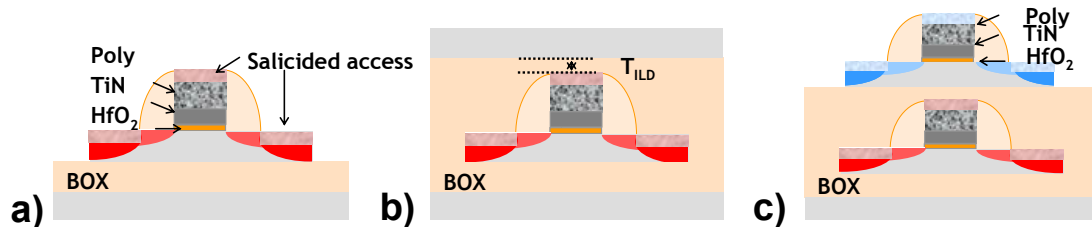


Figure a-g. Cross-sectional view of 3-D monolithic steps – a) optimized bottom FDSOI process b) high quality top film deposition c) low temperature top FDSOI process

The first transistor layer must be optimized in order to improve its robustness to temperature. It is worth noting that the second active transistor layer will have an impact on the first layer, implying that the thermal robustness of the bottom layer is critical. The maximal thermal budget for subsequent steps is limited by salicidation. In fact, standard NiSi dewetting occurs in only 3 minutes at 600°C. To obtain a low resistance access with processes around 600°C, a salicide stabilization procedure has been demonstrated in [97]. This salicidation approach is based on an optimized W-NiPtSi-F. With this technology, the salicide is stable at 600°C, which is considered to be the maximal thermal budget for all the subsequent steps.

The second step corresponds to the realization of a high quality top film. While seed window techniques were first used in the literature [98], they lead to crystalline defects, poor thickness control and density loss. It is thus better to use molecular bonding [97], where a blanket Silicon-On-Insulator wafer is transferred on top of the processed MOS wafer. This step plays a major role in alignment accuracy, since the alignment occurs after bonding, conversely to parallel integration [99].

A low temperature Fully Depleted Silicon-On-Insulator is then processed on the top film. Due to the limited thermal budget (600°C), it is not possible to perform thermal activation of dopants (around 1000°C-1100°C). The process uses a solid phase epitaxial re-growth [100], which consists of pre-amorphization of the top film, followed by dopant implantation and finishing by a re-crystallization at 600°C.

Finally, 3-D contacts are realized using the same contact techniques as in standard processes. Only one lithography step is required for all contacts.

#### A.5 Device Change: 1-D Active Element Structures

For several decades, the micro-electronics industry has been known to scale the transistor dimensions, in order to improve the performance. Nevertheless, this increase of performance is not only due to the scaling, but also due to several structural improvements.

Figure a-h shows the evolution of the silicon structures from bulk to nanowires. 1-D structures demonstrate good electrostatic controllability, which is fundamental to ensure good performance levels for scaled devices. While the trend in industry is towards the *Fully Depleted Silicon-On-Insulator* (FDSOI) process, we can note that nanowires tend to be the optimal solution, combining good electrostatic control with low *Off* current ( $I_{\text{off}} < 1\text{nA}/\mu\text{m}$ ) [152].

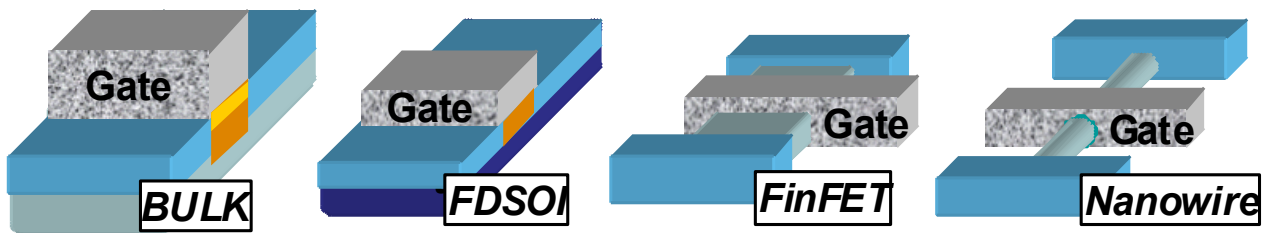


Figure a-h. Silicon electronics evolution from bulk to nanowires

Furthermore, in addition to their good expected properties, we can note that the 1-D structures make the active dimensions very small. This leads to potentially high-density integrated devices, especially when a crossbar organization is proposed.

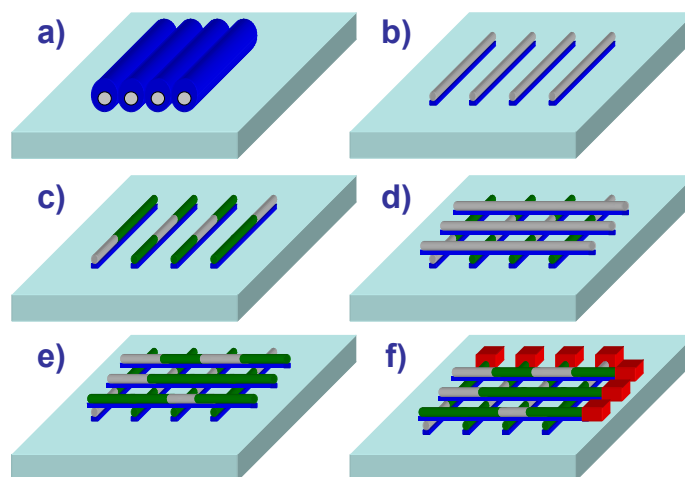
## A.5.1 1-D Dense NanoWire Crossbar Integration

### A.5.1.1 *Brief Literature Overview*

Thanks to bottom-up fabrication processes, nanowires are useful to build regular crossbar structures at sublithographic scales and open the way towards enormous device density improvements over CMOS. In the literature, different active elements have been envisaged at the cross points, such as  $p-n$  junctions [153], molecular programmable switches/diodes [154] or FETs [155]. Several architectures have been proposed, based on these different devices. In [154] molecular diode-switches at crossbar intersections are proposed. This structure forms a diode logic grid, as in a Programmable Logic Array. However, diode logic requires level-restoring circuitry and addressing of individual points, which leads to complex interfaces. Considering these limitations, solutions based on FET logic, as introduced in [156], are explored. The NASIC extends this concept in [157] by creating double-stage combinational logic on Crossbars of NanoWires FETs (CB-NWFETs). The active element is therefore a FET.

### A.5.1.2 *Technological Assumptions*

The crossbar can be manufactured with silicon nanowires grown by Chemical Vapor Deposition, using metallic nanoparticles as catalyst and the Vapor-Liquid-Solid mechanism [158]. This technique allows the achievement of nanowires with diameters controlled by the catalyst size [159] and diameter values around 3 nm [160]. *In situ* doping during nanowire growth can be used to obtain n-type or p-type nanowires [161]. Nanowires can subsequently be thermally oxidized to obtain a core-shell structure [162] and deposited on a substrate (Figure a-i-a). The well-known Langmuir-Blodgett technique [163] can be used to align the nanowires and to obtain a first layer of nanowires with spacing controlled by the oxide thickness of the shell [162]. For better alignment when nanowires have very small diameters (i.e. below 3 nm), their lengths cannot be longer than a few micrometers to avoid gradual bending observed on nanowires with these diameters [160]. Nanowire pitch has also to be equal to (or greater than) the corresponding photolithography half pitch of a given technology node plus the nanowire radius for subsequent photolithography steps. The oxide shell should be removed except for along the bottom contact between the nanowires and the substrate, in order to avoid removing the nanowires at the same time (Figure a-i-b). Salicidation of some parts of the nanowires can be achieved using nickel (or platinum) physical vapour deposition, photolithography and etching to define the regions where the silicon nanowire will be the channel of transistors, annealing to form NiSi (or PtSi) regions and chemical removal of the unreacted deposited metal [155] (Figure a-i-c). The second array of nanowires (crossing the first) can be obtained with the same process sequence (Figure a-i-d). The oxide shell along the bottom of the second layer of nanowires serves as the gate dielectric and the NiSi (or PtSi) region serves as the metallic gate in the MOSFET structure (Figure a-i-e). In the same way, some nanowire regions of the second layer can serve as the channel of the MOSFET, controlled by the nanowire metallic regions of the first layer. The process ends with a final metallization to contact all nanowires around the crossbar area by a conventional sequence of deposition, photolithography and etching steps (Figure a-i-f).



**Figure a-i.** Manufacturing process to build a nanowire crossbar. *a)* Deposition of silicon nanowire after growth, oxidation and Langmuir Blodgett alignment. *b)* Etching of the oxide shell. *c)* Salicidation of the nanowire regions which will not be transistor channels. *d)* Deposition, alignment and etching of the oxide shell of the second layer. *e)* Salicidation as step *c)*. *f)* metallization to contact the nanowires around the crossbar.

## A.5.2 Lithographic Crossbar Integration

### A.5.2.1 Brief Literature Overview

Previous research works addressing the crossbar architecture generally considered bottom-up nanowire fabrication techniques. The low maturity of the technology implies several limitations on the assessment of the architecture. On the one hand, bottom-up fabrication techniques yield dispersed nanowires with highly variable position and spacing. Thus, mean values with respect to the geometry have to be assumed in order to characterize the architecture. On the other hand, the ability to define two layers of perpendicular nanowires is a very challenging technological task. To date, it has been demonstrated only for metallic nanowires in a top-down approach [170], while silicon-based nanowire crossbars are dispersion-based, micrometer-scale demonstrators without any functionalization of the cross points [166]. Capacitive coupling between parallel nanowires has been addressed only in routing applications [169].

In this proposal, we will assess the questions related to the technology limitations and to the impact of the coupling effects between semi-conducting nanowires. Our approach is based on the choice of a fully characterized industrial technology. We propose a process flow to fabricate lithography-based crossbars with FETs as active devices at the cross points. Then, we electrically simulate the circuit in order to assess its performance and the potential loss due to the capacitive coupling and the resistance of the long nanowires.

### A.5.2.2 Technological Assumptions

Using a *Fully-Depleted Silicon-on-Insulator* (FDSOI) process, wires can be manufactured with ultra-regular lines as demonstrated in [171]. In this section, we conceptually complete the already established process for a single layer of parallel nanowires, with a perpendicular top layer of parallel nanowires. In the proposed process, the bottom-most nanowires are defined using photo-lithography at the lithographic pitch, where their dimensions can be controlled through oxidation and etching below the lithographic limit down to 15nm [172]. Thereby, both n- and p-type dopings are allowed. On the other hand, the topmost lines are defined as *polycrystalline silicon* (poly-Si) stripes at the lithographic scale. These two perpendicular layers of parallel lines form a crossbar whereby the intersections are called cross points. In such a crossbar, the top lines can electrostatically control the nanowires underneath at the cross points in a FET fashion, when the ladders are covered by a gate oxide. Moreover, the top nanowires can form an ohmic contact to those lying underneath when a via is defined at the cross point.

Figure a-j shows the associated process flow. A P-type SOI substrate is patterned by lithography to form parallel ridges that are subsequently etched into nanowires. *Plasma Doping* (PLAD) is used to softly define N-type wires [173]. Then, the nanowires are passivated in oxide and planarized (Figure a-j-a). Following this step, the passive regions, *i.e.*, the parts of the nanowires connecting every series FET, are defined by n- and p-type PLAD on the p- and n-type nanowires respectively (Figure a-j-b). It is worth noticing that the implantation step is performed softly because of the small dimensions of the nanowires, such that dopant migration is limited. Moreover, the nanowires are separated by an oxide, further limiting the diffusion of dopants. This reduces the requirements on spacing, which are generally included in the design rules. This allows the smallest lithographic dimensions for all operations of patterning and doping to be reached. Then, the gate stack is defined by depositing the gate insulator, followed by the poly-silicon gate deposition and etching steps (Figure a-j-c). The poly-silicon lines carrying the gates are defined with regular parallel lines. At this level, the active devices are defined and the east-west connections between them are established through the passive parts of the nanowires, operating as resistances.

The north-south connections are composed of the poly-silicon lines and require the definition of vias between them, as well as the nanowires underneath them. The vias are defined by etching the poly-silicon lines and filling them with metal. In order to decrease the resistance of the north-south poly-silicon lines and the passive parts of the east-west silicon NW lines, it is possible to sputter a thin layer of nickel (or platinum) over the whole structure, which diffuses into the silicon and poly-silicon and forms a low resistance silicon silicide (Figure a-j-d). For this reason, it is important to first etch the oxide covering the passive regions before sputtering the metal. The remaining metal after the diffusion can be removed by wet etching [156]. Finally, the contacts between the crossbar and the outer circuit are implemented through conventional metallization steps (Figure a-j-e).

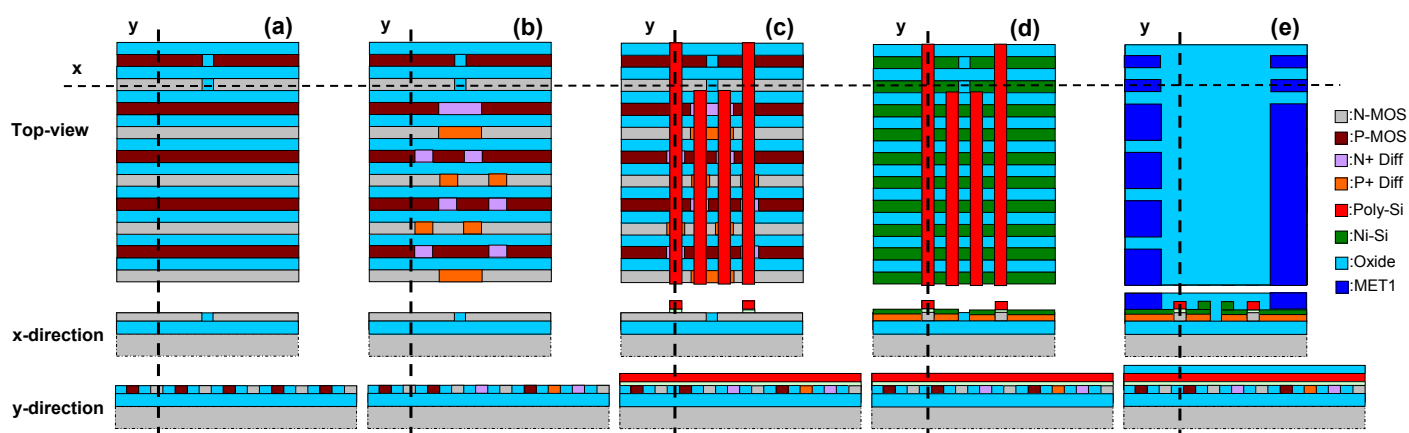


Figure a-j. Manufacturing process to build a FDSOI crossbar: a) grating patterning and active regions doping. b) Passive regions definition. c) Gate deposit. d) Passive regions finalization and salicidation. e) Metallization to contact passive regions.

## A.5.3 Vertical NanoWire FET Integration

### A.5.3.1 Brief Literature Overview

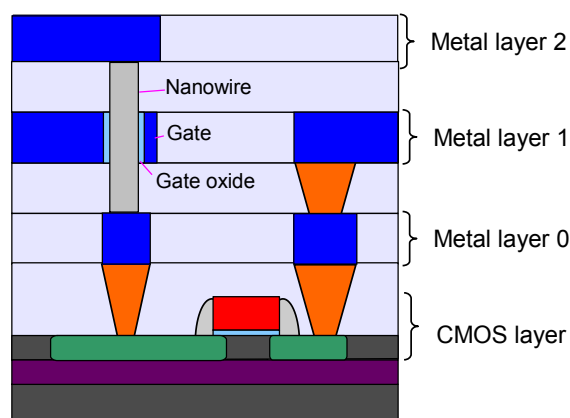
As previously discussed, the 3-D integration of transistors is an attractive solution to pursue the increase of circuit performance, while limiting the cost, as opposed to the continued single use of scaling. Stacking technologies, whether the traditional sequential technology (where wafers are processed separately and then stacked) or advanced monolithic integration (where transistors are processed step by step on the same wafer) only deal with stacks of planar transistors. The transistor itself (or more specifically, the transistor channel) does not exploit the vertical direction in these approaches.

Meanwhile, semiconducting nanowires have recently attracted considerable attention. To further miniaturize the transistor while still maintaining control over power consumption, alternative transistor geometries have been considered [106]. With their unique electrical and

optical properties, they offer interesting perspectives for basic research as well as for technology. A variety of technical applications, such as using nanowires as parts of sensors [109], and electronic [108, 171] and photonic devices [111] have already been demonstrated. In particular, electronic applications are increasingly coming into focus, as ongoing miniaturization in microelectronics demands new innovative solutions. Typically, silicon nanowire transistors have a horizontal, planar layout with either top or back gate geometry [110]. However, the amount of energy and time required to align and integrate these nanowire components into high-density planar circuits remains a significant hurdle for widespread application. More advanced works show a *Gate-All-Around* (GAA) organization in a planar topology [108]. In-place growth of vertically aligned nanowires, on the other hand, would in principle significantly reduce the processing and assembly costs of nanowire-based device fabrication, while opening up opportunities for "true" 3-D. Some research works have demonstrated the possibility of fabricating transistors directly between two metal lines, within the back-end levels [112, 113].

#### A.5.3.2 Technological Assumptions

Recent studies have demonstrated the possibility to grow single crystalline silicon nanowires on a metallic line, into a CMOS compatible process [129]. This work represents a great opportunity to build FET devices in the interconnect levels [112, 113]. We propose to co-integrate standard CMOS with vertical Nanowires Field Effect Transistors. The cross-sectional view is shown in figure a-k.



**Figure a-k. Cross sectional schematic showing a BEOL FET and standard CMOS FET co-integration**

First of all, standard transistors are processed using the specified technology, which could be very versatile, such as Bulk, Silicon-on-Insulator, Fully Depleted SOI, Thin Box FDSOI, among others. Then, silicon nanowires can be grown in a CVD reactor using the VLS mechanism. Even on a metallic line, they have a single crystalline structure and semiconducting properties. Taking advantage of this, and respecting low temperature processes under 400°C, it is possible to make vertical transistors between two interconnecting lines. After etching a hole through the oxide to the metallic bottom line, a catalyst can be deposited at the bottom. Nanowires can be grown from the metallic line using the oxide hole as template and a deposited metallic catalyst. Using diborane or phosphine, the nanowires can be doped to form P-N junctions. Nanowires for p-MOS and n-MOS should be grown during two distinct sequences comprising template formation and growth. After growth, a chemical etching can be used to remove a part of the oxide template. A multilayer gate stack can be achieved thanks to ALD and CVD deposit of the dielectric ( $\text{Al}_2\text{O}_3$ ,  $\text{HfO}_2$  ...), and the metal gate (TiN ...) respectively. An oxide can then be deposited before performing a CMP step on the top. Isotropic etching allows the removal of a part of the metal gate and defining the gate length. The space left by the metal gate can be filled by oxide deposition. The top contact is achieved by top line formation using a conventional damascene process.



## A.6 Overall Comparison

In order to compare the different technologies briefly with regard to the others, we use four global metrics: the cost (which globally includes the development and research costs, as well as the cost per device), the CMOS technological compatibility (i.e. the capabilities of existing equipments to be reused), the maturity (which translates the number of developments that are still required before industrialisation), and the operational reliability (including technological dispersion as well as the specific techniques that could be used to tolerate the defects). In this appendix, we are not seeking to compare performances, power consumption or area. Even if there is an obvious impact on the technology on these metrics, they are mostly influenced by the circuit and the architecture. Hence, they are only presented in the thesis work.

Table a-b presents the global results. Firstly, we stress that all the different metrics are closely related. Indeed, a mature process is generally reliable enough for its industrial use while the costs remain moderate. Hence, we should remark that the PCM technology which is based on back-end elements is less expensive than the technologies that have a direct influence on front-end. Lithography and depositions in back-layers are much simpler than in the front-end level. Technologies derived from standard technologies are globally interesting for all the different metrics. Hence, the 3-D Monolithic integration process and lithographic crossbars can be considered to be accessible from CMOS facilities and this could lead to significant cost lowering. Conversely, the technologies that address direct modifications to the front-end levels are the least mature. While 1-D structures such as DG-CNFETs and Crossbars of NWFETs are promising for the reduction of area used by circuits, they require several novel techniques compared to CMOS. The complex process flow required is also costly and unlikely to achieve industrial compatibility in the near term. Finally, it is interesting to make a special comment on the vertical NWFET technology. This technology is highly prospective and thus appears highly immature for near term. Nevertheless, the associated assumptions are compatible with low cost back-end process and we thus could expect that this solution could quickly achieve good maturity.

**Table a-b. Emerging technologies global metrics comparisons**

	<b>Reference publications</b>	<b>Cost</b>	<b>Maturity</b>	<b>Operational reliability</b>	<b>CMOS fab compatibility</b>
<i>PCM</i>	[70, 83, 73]	+++	+++	+	++
<i>DGCNFET-</i>	[142, 145]	-	-	+	-
<i>Monolithic 3-D</i>	[100, 101, 107]	+	++	++	+++
<i>Xbar NWFET</i>	[155, 162]	---	---	--	-
<i>Litho NWFET</i>	[171, 172]	++	+	+	+++
<i>Vertical NWFET</i>	[112, 113]	+	-	-	--

## **APPENDIX B** *Résumé Etendu*

---

B.1	Introduction .....	xxxv
B.2	Etat de l'art et Motivations.....	xxxvii
B.2.1	L'Architecture FPGA.....	xxxvii
B.2.1.1	Organisation du FPGA.....	xxxvii
B.2.1.2	Limitations Architecturales.....	xxxviii
B.2.2	Architectures Reconfigurables à base de Composants Emergents.....	xxxviii
B.2.3	Gabarit Architectural.....	xl
B.3	Structures Innovantes pour le Routage et la Configuration.....	xli
B.3.1	Proposition 1 : Technologie de Mémoires Résistives .....	xli
B.3.1.1	Hypothèses Technologiques .....	xli
B.3.1.2	Mémoire de Configuration.....	xli
B.3.1.3	Estimation des Performances.....	xlii
B.3.1.4	Circuits de Routage.....	xliii
B.3.2	Proposition 2: Procédés d'Intégration 3-D Monolithique .....	xliv
B.3.2.1	Hypothèses Technologiques .....	xliv
B.3.2.2	Tables de Correspondance .....	xliv
B.3.2.3	Éléments de Routage.....	xlvi
B.3.3	Proposition 3: Transistors à Nanofils Silicium Verticaux.....	xlvii
B.3.3.1	Hypothèses Technologiques .....	xlvii
B.3.3.2	Circuits Logiques Verticaux .....	xlviii
B.3.3.3	Evaluation des Performances et Méthodologie.....	xlviii
B.3.4	Conclusion.....	xlix
B.4	Impact Architectural des Circuits de Routage et de Configuration 3-D.....	li
B.4.1	Méthodologie et Outil d'Evaluation.....	li
B.4.2	Répercussion Architecturales des Mémoires à Changement de Phase .....	li
B.4.3	Répercussion Architecturale de l'Intégration 3-D Monolithique.....	lii
B.4.3.1	Etude de Surface .....	lii
B.4.3.2	Etude de Délai.....	liii
B.4.4	Répercussion Architecturale des Transistors à Nanofils Verticaux .....	liii
B.4.5	Conclusion.....	liv

B.5	Blocs Logiques en Ruptures .....	lv
B.5.1	Proposition 1: Electronique Carbone Ambipolaire .....	lv
B.5.1.1	Hypothèses Technologiques .....	lv
B.5.1.2	Logique Reconfigurable.....	lv
B.5.1.3	Amélioration des Circuits Logiques Dynamiques .....	lvii
B.5.2	Proposition 2: Structures Entrecroisées à base de Composants Silicium Monodimensionnels .....	lviii
B.5.2.1	Hypothèses Technologiques .....	lviii
B.5.2.2	Cellule Logique Reconfigurable à base de Nanofils.....	lix
B.5.3	Proposition 3: Structures Entrecroisées aux Dimensions Lithographiques .....	lx
B.5.3.1	Hypothèses Technologiques .....	lx
B.5.3.2	Méthode de Conception .....	lx
B.5.3.3	Evaluation des Performances dans le Cas Idéal.....	lxi
B.5.3.4	Evaluation des Performances dans le Cas Réel .....	lxii
B.5.4	Conclusion.....	lxiv
B.6	Propositions Architecturales en Ruptures et Analyse de Performances.....	lxv
B.6.1	Proposition Architecturale.....	lxv
B.6.1.1	MCluster .....	lxv
B.6.1.2	Organisation de la Hiérarchie BLE/CLB/FPGA.....	lxvi
B.6.2	Outils d'Evaluation .....	lxvii
B.6.2.1	Flot .....	lxvii
B.6.2.2	MPack .....	lxvii
B.6.3	Evaluation des Topologies d'Interconnexions .....	lxx
B.6.4	Evaluation de l'Architecture Complète.....	lxxi
B.6.5	Conclusion.....	lxxii
B.7	Conclusion.....	lxxiv
B.7.1	Méthodes et les Outils .....	lxxiv
B.7.2	Mémoires et Circuits de Routage pour FPGA .....	lxxiv
B.7.3	Blocs Logiques de Traitements pour FPGA .....	lxxv
B.7.4	Architectures .....	lxxv



## B.1 Introduction

Durant les quatre dernières décennies, l'industrie des semi-conducteurs a connu une croissance exponentielle en accord avec la loi de Moore. A mesure de l'approche vers le nanomètre, les promesses sont énormes grâce à des composants réduits à leurs limites physiques et économiques ultimes (Figure b-a). Néanmoins, l'*International Technology Roadmap for Semiconductors (ITRS)* a pointé de nombreux écueils aux niveaux des dispositifs tels que les courants de fuites, la puissance consommée et les effets quantiques. Ces limites poussent la recherche à explorer l'emploi de nouveaux matériaux et composants susceptibles de compléter, voire même de remplacer le transistor CMOS silicium utilisé dans les systèmes sur puces d'ici la fin de la décennie et ce, avant que les technologies silicium n'atteignent leur limite prévue autour de 2020 (c'est-à-dire lorsque la longueur physique du canal des transistors MOSFET descendra bien en dessous de la dizaine de nanomètres).

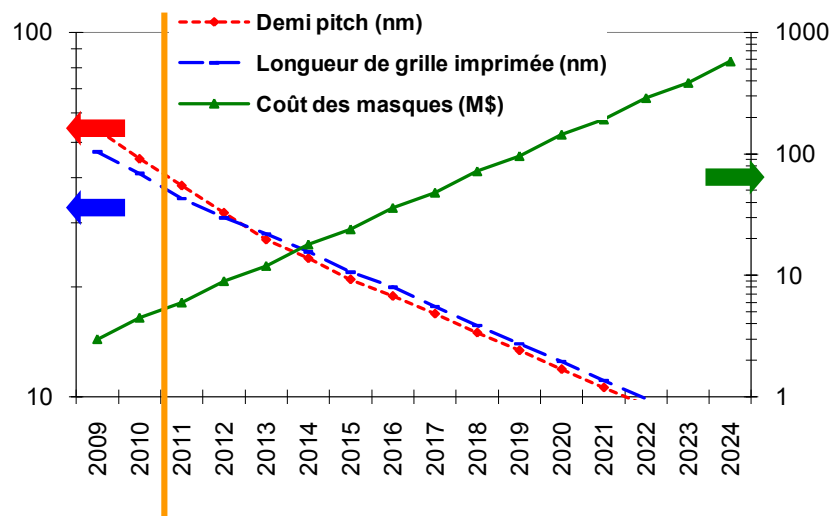


Figure b-a. Projection de l'évolution des dimensions et du coût des masques [6]

De nombreux composants issus des nanotechnologies sont actuellement à l'étude. La figure b-b présente la taxonomie associée aux nouvelles technologies. De nombreux phénomènes physiques peuvent être employés au travers de nombreux substrats, et ainsi réaliser des systèmes de traitement en rupture. Les systèmes émergents sont principalement formés par des structures unidimensionnelles ou quantiques. De manière générale, les recherches sont principalement centrées sur l'étude des nouveaux phénomènes physiques et des technologies de fabrication nécessaires à la réalisation à grande échelle de ces composants. Toutefois, un travail de recherche en avance de phase traitant de la conception de nouveaux circuits apparaît d'une importance stratégique capitale, afin de guider au mieux les développements technologiques vers des applications en rupture.

Cette large gamme technologique ainsi que les circuits élémentaires existants dans la littérature peuvent être utilisés pour construire de nouveaux blocs de base adaptés aux futures architectures de traitement. Néanmoins, l'utilisation efficace de dizaines de milliards de composants nanométriques présentant un fort taux de défauts va probablement aboutir à l'émergence de plate-forme reconfigurable en tant que principale fabrique de traitement. L'approche reconfigurable permet plus simplement une fabrication en volume, réduisant ainsi l'impact de l'évolution du coût des masques, tout en assurant la régularité nécessaire à l'amélioration de la robustesse.

Bien que l'utilisation de composants émergents apparaisse comme nécessaire pour aller au delà des limites de la technologie CMOS, les coûts exorbitants et la faible maturité des procédés de fabrication rendent difficile la recherche du meilleur candidat. En effet, l'évaluation d'une technologie nouvelle est à l'interface entre différents domaines tels que la technologie, l'architecture et les modèles.

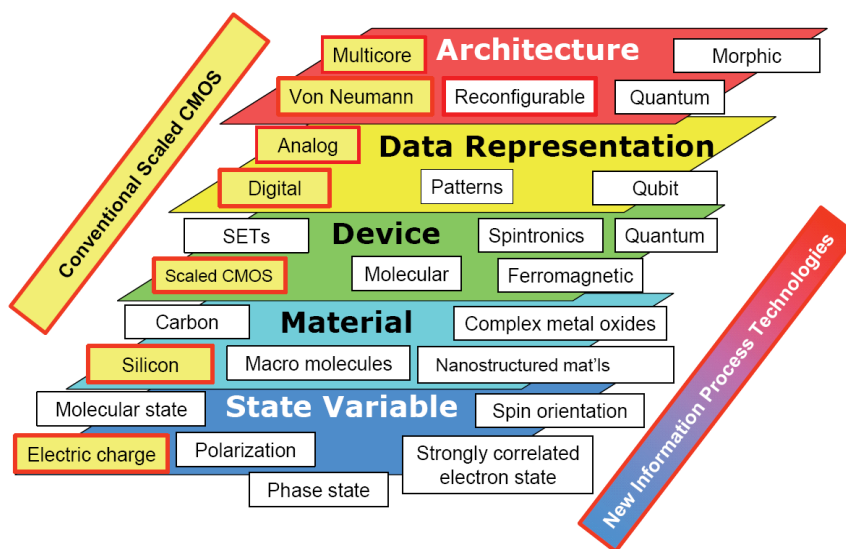


Figure b-b. Organisation hiérarchique des systèmes électroniques émergents [7]

Dans ce travail, nous proposons une méthodologie d'évaluation compatible avec un large éventail de technologies. La méthodologie est présentée dans la figure b-c. Les technologies seront évaluées au travers d'un modèle d'architecture générique. Le gabarit architectural est basé sur une architecture reconfigurable de type mer de porte. Selon la maturité, différents niveaux de modèles pourront être utilisés pour émuler la technologie. Les applications employées sont des circuits de tests standards, dont l'objectif est de simuler une large gamme de domaines applicatifs. Une boucle d'optimisation rapide est disponible au sein de la méthode. Elle permet d'identifier le meilleur compromis architectural pour une technologie donnée. Enfin, les résultats d'évaluation seront rétrocedés aux technologues, afin de guider les développements vers les caractéristiques technologiques les plus intéressantes du point de vue applicatif.

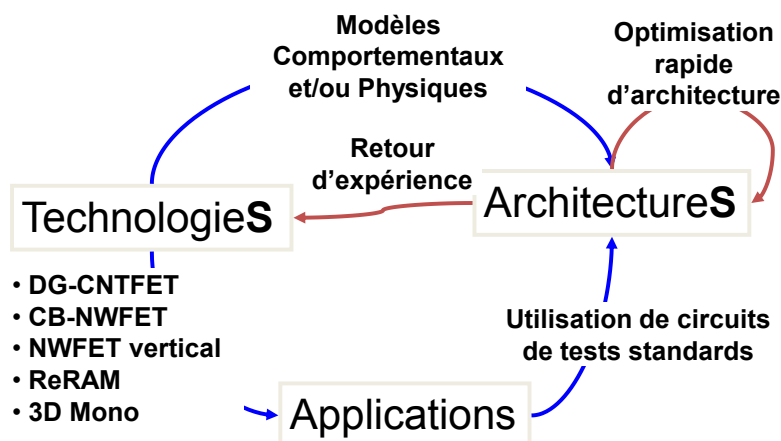


Figure b-c. Méthode d'évaluation technologique rapide basée sur l'utilisation d'un gabarit architectural

La thèse est organisée en cinq chapitres principaux, sans inclure l'introduction et la conclusion. Le chapitre B.2 présente un état de l'art relatif aux architectures reconfigurables conventionnelles, mais aussi à base de technologies émergentes. Il a aussi pour objet de définir le gabarit d'évaluation architectural. Le chapitre B.3 se concentre sur l'amélioration de structures de routage et de configuration pour les systèmes reconfigurables standards. Alors que ce chapitre centre son étude sur les circuits de base, le chapitre B.4 évalue l'impact de ces différentes améliorations du point de vue de l'architecture complète. Au sein du chapitre B.5, un travail est effectué sur les blocs logiques de base, afin d'obtenir de nouvelles graines pour l'organisation de l'architecture. Ces nouveaux blocs sont ensuite utilisés dans le chapitre B.6, où une nouvelle approche architecturale pour les systèmes reconfigurable est présentée et évaluée. Cette évaluation passe par la mise en place d'outils spécifiques. Enfin, la thèse est conclue et les contributions sont détaillées.

## B.2 Etat de l'art et Motivations

Comme introduit, les architectures logiques reconfigurables sont génériques et représentent un excellent compromis entre les coûts, le temps de développement et les performances. Dans ce chapitre, on se propose d'étudier les principales caractéristiques des *Field Programmable Gate Arrays* (FPGA) qui sont actuellement les composants configurables les plus répandus. Ensuite, nous aborderons les architectures reconfigurables utilisant des composants émergents, puis nous décrirons succinctement l'organisation choisie pour le gabarit d'architecture générique.

### B.2.1 L'Architecture FPGA

#### B.2.1.1 Organisation du FPGA

Un FPGA est un assemblage de nombreux blocs logiques programmables, capables d'implémenter de nombreuses fonctions logiques. Entre ces blocs logiques se trouve des connexions programmables permettant de relier les entrées et les sorties des blocs logiques entre elles. Cette structure FPGA simplifiée est illustrée par la figure b-d. Cette matrice symétrique de blocs logiques entourés par des canaux de routage est nommée architecture en îlots... Les canaux de routage sont connectés aux blocs logiques par des Boîtes de Connexions (CB – *Connection Box*) et ils sont interconnectés entre eux par des Boîtes de Brassage (SB – *Switch Box*). Les entrées et sorties du circuit sont quant à elles distribuées autour du périmètre du FPGA [17].

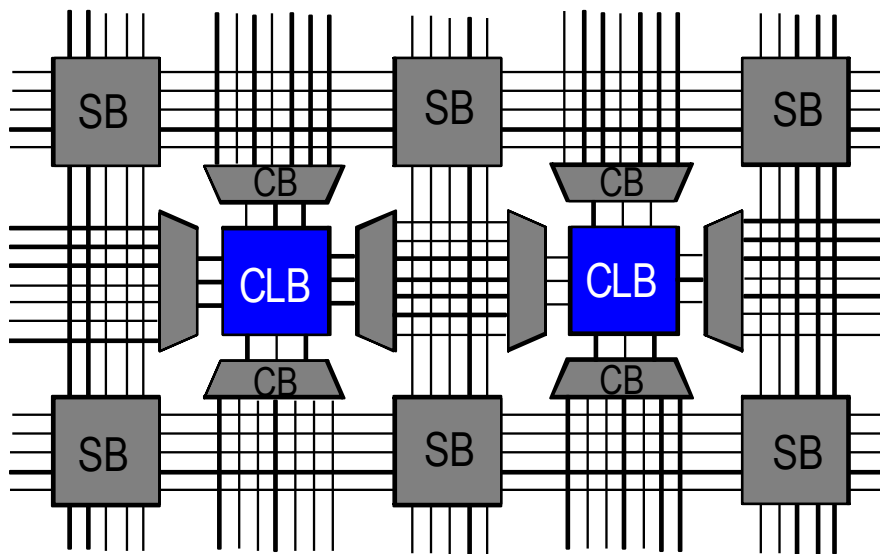


Figure b-d. Organisation d'un FPGA de type îlot

La figure b-e présente un groupement logique typique, appelé Bloc Logique Complexe (CLB – *Complex Logic Block*). Ce bloc consiste en un ou plusieurs Eléments Logiques de Base (BLE – *Basic Logic Elements*) groupés ensemble. Chaque BLE est composé d'une Table de Correspondance (LUT – *Look-Up Table*) et d'une bascule D qui permet de synchroniser le signal de sortie sur une horloge externe. Un multiplexeur permet la sélection entre le signal d'origine et le signal synchronisé. La structure d'un BLE est présentée dans la figure b-f. Les BLEs au sein du CLB sont entièrement connectés. Cela signifie que des multiplexeurs de routage permettent la liaison de n'importe quelle sortie et n'importe quelle entrée du CLB vers les BLEs. On peut noter que l'interconnexion complète interne au CLB simplifie grandement le routage physique par l'outil de conception (CAD – *Computer Aided Design*) et induit une réduction du nombre de canaux de routage utilisés. En revanche, elle nécessite un surcoût en termes de surface et de délai dû aux multiplexeurs internes utilisés pour le routage. Pour des groupements aux dimensions importantes, la hausse de surface et de délai peut être significative [17]. Chaque CLB contient  $N$  BLEs, qui peuvent être reliés par  $I$  entrées externes. Le BLE contient une LUT dont la taille est donnée par le nombre d'entrées  $K$ .

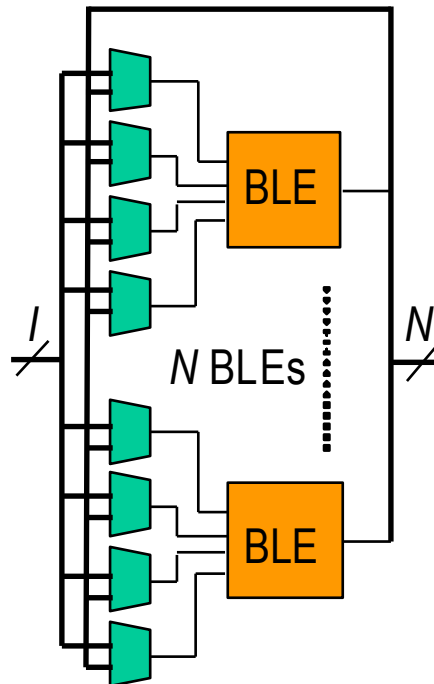


Figure b-e. Architecture d'un Block Logique Configurable [17]

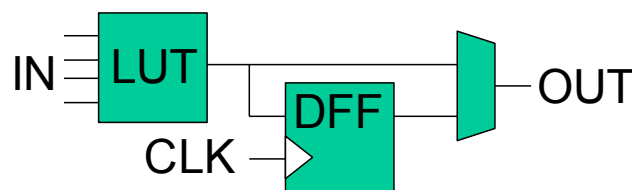


Figure b-f. Architecture d'un Élément Logique de Base [17]

### B.2.1.2 Limitations Architecturales

L'architecture FPGA souffre de nombreuses limitations. La figure b-g montre la répartition de la surface occupée par les différents blocs fonctionnels. On peut remarquer que les mémoires de configuration (pour la logique et le routage) occupent environ la moitié de la surface. Les blocs logiques occupent seulement 22% de la surface totale incluant leurs propres mémoires de configurations et seulement 14% de la surface est utilisée pour le calcul. En sus de l'occupation de surface, le routage programmable contribue significativement au délai total du FPGA (80% du délai total [21, 22]) et à la consommation dynamique (60% de la consommation totale [23, 24, 25]).

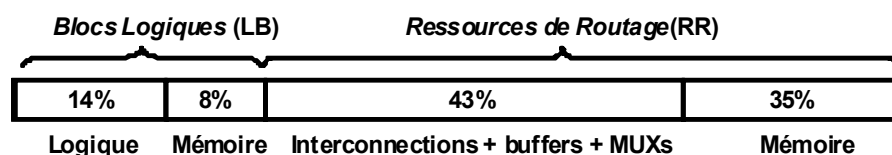


Figure b-g. Répartition de la surface d'un FPGA par rapport aux fonctions [38]

## B.2.2 Architectures Reconfigurables à base de Composants Emergents

La Feuille de Route Internationale pour la Recherche sur les Semi-conducteurs (ITRS – *International Technology Roadmap for Semiconductors* [7]) identifie les applications possibles des technologies émergentes. Deux familles principales sont alors définies: les architectures morphiques et les architectures hétérogènes.

Tel qu'indiqué par l'étymologie, l'approche morphique développe des architectures fortement en rupture par rapport aux circuits classiques. Souvent inspirées par la biologie, ces systèmes sont très performants en regard de leur faible consommation et de leur tolérance naturelle aux

défauts. L'apport de nouveaux composants et de nouveaux concepts d'intégration va certainement conduire à de nombreuses améliorations de ce type de circuit, et permet d'envisager de toutes nouvelles techniques de fabrication. Néanmoins, la transition vers une technologie en rupture ne peut arriver brusquement en raison des contraintes économiques mises sur les équipements standards. Ainsi, dans un futur proche, il semble raisonnable de penser que les circuits CMOS standards seront améliorés de façon incrémentale par les nouvelles technologies. Cette approche est appelée l'intégration hétérogène. Deux grandes pistes peuvent être tracées dans cette voie : une amélioration de la fonctionnalité intrinsèque des composants (c'est-à-dire que de nouvelles fonctionnalités peuvent être utilisées pour une surface identique) et une amélioration venant de la densité d'intégration (c'est-à-dire avec l'intégration de composants de plus en plus petits).

Les composants monodimensionnels améliorent la performance du canal des transistors. Néanmoins, l'obtention de dimensions ultra-compactes représente un réel challenge à l'intégration. En particulier, le manque de fiabilité de la photolithographie recommande l'utilisation de circuits hyper-réguliers.

Parmi les composants émergents, les nanofils semi-conducteurs apparaissent comme les plus prometteurs pour les systèmes denses [174]. Une application potentielle des nanofils est l'architecture *crossbar*, au sein de laquelle des fils parallèles sont organisés en deux couches perpendiculaires entrecroisées. Ces deux couches réalisent des points de croisement à leurs intersections, qui peuvent être fonctionnalisés pour réaliser une fonction. Ce type d'architecture promet des améliorations significatives en termes de surface, mais également une régularité intrinsèque comparée à la technologie CMOS [174]. Historiquement, la régularité est utilisée dans les circuits reconfigurables de type *Programmable Logic Array* (PLA). Dans les approches modernes, les crossbars sont utiles pour réaliser des circuits reconfigurables ultra-denses à base de composants nanométriques. Dans [154], un couple diode-interrupteur programmable moléculaire est employé aux points de croisements. Cette structure forme une grille de logique à diode, se comportant comme un PLA. Néanmoins, la logique à diode nécessite d'assurer la restauration des niveaux logiques, aboutissant à de complexes interfaces. Considérant les limitations de la logique à diode, des solutions employant des transistors sont explorées dans [156]. Le NASIC étend ce concept dans [157] en créant une logique combinatoire à deux étages à partir d'un crossbar de transistors à nanofils. Cette structure est employée comme une fabrique généraliste permettant d'implémenter des blocs logiques suffisamment complexes pour former le nanoprocesseur WISP-0 [157].

Toutes les architectures proposées précédemment emploient l'amélioration de la densité d'intégration autorisée par l'utilisation de procédés de fabrication sublithographiques tandis que la fonctionnalité originelle du transistor reste inchangée vis-à-vis du traditionnel MOS. De récentes percées technologiques du côté de l'électronique carbone ont montré de nouvelles propriétés à l'échelle du composant. En effet, le Transistor à Nanotube de Carbone Double-Grille (DG-CNFET – *Double Gate Carbon Nanotube Field Effect Transistor*), réalisé à partir d'un unique nanotube et contrôlé par deux grilles indépendantes présente une propriété d'ambipolarité [142]. Les deux grilles permettent de contrôler les caractéristiques du transport dans le canal du transistor. L'ambipolarité signifie, dans un contexte double-grille, que les comportements d'un type  $n$  et d'un type  $p$  peuvent être observés au sein du même composant dépendant simplement de la tension appliquée sur la grille arrière. Dans [51, 52], cette propriété est utilisée pour concevoir des fonctions logiques très compactes utilisant à la fois les grilles avant et arrières des DG-CNFET. Dans [122], une cellule logique reconfigurable dynamiquement et utilisant seulement 7 transistors pour réaliser 14 fonctions Booléennes de façon reconfigurable est présentée.

Afin de comparer les différentes architectures introduites précédemment, cinq métriques sont employées : la surface, la performance, la puissance consommée, la maturité technologique et la capacité de tolérance aux fautes. Le tableau b-a présente les résultats de l'estimation. On

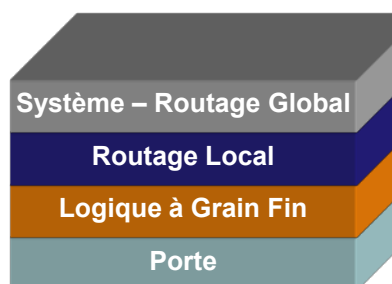
peut remarquer que les solutions utilisant des nanofils ultra-dense améliorent radicalement la surface du circuit final. Ceci vient de l'augmentation de la densité d'intégration de chaque composant élémentaire. Néanmoins, ces solutions sont plus discutables au regard des autres métriques. En fait, les solutions à nanofils entrecroisés nécessitent souvent l'utilisation de signaux d'horloge complexe afin de restaurer les niveaux logiques. Cela impacte évidemment la performance et la consommation du fait de ces besoins « externes ». Dans ce contexte de puissance consommée, on peut noter que l'électronique carbone est prometteuse du fait de ses bonnes propriétés intrinsèques. Grâce à l'électronique carbone, il est possible de bâtir des solutions qui globalement améliorent toutes les métriques jusqu'à la maturité et la tolérance aux défauts. En effet, de nombreux travaux ont été menés afin d'améliorer les procédés de fabrication associés et de rendre la technologie tolérante à ses imperfections [55, 56].

**Tableau b-a. Comparaison globale des architectures reconfigurables émergentes**

	<b>Publications de références</b>	<b>Surface</b>	<b>Performances</b>	<b>Puissance consommée</b>	<b>Maturité</b>	<b>Tolérance aux fautes</b>
<i>NanoFabric</i>	[47, 48]	++	-	-	-	-
<i>NanoPLA</i>	[154]	++	-	-	-	--
<i>NASIC</i>	[64]	++	+	-	-	-
<i>DG-CNFET</i>	[52, 122]	+	+	++	-	+

### B.2.3 Gabarit Architectural

Dans ce travail de thèse, le point est fait autour des architectures reconfigurables et de leur amélioration par les technologies émergentes. On se propose ainsi d'utiliser un gabarit d'architecture générique. Le gabarit est présenté dans la figure b-h. Le gabarit est organisé en quatre niveaux : le Niveau Porte, le Niveau Logique à Grain Fin, le Routage Local et le Routage Global. Le Niveau Porte correspond aux portes élémentaires réalisant les opérations combinatoires et séquentielles. Dans le cas d'un FPGA, il s'agit des LUTs, des Bascules et des Multiplexeurs. Le Niveau Logique à Grain Fin correspond au bloc élémentaire le plus simple capable d'effectuer une opération complète. Ce bloc contient sa propre mémoire de configuration. Dans le cas du FPGA, ce niveau est réalisé par un BLE. Le Routage Local permet l'arrangement de la Logique à Grain Fin, ainsi que son interconnexion totale. Il a pour but de réaliser des blocs logiques avec une granularité plus grande, réalisant ainsi des opérations plus complexes. Il s'agit du CLB dans le contexte FPGA. Enfin, au niveau du Routage Global, c'est-à-dire du système. Le circuit est organisé avec une macro-régularité. Les blocs logiques définis précédemment sont arrangés de façon régulière et interconnectés de manière programmable. C'est typiquement le niveau de la structure en îlot caractéristique des FPGAs.



**Figure b-h. Organisation par niveau du gabarit d'architecture générique**

## B.3 Structures Innovantes pour le Routage et la Configuration

Dans le chapitre précédent, nous avons pu constater qu'environ 45% de la surface d'un FPGA est utilisée par les mémoires de configurations, et que 78% de la surface est occupée par les ressources de routage. Nous allons, au sein de ce chapitre, utiliser des technologies d'intégration 3-D, afin de placer des composants dans les niveaux supérieurs des circuits. Trois différentes technologies seront employées : les Mémoires à Changement de Phase (PCM – *Phase Change Memory*), l'intégration 3-D monolithique et l'utilisation de Transistors à Nanofils (NWFET – *NanoWire Field Effect Transistors*) verticaux.

### B.3.1 Proposition 1 : Technologie de Mémoires Résistives

Cette première proposition a pour objectif de placer des composants mémoires passifs dans les premiers niveaux métalliques d'un circuit microélectronique. On envisage ainsi de remplacer les mémoires de configuration par des circuits mémoires élémentaires placés directement au-dessus des portes en ayant besoin, mais également de réaliser des points de routage performants.

#### B.3.1.1 Hypothèses Technologiques

Une mémoire à changement de phase est basée sur les propriétés des alliages à base de chalcogénures, utilisés en tant que matériaux actifs dans le point mémoire. Un alliage de chalcogénure est un verre semi-conducteur réalisé à partir d'éléments du groupe VI de la table des éléments, tels que le sulfure, le sélénium ou le tellure. Ces verres présentent la capacité d'une transformation de phase réversible. En effet, par le biais d'un contrôle précis de la température (par effet Joule) dans la cellule mémoire, il est possible de changer l'état du matériau entre deux configurations stables, soit un état polycristallin hautement conducteur soit un état amorphe faiblement conducteur.

La technologie PCM est compatible avec la technologie CMOS. En effet, la fabrication de l'élément mémoire peut être faite juste après la réalisation des transistors standards sous-jacents. Il est ainsi possible de réaliser le point mémoire dès les premières étapes métalliques ou après dans les niveaux supérieurs [75]. Une vue en coupe d'un élément mémoire est présentée par la figure b-i. Sur cette figure, la PCM est constituée d'une couche de matériaux à changement de phase (PC) avec deux électrodes de contacts (BEC et TEC). Le matériau est intégré entre les interconnexions métalliques M0 et M1. Le transistor d'accès est réalisé sous le point mémoire par des techniques conventionnelles.

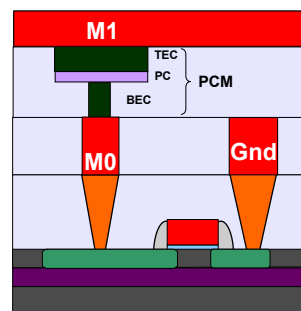


Figure b-i. Vue en coupe d'intégration de composant PCM.

#### B.3.1.2 Mémoire de Configuration

##### i) Description

L'élément mémoire est présenté dans la figure b-j. Le circuit consiste en 2 mémoires résistives connectées selon une structure de diviseur de tension entre 2 lignes de potentiels fixes. Un transistor est également présent entre la masse et le nœud de sortie de la cellule. Il est utilisé pour sélectionner le nœud durant la phase de programmation. La sortie se fait directement en tension, et peut piloter nativement une entrée de porte logique. L'opération de



lecture est intrinsèque avec cette structure, tandis que la programmation est une opération externe.

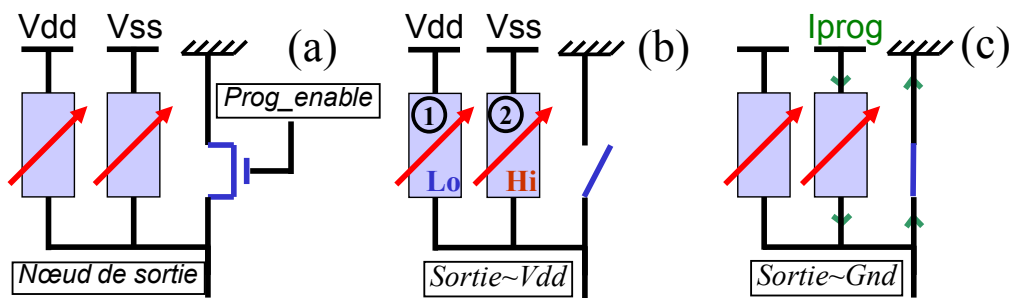


Figure b-j. (a) Point mémoire élémentaire, (b) en lecture et (c) en écriture

### ii) Phase de Lecture

Un arrangement en diviseur de tension est utilisé dans ce circuit pour réaliser intrinsèquement la conversion de la donnée stockée sous forme de résistance en une tension. La figure b-j-b présente un exemple de configuration, où le nœud stocke la valeur logique « 1 ». Le transistor de programmation est dans l'état *off* (signal *Prog\_enable* inactif), déconnectant ainsi la masse de la sortie. La mémoire résistive (1), connectée à la ligne  $V_{dd}$ , est configurée dans son état cristallin, de sorte que sa résistivité soit basse (quelques  $k\Omega$ ). L'autre mémoire (2), connectée à  $V_{ss}$ , est dans son état amorphe de haute résistivité (proche du  $M\Omega$ ). Par conséquent, le nœud de sortie est chargé au potentiel proche par la branche de basse résistivité. Les niveaux logiques dépendent des valeurs de  $R_{ON}$  et de  $R_{OFF}$  selon les relations suivantes:

$$'1' = V_{dd} - \frac{R_{ON}}{R_{ON} + R_{OFF}} (V_{dd} - V_{ss}) \quad '0' = \frac{R_{ON}}{R_{ON} + R_{OFF}} (V_{dd} - V_{ss})$$

Il est également bon de noter que dans le cadre d'une lecture continue, un courant s'établit au travers des résistances. Ceci conduit à une consommation passive au travers de la structure, dépendant de la relation suivante:

$$I = \frac{V_{dd} - V_{ss}}{R_{ON} + R_{OFF}} \approx \frac{V_{dd} - V_{ss}}{R_{OFF}}$$

Ce courant statique peut être réduit par un choix technologique, utilisant des matériaux maximisant la valeur de  $R_{OFF}$ .

### iii) Phrase d'écriture

La figure b-j-c illustre la phase de programmation du nœud. Tandis que le transistor de sélection est placé dans son état conducteur (signal *Prog\_enable* actif), les tensions fixes de lecture sont déconnectées des rails supérieurs et remplacées par l'unité de programmation. Par la suite, les courants de programmation sont appliqués séquentiellement au travers des résistances afin de commuter leur état. Ces courants de programmation s'écoulent vers la masse. Chaque cellule disposant de son transistor de sélection, les lignes de programmation peuvent être partagées de la même façon que dans les architectures de mémoires conventionnelles (Figure b-k).

#### B.3.1.3 *Estimation des Performances*

Le tableau b-b recense les résultats de caractérisation en termes de surface, temps d'écriture et énergie de programmation. Ces métriques ont été évaluées pour la solution proposée et pour des points mémoires traditionnels de FPGAs. On peut constater que la solution à base de PCM est la plus compacte avec un gain d'un facteur 1,5, malgré l'impact du large transistor de programmation. Ce gain est dû à la réduction de la surface du point mémoire (en projection sur la surface active) à seulement un transistor, en comparaison de 5 pour la cellule SRAM et de 2 pour la cellule flash. Il est également important de remarquer que les mémoires PCM offrent une intéressante réduction du temps de programmation au détriment d'une énergie de programmation plus importante.



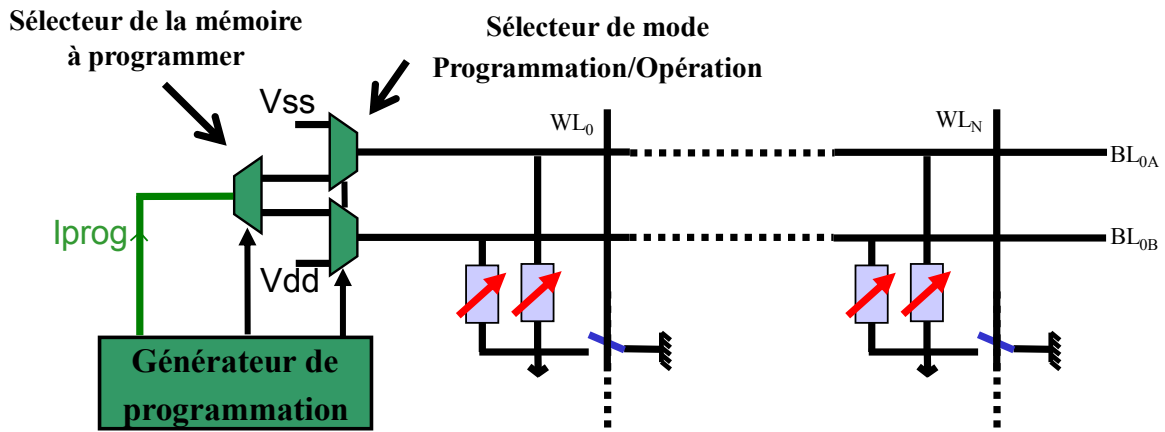


Figure b-k. Illustration du partage des lignes d'alimentations

Tableau b-b. Comparaison des technologies pour les mémoires de configuration

	Topologie de la cellule	Surface (F <sup>2</sup> )	Temps d'écriture (ns) [6]	Energie de programmation (pJ) [6]
SRAM	5T	115	0,3	7.10 <sup>-4</sup>
Flash	2T	46	1000	0,1
PCM	1T2R	30	60	12
Flash vs. PCM	-	x 1,5	x 16,6	x 0,08

B.3.1.4 Circuits de Routage

i) Description

Comme abordé précédemment, les blocs le plus influents dans une structure de routage FPGA sont les boîtes de brassage. Celles-ci assurent en effet les interconnexions entre les canaux de routage de façon programmable. On se propose d'améliorer leur structure en remplaçant les habituelles portes à transmission par les structures illustrées dans la figure b-l. La figure b-l-a représente le diagramme fonctionnel d'une boîte de brassage 2x2 voies à base de PCMs. A chaque point de croisement, des éléments de routage configurables sont disposés. Ils permettent de réaliser n'importe quel chemin entre les terminaux Nord/Sud/Est/Ouest. La structure des points de croisement est similaire au CMOS, et est basée sur des mémoires PCMs à 2 terminaux qui remplacent les transistors de passage (Figure b-l-b). Ainsi, une connexion à l'état passant est réalisée en programmant la mémoire associée (c'est-à-dire celle connectant les deux fils à relier) dans son état de faible résistivité.

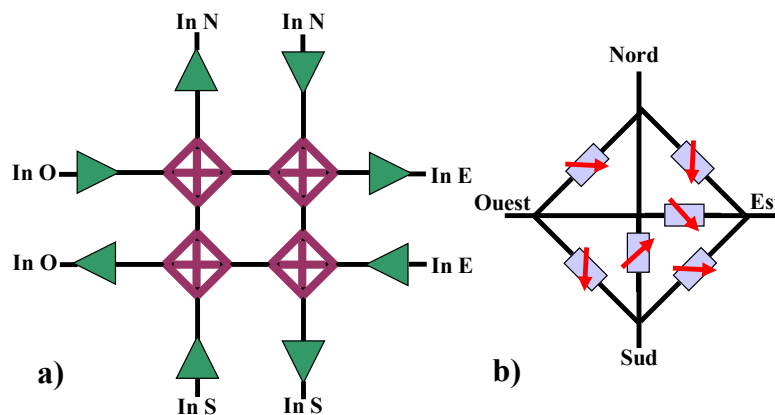


Figure b-l. Architecture d'une boîte de brassage 2x2 voies à base de PCM

Une PCM est programmée par application d'un pic de courant entre ses deux terminaux. Ainsi, il est nécessaire d'adresser chaque mémoire séquentiellement. La sélection du nœud est faite au travers des interfaces d'entrée/sorties (Figure b-l-a). Ces interfaces permettent de commuter entre les canaux de routage et l'unité de programmation, mais également de laisser

certaines terminaux flottant (afin de supprimer les programmations parasites). Une fois la sélection, l'unité de programmation applique le signal adéquat pour programmer l'état désiré. Un exemple de programmation séquentielle est démontré dans la figure b-m.

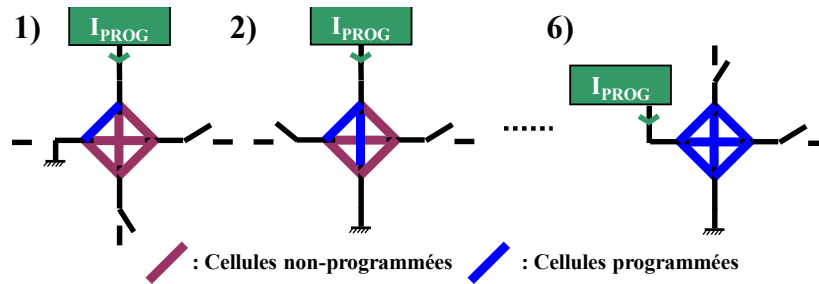


Figure b-m. Exemple d'une séquence de programmation d'un point d'intersection

## ii) Evaluation des Performances

Le tableau b-c présente les résultats de caractérisation en termes de surface, temps d'écriture et résistance du chemin de donnée pour les circuits proposés et pour les technologies FPGA traditionnelles. Une fois de plus, la technologie PCM est la plus compacte avec un gain d'un facteur 3,4. De plus, on se doit remarquer que les résistances du chemin de données sont décalées vers des valeurs faibles. En effet, tout en conservant un rapport *on/off* convenable, la résistance de *on* passe sous la centaine d'Ohms. Cette propriété est particulièrement intéressante car elle permet d'améliorer les performances de l'interrupteur réalisé avec cette technologie.

Tableau b-c. Comparaison des technologies pour les circuits de routage (Brassage 2x2 voies)

	Topologie de la cellule	Surface (F <sup>2</sup> )	Temps d'écriture (ns) [6]	Résistance On (Ω)	Résistance Off (MΩ)
<i>SRAM</i>	72T	2576	2,4	9100	144
<i>Flash</i>	48T	1104	24000	9100	144
<i>PCM</i>	8T24R	321	1440	50	1
<i>Flash vs. PCM</i>	-	X 3,4	x 16,6	x 182	x 0,007

## B.3.2 Proposition 2: Procédés d'Intégration 3-D Monolithique

Au travers de la technologie précédente, nous avons placé un composant passif dans les niveaux métalliques. Il apparaît pertinent d'étudier à présent une technologie permettant de placer des éléments actifs de façon 3-D.

### B.3.2.1 Hypothèses Technologiques

L'intégration 3-D monolithique permet d'intégrer plusieurs couches de silicium actifs les unes sur les autres, tout en préservant une grande qualité d'alignement et une grande densité d'interconnexion. En effet, contrairement aux techniques 3-D conventionnelles, les couches actives sont traitées de façon séquentielle, comme montré par la figure b-n.

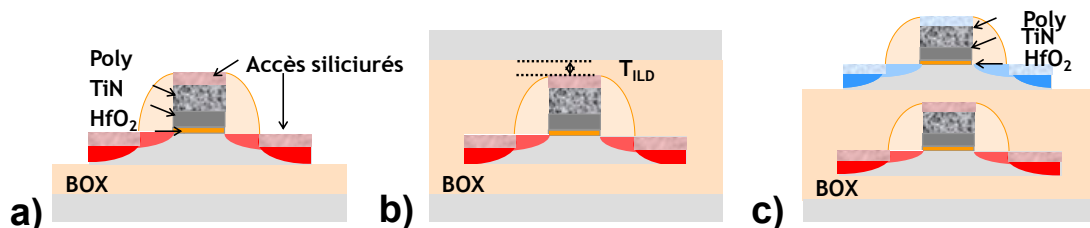


Figure b-n. Vue en coupe des étapes d'intégration 3-D monolithique – a) Couche FDSOI optimisée basse température b) Dépôt d'une couche haute de silicium haute qualité c) Couche FDSOI haute structurée par un procédé à basse température

Les principales étapes du procédé sont (a) la réalisation des transistors bas, (b) le dépôt d'un film mince de silicium sur les transistors existants et (c) la réalisation des transistors haut. La

première couche de transistor est identique à tous procédés conventionnels tels que le FDSOI (*Fully Depleted Silicon-On-Insulator*). Néanmoins, ces transistors doivent être robustes en regard du budget thermique. En effet, il ne faut pas que les performances de ces transistors soient dégradées par la réalisation des transistors supérieurs. Ainsi, une attention particulière est portée sur la siliciuration des contacts avec un procédé optimisé W-NiPtSi-F [97]. Le dépôt de la couche mince se fait par report de plaque. Les transistors hauts sont enfin réalisés avec un budget thermique maximal de 600°C, ce qui réclame l'emploi de techniques particulières d'activation des dopants [97].

### B.3.2.2 Tables de Correspondance

La technologie d'intégration 3-D monolithique permet la répartition des blocs logiques à un niveau très fin. En effet, du fait de la grande densité de contacts disponibles, il est possible de répartir verticalement les composants jusqu'au niveau transistor. Dans notre approche, on propose une séparation au niveau porte. Cela va nous permettre de séparer la mémoire de configuration de la logique intégrée au chemin de données. Ce type de séparation est extrêmement intéressant puisqu'il permet de réaliser une optimisation technologique spécifique à la logique haute performance sur une couche et aux mémoires basse consommation sur l'autre couche. On peut noter dès à présent que ce genre de répartition n'est pas possible avec les technologies 3-D standards, puisque la densité des contacts ne permet pas de réaliser les diverses interconnexions.

La séparation des éléments appartenant au chemin de données et des éléments appartenant à la configuration ont un impact immédiat sur l'architecture d'une LUT. La figure b-o présente le schéma ainsi que le dessin d'une LUT 1-bit. Bien que cette structure soit trop simple pour une utilisation fonctionnelle, elle permet la démonstration technologique. Deux points mémoires de configuration sont ainsi placés au dessus d'un circuit multiplexeur.

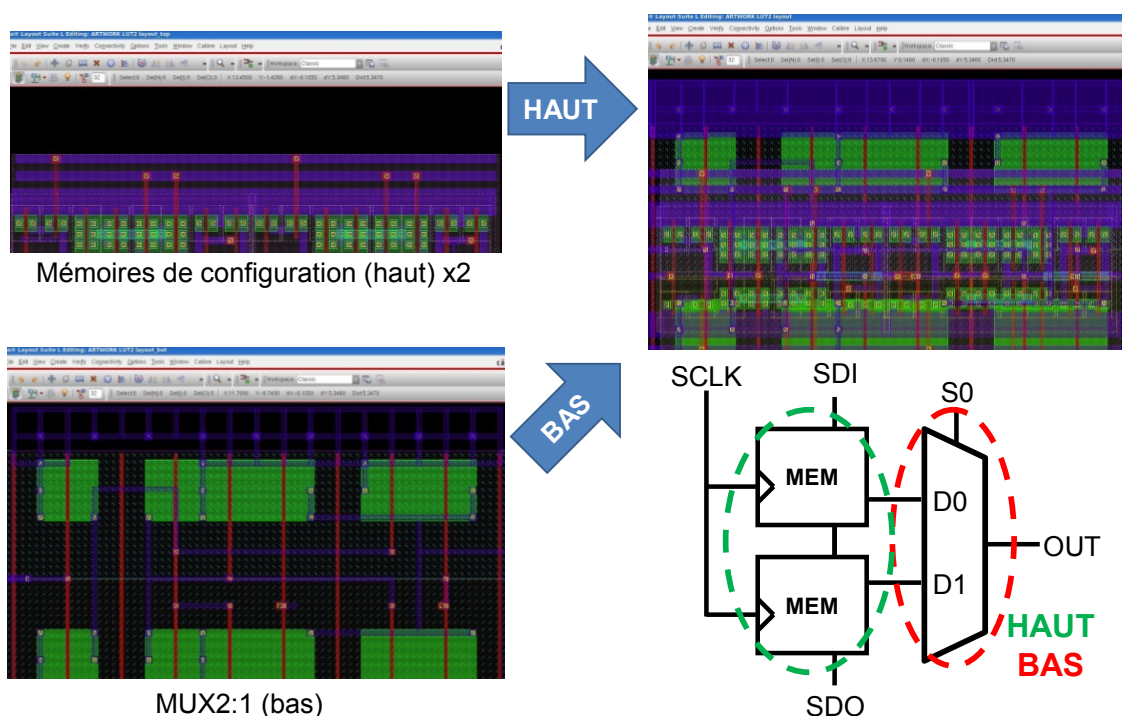


Figure b-o. Table de correspondance 1-bit réalisée sur 2 niveaux par intégration 3-D monolithique – dessin et schéma

Le tableau b-d montre les performances d'une LUT à 2 entrées. Concernant la surface, on constate une amélioration d'un facteur 2 en regard du bulk 2-D. Ce gain vient de la superposition de la mémoire sur le multiplexeur de données. A propos des performances et de la consommation, on observe une réduction du délai intrinsèque d'un facteur 1,62 et de l'influence de la charge d'un facteur 6,1 tandis que la puissance dynamique est réduite d'un facteur 2 à 1GHz. Ces intéressantes performances viennent de deux phénomènes. Effectivement, la technologie utilisée (ici le FDSOI) lors de l'intégration est une technologie

adaptée pour la réalisation de circuits à la fois haute performance et basse consommation. Néanmoins, l'intégration 3-D a aussi un impact non négligeable, car elle permet l'optimisation technologique séparée des différentes couches. Ainsi, les mémoires placées sur la couche haute pourront bénéficier d'améliorations technologiques spécifiques afin de réduire encore la consommation.

Tableau b-d. Evaluation des performances d'une LUT 3-D

LUT2 Nœud 65-nm	Surface ( $\mu\text{m}^2$ )	Délai intrinsèque (ps)	$K_{\text{Load}}$ ( $\text{ps}\cdot\text{fF}^{-1}$ )	Puissance moyenne à 1 GHz ( $\mu\text{W}$ )
2-D LP Bulk	64,48	85,08	36,29	46,90
3-D FDSOI	31,7	52,52	5,94	23,23
2-D Bulk vs. 3-D FDSOI	x 2,03	x 1,62	x 6,11	x 2,02

### B.3.2.3 Éléments de Routage

Dans la même veine, nous avons proposé l'intégration 3-D d'un point de routage programmable. Ce point de routage est réalisé par une porte à transmission pilotée par une mémoire statique, comme illustré par la figure b-p. Cette figure implémente une classique cellule SRAM 6T au dessus de la porte à transmission du chemin de données.

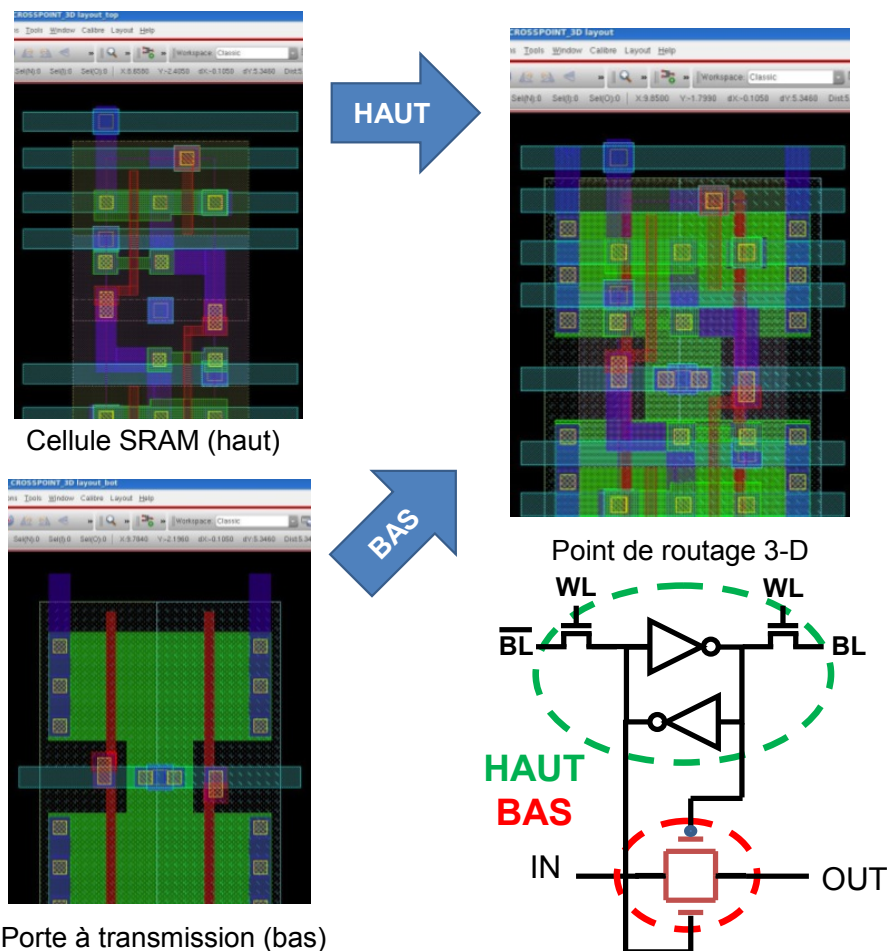


Figure b-p. Point de routage 3-D (base sur une cellule SRAM) – dessin et schéma

Le tableau b-e présente les performances du point de routage proposé. Concernant la surface, on note que ce point de routage améliore l'équivalent bidimensionnel par un facteur 2,93. En effet, contrairement à la LUT où de nombreuses connexions sont nécessaires entre les niveaux métalliques et les niveaux actifs, le point de routage est plus économe en interconnexions ce qui fait que l'utilisation de la 3-D peut être optimisée finement. En termes de performances, le délai intrinsèque est réduit d'un facteur 3,1. Les performances proviennent de la porte à

transmission. La porte à transmission est moins complexe que la cellule mémoire et donc plus compacte. Cette porte étant placée sous la cellule mémoire, il est possible d'augmenter ses dimensions sans pour autant augmenter la taille de l'ensemble. Ainsi, il est possible de dessiner de larges portes à transmission tout en conservant un gain en surface important. Les larges transistors des portes à transmission sont alors directement responsables des bonnes performances en délai. Finalement en termes de puissance dynamique, le point de routage améliore l'équivalent bulk d'un facteur 2 à 2GHz. De la même façon que précédemment, ceci s'explique par la basse consommation de la technologie FDSOI.

Tableau b-e. Evaluation des performances du point de routage

LUT Nœud 65-nm	Surface ( $\mu\text{m}^2$ )	Délai intrinsèque (ps)	$K_{\text{Load}}$ ( $\text{ps}\cdot\text{fF}^{-1}$ )	Puissance moyenne à 2 GHz ( $\mu\text{W}$ )
2-D LP Bulk	15,83	0,44	0,71	0,12
3-D FDSOI SRAM	5,40	0,14	0,66	0,058
2-D Bulk vs. 3-D FDSOI SRAM	x 2,93	x 3,1	x 1,07	x 2,1

### B.3.3 Proposition 3: Transistors à Nanofils Silicium Verticaux

La technologie précédente a permis l'intégration en 3-D de composants actifs. Néanmoins, on peut se questionner sémantiquement sur l'appellation 3-D pour un empilement de composants 2-D. Ainsi, on se propose d'étudier l'intérêt d'une technologie permettant la « vraie » 3-D. En effet, la réalisation de transistors verticaux directement entre deux niveaux métalliques permet de créer des transistors réellement dans la 3ème dimension.

#### B.3.3.1 Hypothèses Technologiques

De récentes études ont montré qu'il était possible de faire croître un nanofil de silicium cristallin entre 2 lignes métalliques, tout en restant compatible avec les procédés CMOS conventionnels [129]. Ce travail représente une grande opportunité pour réaliser des transistors au sein des niveaux métalliques d'interconnexion [112, 113]. Dans cette proposition, on se propose de co-intégrer cette technologie avec les techniques siliciums standards.

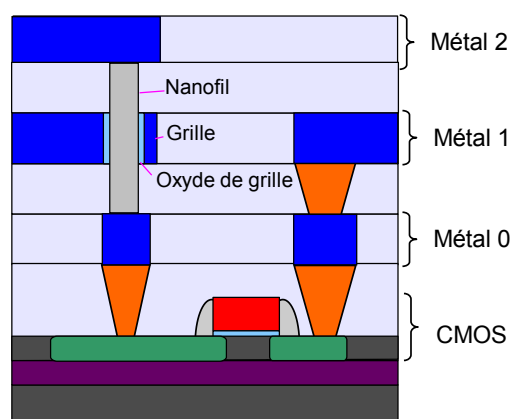


Figure b-q. Vue en coupe illustrant la co-intégration d'un transistor vertical dans les niveaux métalliques et d'un transistor standard CMOS.

La figure b-q présente une vue en coupe de la co-intégration entre transistors planaires et verticaux. Après avoir réalisé le transistor planaire grâce aux technologies VLSI standards, il est possible de faire croître un nanofil dopé par un mécanisme VLS. Les caractéristiques du fil sont contrôlées par l'utilisation d'un catalyseur métallique précédemment déposé par lithographie. Après croissance, le dépôt de l'empilement de grille est réalisé par ALD et CVD du diélectrique ( $\text{Al}_2\text{O}_3$ ,  $\text{HfO}_2$ , ...) et du métal de grille ( $\text{TiN}$ , ...) respectivement. La définition



de la longueur de grille se fait par gravure isotropique. Finalement, le contact supérieur est réalisé par dépôt métallique.

### B.3.3.2 Circuits Logiques Verticaux

Grâce à l'intégration verticale des transistors, il est possible d'implémenter des cellules standards directement en 3-D. Pour illustrer ceci, la figure b-r montre la réalisation d'un contact intelligent entre deux lignes métalliques. Il est effectivement possible de réaliser un contact contrôlé par n'importe quelle fonction combinatoire.

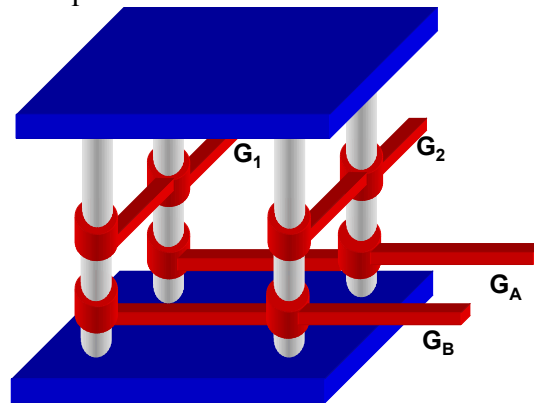


Figure b-r. Contact intelligent à base de transistors verticaux

Dans l'exemple proposé, quatre fils verticaux sont montrés. Deux transistors par fils sont assemblés en série. Quatre grilles de contrôles sous forme de mailles pilotent les fils. Les deux lignes métalliques sont ainsi connectées lorsque tous les transistors d'une branche verticale sont passants. Ceci correspond à l'équation Booléenne suivante:

$$(G_A + G_B).(G_1 + G_2) = 1$$

A partir de ce contact configurable, il est possible de construire une porte logique complémentaire en incluant la branche complémentaire de type  $p$ . La fonction NOT réalisée dans un schéma 3-D est présentée par la figure b-s. La porte est réalisée par deux transistors verticaux de dopage différents. Un transistor de type  $n$  est réalisé entre la masse et la ligne de sortie, tandis qu'un transistor de type  $p$  est placé entre la sortie et  $V_{dd}$ . Les grilles sont connectées au signal d'entrée. Il est bon de noter que différents diamètres de fils peuvent être réalisés afin de dimensionner correctement la structure au regard des spécifications voulues. Tous les circuits peuvent être étendus à cette structure.

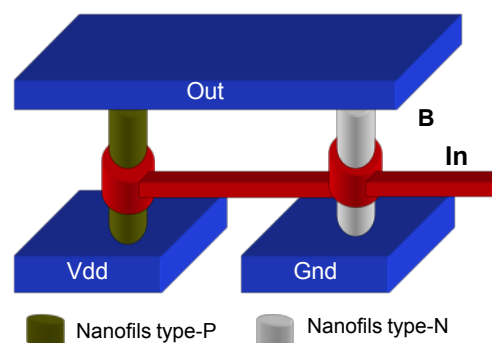


Figure b-s. Fonction NOT réalisée sous forme de cellule 3-D

### B.3.3.3 Evaluation des Performances et Méthodologie

Afin d'évaluer les performances de l'implémentation en véritable 3-D, on étudie la surface, le délai et la consommation du circuit, afin de les comparer à un équivalent CMOS. La technologie étant fortement spéculative, aucuns modèles compacts ne sont à ce jour disponibles pour effectuer les simulations de performances. Ainsi, il a été proposé l'utilisation de simulation technologique TCAD afin d'émuler le comportement d'un circuit simple depuis la modélisation technologique des dispositifs. Alors qu'il est assez standard d'utiliser la simulation TCAD pour extraire des courbes I-V, il est beaucoup plus innovant de réaliser des

simulations transitoires au niveau circuit servant à étudier l'intérêt des composants très avancés.

Le tableau b-f résume les performances atteintes par la porte NOT 3-D et par son équivalent 2-D. On peut constater que la cellule 3-D améliore clairement la surface du circuit par un facteur 31. Ce résultat vient évidemment de l'utilisation de la troisième dimension pour réaliser les composants actifs. En effet, la projection du circuit sur la face active traditionnelle du silicium est considérablement réduite. Le délai, ainsi que les fuites sont également améliorés par un facteur de 2,5 et de 14,5 respectivement. Ceci peut s'expliquer par les bons niveaux de performances atteints par la technologie de grille enrobante utilisée par les transistors verticaux. En effet, elle permet un bon contrôle électrostatique du canal et ainsi la réduction du courant  $I_{off}$  tandis que le courant  $I_{on}$  reste de bonne qualité.

**Tableau b-f. Evaluation de performances relatives à la technologie des transistors verticaux**

	Surface ( $\mu\text{m}^2$ )	Délai intrinsèque (ps)	$K_{Load}$ (ps.fF $^{-1}$ )	Puissance de fuite (pW)
<i>NON (MOS 65nm)</i>	1,6	17,5	3,8	40,1
<i>NON (3-D BE)</i>	0,05	6,9	4,7	2,8
<i>3-D BE vs. CMOS</i>	x 31,2	x 2,5	x 0,8	x 14,5

### B.3.4 Conclusion

Dans ce chapitre, nous avons investigué l'apport des technologies émergentes pour améliorer les circuits de base de configuration et de routage. Ces éléments sont fondamentaux dans les circuits reconfigurables étant donné qu'ils représentent la plus grande part de silicium et sont responsables du goulot d'étranglement limitant l'amélioration des performances. Globalement, l'amélioration de ces circuits passe par l'intégration de composants 3-D. Nous avons étudiés trois différentes technologies, qui respectivement permettent (i) d'utiliser des éléments reconfigurables passifs dans les niveaux métalliques, (ii) d'utiliser plusieurs niveaux de silicium actifs 2-D empilés et (iii) d'utiliser de « vrais » transistors 3-D réalisés verticalement entre les niveaux métalliques.

Les technologies de mémoires résistives, et plus particulièrement les mémoires à changement de phase, montrent un intérêt certain du fait de la non-volatilité, la compatibilité d'intégration avec les niveaux métalliques et la faible résistance à l'état passant (jusqu'à 50 $\Omega$ ). Nous avons ainsi proposé l'utilisation de mémoires PCM afin de réaliser des points mémoires de configuration pour la logique reconfigurable. Un circuit à base de 2 mémoires résistives et d'un transistor de programmation permettant de stocker une information et de la restituer sous forme de tension a été étudié. Nous avons montré que ce point mémoire aboutit à une amélioration en surface d'un facteur 1,5 et d'une amélioration d'un facteur 16 en temps d'écriture par rapport à une technologie flash. Finalement, nous avons étudié la réalisation d'une boîte de brassage, où les connexions programmables sont réalisées à partir d'unicques mémoires résistives. Cette implémentation aboutit à un gain en surface d'un facteur 3,4, tandis que la faible résistance de l'état passant est directement introduite dans le chemin de données.

Les procédés d'intégration 3-D monolithique ouvrent la voie vers l'intégration dense des mémoires de configuration au plus proche des circuits logiques les nécessitant. Ainsi, la réalisation de blocs FPGAs dans une approche 3-D a été faite. Les blocs réalisés sont des structures de routages mais également des LUTs. Il a été montré qu'une réalisation 3-D permet de réduire la surface jusqu'à un facteur 2 par rapport à son équivalent 2-D tandis que les performances sont améliorées d'un facteur 1,6.

Enfin, nous avons examiné les opportunités offertes par une technologie clairement en rupture permettant de réaliser des transistors verticaux. L'intégration verticale de composants actifs a un impact évident sur la surface, tandis que les performances sont conservées par les grandes dimensions accessibles (du fait des contraintes relâchées dans les niveaux métalliques). Dans

cette optique, un nouveau schéma d'implémentation de cellules logiques bénéficiant des avantages de la technologie verticale a été décrit. Les performances ont été évaluées et il a été démontré qu'à nœud technologique équivalent, les cellules 3-D permettent des gains en surface et délai de facteurs 31 et 2 respectivement.



## B.4 Impact Architectural des Circuits de Routage et de Configuration 3-D

Dans ce chapitre, l'impact architectural des circuits 3-D, développés au chapitre précédent, est étudié. Ces circuits, spécialisés dans le routage et la mémoire de configuration, ont pour objectifs d'améliorer l'architecture FPGA traditionnelle. Les diverses technologies étudiées placent des composants dans les niveaux métalliques selon trois options différentes: composants uniquement passifs pour les mémoires résistives, empilement de composants actifs pour l'intégration 3-D monolithique puis composants actifs réalisés dans une réelle configuration 3-D. Ainsi, chaque technologie sera étudiée séparément.

### B.4.1 Méthodologie et Outil d'Évaluation

L'évaluation architecturale repose sur le modèle du FPGA traditionnel. Ainsi, nous considérons l'architecture décrite en B.2.1.1. En particulier, la logique est prise en charge par des CLBs, construits avec 10 BLEs de LUT 4 entrées. 22 entrées connectent les CLBs aux lignes de routage global. Les modifications ciblent donc les transistors de routage, ainsi que les éléments mémoires.

L'évaluation est réalisée par les circuits de test MCNC [204], qui seront dans un premier temps synthétisés par l'outil ABC [183]. L'affectation technologique sera alors réalisée selon une librairie de LUTs 4 entrées, également avec l'outil ABC. Le regroupement logique des BLEs en CLBs est pris en charge par l'outil T-VPACK [186]. Finalement, le placement et le routage sont pris en charge par l'outil VPR [186]. Chaque circuit de test sera évalué deux fois, tout d'abord, avec un FPGA de référence à base de circuits SRAMs conventionnels (basés sur un procédé CMOS 65-nm), puis avec une structure utilisant la technologie en test. Le flot d'évaluation est décrit par la figure b-t.

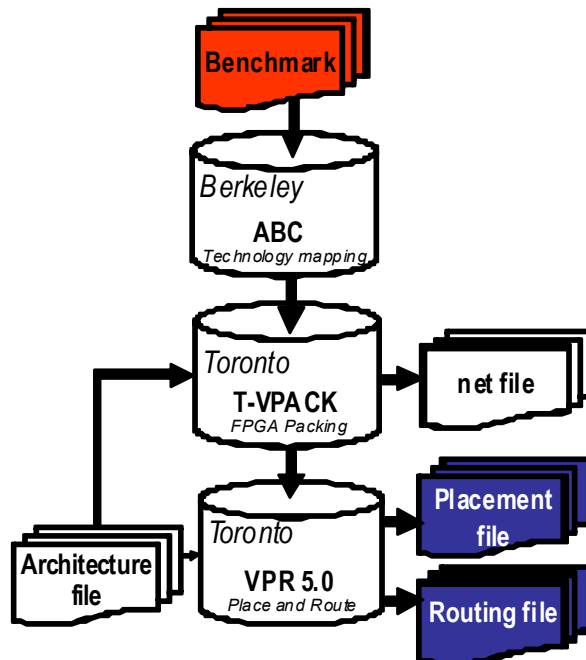


Figure b-t. Flot d'évaluation FPGA

### B.4.2 Répercussion Architecturales des Mémoires à Changement de Phase

Les améliorations potentielles de l'architecture FPGA viennent de deux points. Tout d'abord, considérant les mémoires de configuration, on se propose de remplacer les mémoires SRAMs de grande taille par des circuits compacts. De plus, tandis que cette amélioration conduit à la réduction de taille du circuit FPGA associé, il est également de grand intérêt d'utiliser la faible résistance à l'état passant des mémoires résistives pour réduire le délai de propagation

dans les circuits de routage, et ainsi améliorer les performances globales du système. Cette amélioration est faite par remplacement des boîtes de brassage CMOS par des boîtes à base de PCMs.

Le résultat d'évaluation architectural est présenté par la figure b-u. Cette figure montre une réduction du délai allant de 22% à 51%, avec une moyenne de 40%. On peut également noter un gain en surface d'environ 13% sur tous les échantillons, du à la diminution de la taille occupée par la mémoire. Dans cette application, le principal avantage en faveur de l'utilisation des mémoires résistives en lieu et place de circuits CMOS standard est la faible résistance interne à l'état passant. En effet, comme extrait à partir d'un design kit CMOS 65-nm industriel, la résistance à l'état passant d'un transistor de type  $n$  est de l'ordre de  $9\text{k}\Omega$  tandis que les mémoires PCMs sont de l'ordre de  $3,7\text{k}\Omega$  [115] en résistance passante. Ceci rend le FPGA basé sur la technologie PCM potentiellement plus rapide que son équivalent SRAM, car la résistance des éléments de routage influence directement le délai du chemin de données.

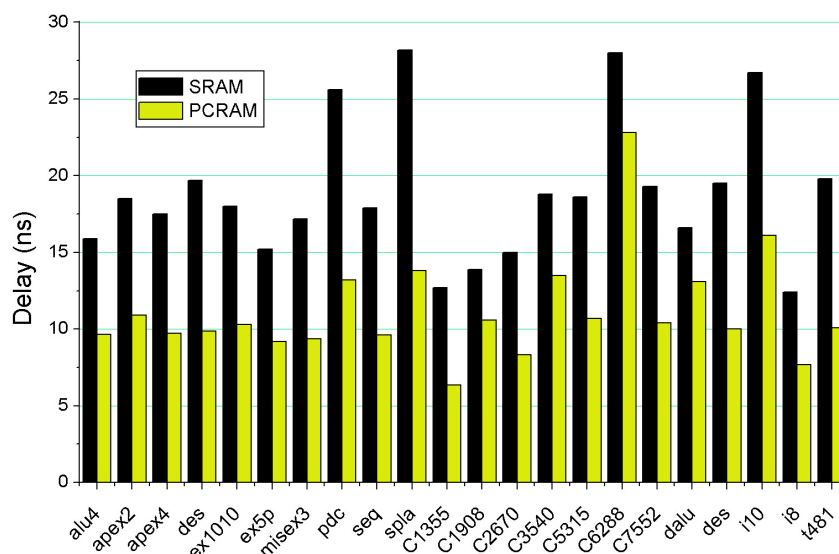


Figure b-u. Estimation du délai pour des FPGAs utilisant la technologie PCM et la technologie SRAM (Logique LUT et Routage)

### B.4.3 Répercussion Architecturale de l'Intégration 3-D Monolithique

Dans un circuit FPGA CMOS, les mémoires de configuration sont distribuées à travers toute la logique. En effet, les mémoires de configuration sont utilisées pour piloter les entrées des LUTs, mais aussi pour configurer les circuits de routage. En sus des boîtes de brassage, de nombreux multiplexeurs de routage sont présents, depuis les BLEs jusqu'aux boîtes de connexions. Chaque multiplexeur de routage utilise des mémoires 3-D pour piloter les entrées de configuration. Dans cette évaluation, nous considérons les architectures des tables de correspondance et des boîtes de brassage décrites dans les sections B.3.2.2 et B.3.2.3 respectivement.

#### B.4.3.1 Etude de Surface

A l'issue de l'affectation des différents circuits de test sur une structure 3-D intégrée par un procédé monolithique et sur un circuit de test 2-D, la surface totale du FPGA a été évaluée. L'estimation de surface normalisée est présentée par la figure b-v. Une réduction de surface allant de 20% à 25%, avec une moyenne à 21% a été identifiée. Le principal intérêt de la technologie 3-D monolithique est évidemment la diminution de l'impact de la mémoire sur la surface du FPGA. En effet, les circuits mémoires sont placés au dessus de la structure FPGA, tandis que les circuits relatifs au chemin de données sont placés sur la couche inférieure. Néanmoins, il est bon d'envisager l'impact des circuits empilés sur les capacités de routage.

En effet, les transistors des niveaux supérieurs gênent les connexions entre les niveaux inférieurs et les niveaux métalliques.

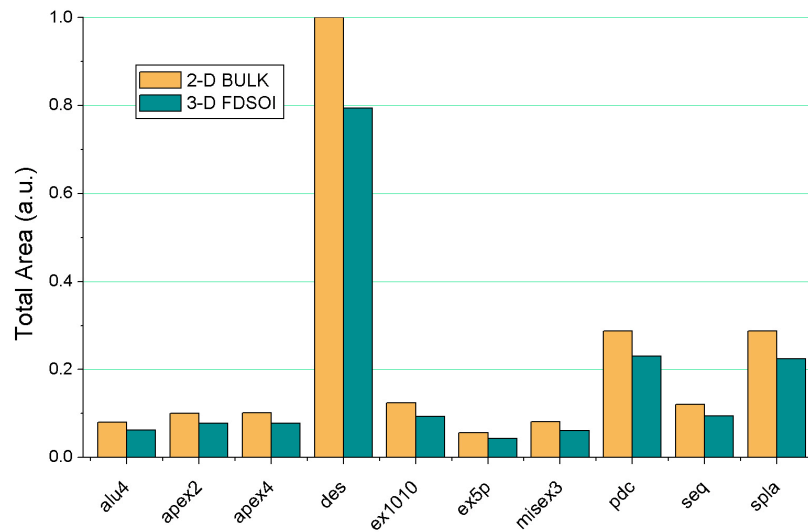


Figure b-v. Estimation de surface pour un FPGA réalisé en technologie standard et en intégration 3-D monolithique

#### B.4.3.2 Etude de Délai

L'estimation du délai critique est quant à elle présentée dans la figure b-w. On peut noter une réduction du délai de 10% à 45%, avec une moyenne à 22%. Les principaux avantages vis-à-vis de l'utilisation de l'intégration 3-D de transistors FDSOI est double. Premièrement, les performances intrinsèques du FDSOI en comparaison du bulk contribuent à l'amélioration du délai des portes et ainsi réduit le temps de calcul. Deuxièmement, la diminution des surfaces occupées par les circuits logiques et les circuits de routage entraînent naturellement une réduction des longueurs de fils et ainsi des parasites devant être piloté par les diverses portes du chemin de données. Ainsi la technologie d'intégration 3-D permet d'accélérer le fonctionnement des circuits en comparaison de son équivalent planaire.

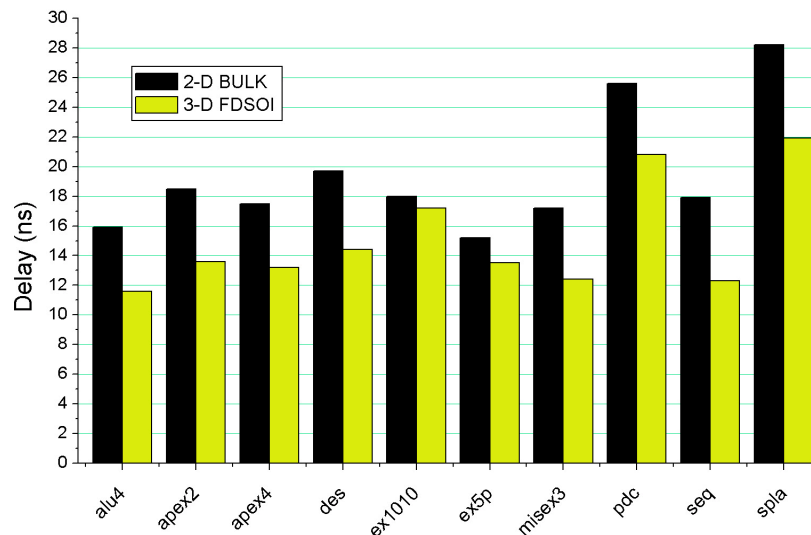


Figure b-w. Estimation du délai pour un FPGA réalisé en technologie standard et en intégration 3-D monolithique

### B.4.4 Répercussion Architecturale des Transistors à Nanofils Verticaux

La cible principale de l'intégration verticale des transistors est la réalisation de circuits de routage compacts et performants. Les ressources de routage sont formées de nombreux sous-circuits. Les connexions entre les segments peuvent être réalisées par les contacts intelligents proposés en B.3.3.2. Pour les connexions de segments, on considère le cas le plus simple où

un unique transistor connecte deux lignes métalliques. Au sein des boîtes de brassage, les points d'intersections sont réalisés en remplaçant chaque transistor de passage par un contact intelligent et par l'intégration 3-D de la mémoire de configuration associé. Enfin, on peut également considérer que les répéteurs (pour les signaux comme les horloges), peuvent être réalisées directement au sein des lignes qui les utilisent. Dans l'organisation proposée, les boîtes de connexions et les blocs logiques restent dans les niveaux actifs bas. Par cette répartition, il est possible de distribuer en part égale la complexité entre les circuits standards et les circuits placés dans les niveaux métalliques.

L'évaluation des performances a été effectuée en termes de surface et de délai critique. En comparaison d'une architecture FPGA planaire, on observe un gain en surface de 46% sur tous les circuits de tests. Cette performance est due à la diminution de la surface allouée au routage, grâce à l'implémentation verticale des circuits et des mémoires associées. Le délai critique est lui présenté par la figure b-r. Cette évaluation montre une réduction du délai allant de 37% à 48%, avec une moyenne de 42%. En fait, il est possible de réaliser des transistors de grandes dimensions (c'est-à-dire avec un grand volume de silicium actif) tout en conservant une très faible surface projeté sur les niveaux actifs standards. Les larges dimensions de ces transistors, couplées à leur architecture de type grille enrobante permettent d'obtenir de bonnes propriétés électriques et de conduire ainsi à une amélioration globale des performances du circuit.

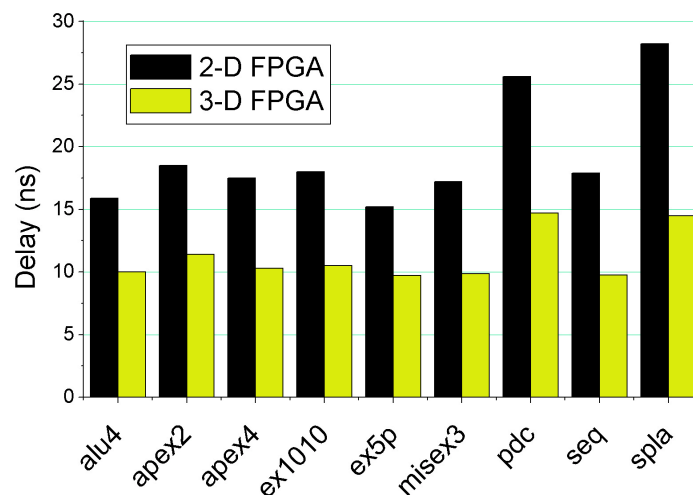


Figure b-x. Estimation de délai pour un FPGA en 2-D et en "vrai" 3-D

#### B.4.5 Conclusion

Dans ce chapitre, nous avons exploré l'impact sur l'architecture FPGA des technologies émergentes présentées au chapitre précédent. Ces technologies ont été employées pour améliorer les ressources de routage et de configuration, grâce à l'utilisation généralisée de la troisième dimension. En comparaison d'un équivalent CMOS, il a été montré qu'une technologie mémoire résistive telle que les PCMs permet une amélioration du chemin critique de 44%, grâce à la faible valeur de la résistance à l'état passant. La réalisation de transistor verticaux quand à elle permet de réduire la surface occupée par le circuit de 46%, mais également de réduire le délai critique de 42%. Ces améliorations sont dues à l'utilisation de la verticalité (pour laquelle il est possible d'obtenir une surface projetée faible au regard du volume de silicium actif). Enfin, on peut remarquer que la technologie d'intégration 3-D monolithique conduit à une amélioration en surface et en délai de 21% et 22% respectivement. Ce bon compromis de gain, allié à la maturité technologique, font de cette technologie un choix intéressant pour les systèmes reconfigurables dans un proche avenir.

## B.5 Blocs Logiques en Ruptures

Dans les deux chapitres précédents, nous avons étudié l'impact des technologies émergentes sur les circuits de configuration et de routage pour les circuits reconfigurables. Tandis que ces améliorations sont légitimes au regard des limitations des FPGAs, on peut toutefois remarquer qu'il ne s'agit que d'une amélioration incrémentale de l'architecture d'origine. Il est alors de bon ton de s'intéresser aux limites du FPGA de façon plus fine. En effet, dans une architecture classique FPGA, seulement 14% de la surface est utilisée par les blocs logiques [38]. Les blocs logiques étant le cœur du calcul, on peut qualifier l'architecture FPGA d'inefficace, vis-à-vis du rapport entre la part de calcul et la part de périphérie. Il est ainsi pertinent d'explorer de nouvelles directions pour l'organisation architecturale. Dans ce chapitre, on se propose d'étudier des briques logiques de base, dont l'objectif sera d'être la nouvelle pierre à l'édifice de l'architecture. Deux concepts principaux issus des technologies émergentes vont être utilisés: l'amélioration de la fonctionnalité d'un composant unitaire et l'amélioration de la densité d'intégration des composants de base, en utilisant des technologies sub-lithographiques comme lithographiques.

### B.5.1 Proposition 1: Electronique Carbone Ambipolaire

Depuis la découverte des Nanotubes de Carbone (CNT - *Carbon Nanotubes*) au début des années 1990 [123] et l'isolation de feuillets de graphène en 2004 [124], l'électronique carbone a connu un intérêt croissant. En particulier l'étude des propriétés intrinsèques du carbone et de ses applications potentielles pour améliorer les composants conventionnels a été dominante. En plus des niveaux de performances prometteurs, il est intéressant de remarquer que la technologie carbone est ambipolaire. Cela signifie que les transistors peuvent conduire pour des tensions de grille positives comme négatives. Traditionnellement, ce type de comportement n'est pas adapté aux applications micro-électroniques, pour lesquelles des composants unipolaires sont préférés. Alors que les technologues tentent de masquer ce phénomène, il est sensé d'étudier comment cette propriété peut être intéressante au niveau circuit. Précisément, l'ambipolarité peut être contrôlée au niveau du composant par l'ajout d'une seconde grille de contrôle. Il est alors possible d'obtenir un composant à fonctionnalité améliorée dont la polarité peut être choisie par contrôle de la tension [142]. Ce composant est le Transistor à Nanotube de Carbone Double-Grille (DG-CNFET - *Double-Gate Carbon Nanotube Field Effect Transistor*)

#### B.5.1.1 Hypothèses Technologiques

Dans [142], les auteurs utilisent un procédé d'intégration basé sur une électrode de grille arrière commune pour tous les dispositifs. Ce procédé simple pour la caractérisation n'est pas utilisable dans un contexte circuit, pour lequel chaque composant doit pouvoir être adressé unitairement, comme requis par les propositions [52, 121].

Afin d'assurer cette unicité de contrôle, le procédé de fabrication est basé sur la technologie Silicium-Sur-Isolant (SOI - *Silicon-On-Insulator*). Il est en effet possible avec cette technologie de fabriquer des îlots de silicium entourés d'oxyde. Ces îlots seront les grilles arrière des composants DG-CNFET et seront du fait isolés des autres. Les îlots de silicium peuvent être fortement dopés ou siliciurés afin d'assurer une bonne conductivité de l'électrode. L'oxyde de grille arrière est réalisé par une couche de  $\text{SiO}_2$  dont le dépôt précède le transfert de nanotubes de carbone intrinsèques. Enfin, l'oxyde de grille avant ( $\text{HfO}_2$ ) et le métal de grille (Al) avant sont déposés et structurés. Finalement, les contacts des grilles, de source et de drain sont réalisés par gravure et dépôt de métal. Le composant ainsi obtenu est présenté par la figure b-y.

#### B.5.1.2 Logique Reconfigurable

La propriété de reconfigurabilité des DG-CNFETs permet de concevoir des circuits avec un faible nombre de transistors, en regard d'une approche conventionnelle. Le contrôle de l'ambipolarité a donc été utilisé par l'Institut des Nanotechnologies de Lyon pour concevoir

une cellule logique reconfigurable [121]. La cellule logique est présentée par la figure b-z-a. Basée sur la logique dynamique, la cellule utilise sept transistors arrangés en deux étages logiques: un étage de fonction et un étage de suivi/inversion. Un schéma d'horloge à 4 phases (2 signaux de précharge  $PC_1$ ,  $PC_2$  et 2 signaux d'évaluation  $EV_1$ ,  $EV_2$ ) est utilisé pour réaliser l'opération logique. Ces signaux d'horloge sont non-recouvrant. La polarité (type  $n$  / type  $p$ ) des composants Double-Grille  $T_1$ ,  $T_2$  et  $T_3$  est contrôlée par les grilles arrières correspondantes  $V_{bgA}$ ,  $V_{bgB}$  et  $V_{bgC}$ . La cellule peut ainsi être configurée parmi un des quatorze modes d'opération accessibles, comme montré par la figure b-z-b.

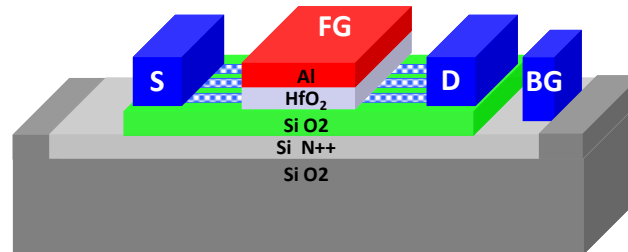


Figure b-y. Vue d'artiste d'un composant DG-CNFET utilisant le procédé de fabrication proposé et montrant les contacts de source (S), drain (D) et grilles avant- (FG) et arrière (BG)

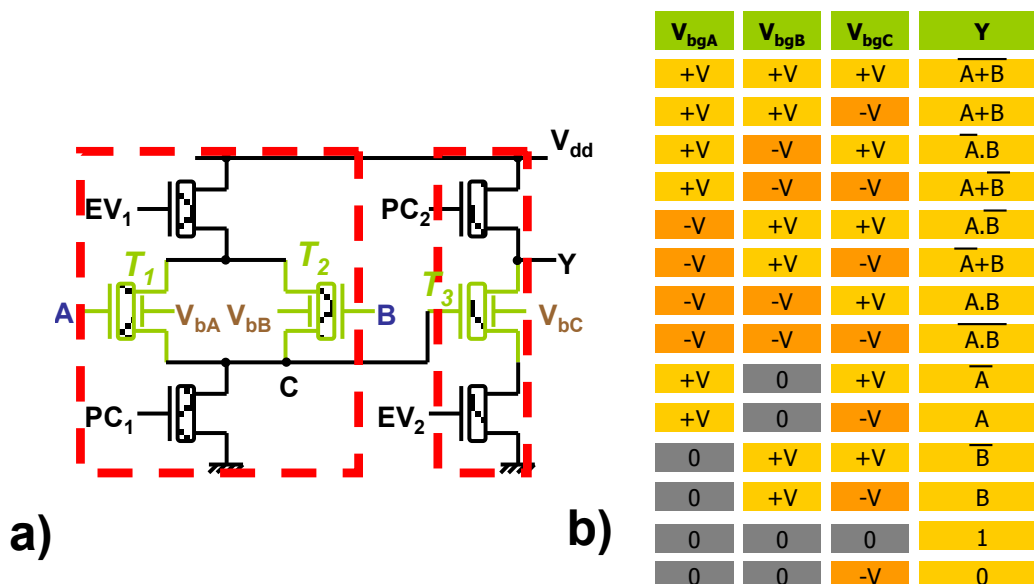


Figure b-z. a) Schéma niveau transistor et b) table de configuration pour une cellule reconfigurable à base de CNT [122]

Cette structure apparaît prometteuse pour la compacité des applications reconfigurables. Les performances de la structure ont été évaluées au travers de la surface, du délai et de la puissance consommée et comparées à un équivalent CMOS. Afin d'évaluer objectivement les performances, on considère l'utilisation d'un nœud technologique avancé 22-nm. Ce nœud est choisi comme étant le nœud standard disponible dès lors que la technologie carbone sera suffisamment mature pour une application industrielle. L'estimation des performances pour une technologie CMOS 22-nm a été faite par extrapolation des performances obtenues en technologie 65-nm vers le nœud 22-nm.

Les résultats sont centralisés dans le tableau b-g. On constate une réduction de la surface occupée d'un facteur 3,1. En effet, en comparaison d'un multiplexeur CMOS de 26 transistors présentant ainsi les mêmes propriétés, la cellule reconfigurable est réalisée avec seulement 7 transistors.

Concernant les résultats obtenus en termes de délai et de puissance consommée, on peut noter que le délai est dégradé d'un facteur 2 tandis que la puissance moyenne est améliorée d'un même facteur 2. En fait, il est bon de considérer le Produit Délai-Puissance (PDP - *Power-Delay Product*). En effet, le PDP reste constant par rapport aux transistors conventionnels. Ce

résultat est surprenant si l'on considère les bonnes propriétés intrinsèques de l'électronique carbone sur le plan de la consommation et des performances. Dans le cas de notre application logique, on peut remarquer que la puissance consommée est essentiellement dynamique et ainsi dépend principalement de la fréquence, de la charge et des capacités parasites. Ces différents paramètres restent dominants dans les performances du circuit et ainsi le PDP ne varie pas.

Tableau b-g. Evaluation globale des performances de la cellule à base de DG-CNFET

	Surface ( $\mu\text{m}^2$ )	Délai intrinsèque (ps)	$K_{\text{Load}}$ (ps.fF <sup>-1</sup> )	Puissance moyenne à 4 GHz ( $\mu\text{W}$ )
<i>MUX MOS</i>	1,191	79	5,7	3,56
<i>DG-CNFET</i>	0,39	149	124,5	1,78
<i>Gain</i>	x 3,1	x 0,53	x 0,04	x 2

### B.5.1.3 Amélioration des Circuits Logiques Dynamiques

La hausse des fonctionnalités permise par la seconde grille des DG-CNFETs peut être utilisée de nombreuses façons. Tout particulièrement, il est possible de réduire le nombre de transistors utilisés dans les circuits logiques dynamiques. En sus du gain apporté par la logique dynamique elle-même, il est possible d'employer le comportement *off* (qui permet d'avoir un transistor toujours bloqué quelque soit l'état de la grille avant) des DG-CNFET pour fusionner les transistors d'évaluation avec les transistors de fonctions. De plus, la polarité programmable permet de supprimer tous les circuits d'inversion de signaux. Ainsi, par un choix adapté des tensions de grille arrière, il est possible d'adapter le réseau à n'importe quelle entrée, qu'elle soit normale ou complémentée. La figure b-aa-a nous montre à titre d'exemple la fonction suiveur réalisée suivant le concept présenté. La cellule est basée sur 2 DG-CNFETs. Le transistor de précharge (grille PC sur le chemin vers la masse) est configuré en type *n*, grâce à la tension  $V_{\text{dd}}$  appliquée sur sa grille arrière. Les transistors d'évaluation et de signal (grille IN sur le chemin vers  $V_{\text{dd}}$ ) sont combinés; le signal d'évaluation EV étant connecté à la grille arrière et le signal d'entrée IN sur la grille avant. Un exemple plus complexe d'une fonction XOR est présenté par la figure b-aa-b

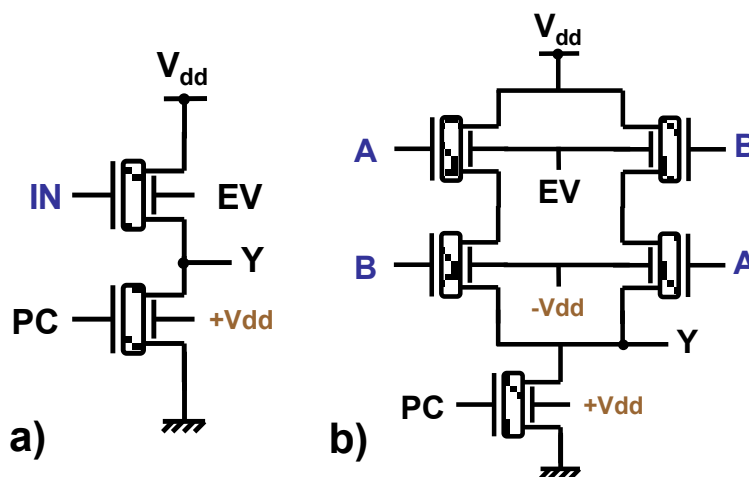


Figure b-aa. a) Schéma d'une cellule buffer et b) d'une cellule XOR

Le tableau b-h présente une comparaison entre l'approche proposée et des cellules logique dynamique CMOS conventionnelles. On peut constater que l'utilisation de transistors DG-CNFET permet une réduction du nombre de transistors jusqu'à 50% en comparaison de la logique dynamique classique.



Tableau b-h. Besoins en transistors pour les cellules dynamiques à base de DG-CNFETs

	Cellule MOS statique	Cellule MOS dynamique	Cellule DG-CNFET	Gain (DGCNTFET vs. MOS dynamique)
<i>Buffer</i>	4T	3T	2T	33%
<i>NOT</i>	2T	3T	2T	33%
<i>AND</i>	6T	4T	3T	25%
<i>OR</i>	6T	4T	3T	25%
<i>XOR</i>	12T	10T	5T	50%

### B.5.2 Proposition 2: Structures Entrecroisées à base de Composants Silicium Monodimensionnels

La technologie étudiée dans la section précédente se base sur l'augmentation de fonctionnalité du transistor afin de réduire la taille des circuits logiques. On se propose à présent d'employer une technologie permettant d'améliorer la densité d'intégration des composants actifs. Ainsi, il apparaît évident que l'augmentation de densité peut aboutir à l'obtention de circuits logiques de forte compacité. Ainsi, il est de grand intérêt de pousser les limites du transistor traditionnel dans ces limites en utilisant des nanofils comme éléments de remplacement du canal, tout en les organisant sous forme de structure entrecroisées à haute densité.

#### B.5.2.1 Hypothèses Technologiques

La réalisation d'une structure entrecroisée est envisagée par l'assemblage de nanofils réalisés par croissance CVD. La méthode CVD permet un bon contrôle des dimensions des fils grâce à l'utilisation de nanoparticules métalliques [159], avec des diamètres de l'ordre de 3-nm [160], tout en assurant un dopage in-situ lors de la croissance [161]. Les nanofils sont ensuite oxydés pour obtenir une structure à coquille [162] et déposés sur un substrat. L'alignement se fait par la méthode de Langmuir-Blodgett [163]. Ainsi, l'épaisseur de l'oxyde permet de régler l'espacement entre les fils (Figure b-bb-a). L'oxyde est ensuite gravé afin de mettre à nu le silicium actif (Figure b-bb-b). Les portions des fils jouant le rôle de simple connectique sont alors siliciurés pour améliorer les propriétés électriques [155] (Figure b-bb-c). Une seconde couche de fils est déposée sur la première en utilisant le même procédé (Figure b-bb-d). La coquille de diélectrique sert ici de diélectrique de grille et les régions siliciurées jouent le rôle de métaux de grilles (Figure b-bb-e). Enfin, les contacts métalliques sont réalisés à la périphérie de la structure (Figure b-bb-f)

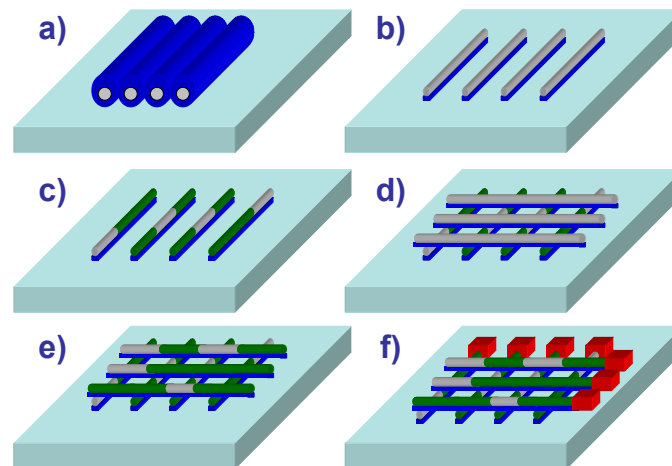


Figure b-bb. Procédé de fabrication d'une structure entrecroisée à base de nanofils. *a)* Dépôt de nanofils après croissance, oxydation et alignement par la méthode de Langmuir Blodgett. *b)* Gravure des coquilles d'oxydes. *c)* Siliciurations des fils sur les régions ne servant pas comme canaux de transistors. *d)* Dépôt, alignement et gravure de la coquille d'oxyde de la seconde couche. *e)* Siliciuration identique à l'étape *c.* *f)* Métallisation des contacts des fils autour de la structure.



### B.5.2.2 Cellule Logique Reconfigurable à base de Nanofils

Le procédé de fabrication ainsi présenté, il est possible de réaliser des structures entrecroisées à base de nanofils, dont les intersections sont des transistors. Ainsi que proposé dans [156], il est possible de construire des portes élémentaires selon un schéma de logique dynamique. Dans [65], ce concept est généralisé avec l'architecture NASIC pour concevoir des circuits complets de façon spécifique. En effet, à partir d'une structure entrecroisée de type  $n$ , il est possible de réaliser de la logique NAND-NAND [157]. Tous les opérateurs logiques peuvent être réalisés à partir de cette logique.

Afin de concevoir une cellule logique de très faible dimension, on se propose de réaliser un multiplexeur, basé sur la méthode du NASIC. Néanmoins, pour des raisons de faisabilité technologique, le routage au sein de l'architecture ne sera pas réalisé avec des nanofils. En effet, il est difficile de réaliser technologiquement des fils alignés sur de grandes dimensions dès lors que le facteur de forme est trop grand. On se propose ainsi de réaliser les longues connexions par des fils métalliques conventionnels. Afin de minimiser la taille associée à ces connexions, on se propose de limiter leur nombre au maximum. Pour ce faire, les signaux complémentaires sont rendus superflus par l'ajout d'inverseurs directement au sein de la structure logique. Un étage de suiveur est aussi ajouté afin d'assurer la synchronisation des données. La figure b-cc représente la cellule logique. Celle-ci est formée de trois étages cascades:

Entrée  $\rightarrow$  Inversion/Suivi  $\rightarrow$  NAND  $\rightarrow$  NAND  $\rightarrow$  Sortie

Les trois étages fonctionnent de façon séquentielle par un jeu d'horloges dynamiques non-recouvrant. Du fait de la structure entrecroisée et des procédés de fabrication, les deux étages de NAND ne sont pas identiques. Pour le premier étage, les entrées sont faites par les fils de l'axe  $y$  et les sorties se font sur les fils d'axe  $x$ . Pour le deuxième étage, les entrées et les sorties sont inversées. Ainsi, des transistors sont formés à la fois sur les couches de nanofils inférieure et supérieure.

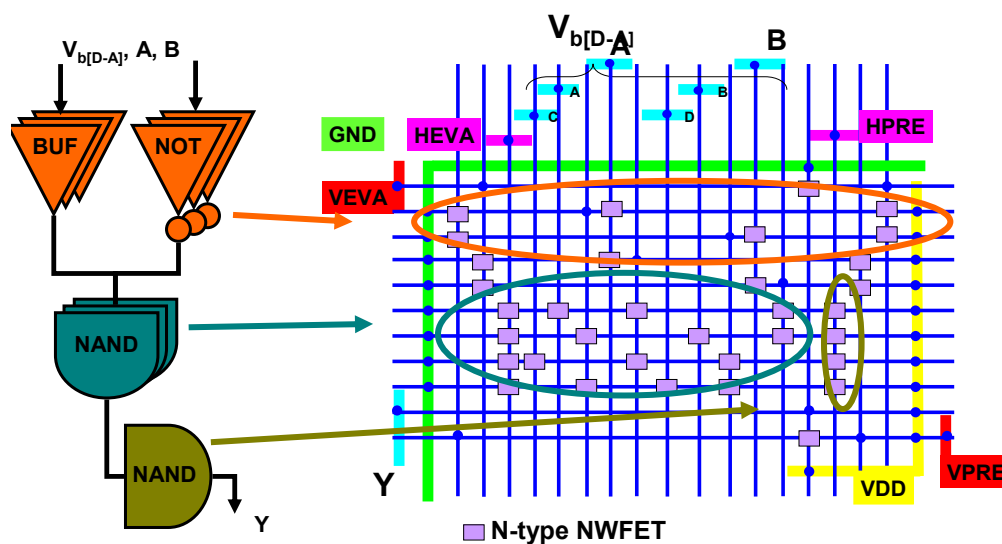


Figure b-cc. Illustration de la structure de la cellule reconfigurable dynamiquement à base de CB-NWFET

L'évaluation des performances de la cellule a été menée au travers de la surface, du délai et de la puissance consommée. Afin d'assurer une comparaison juste avec l'approche précédente, le nœud lithographique considéré est le nœud 22-nm.

Le tableau b-i résume les résultats. Du fait des dimensions sublithographiques atteintes par les fils (l'espace entre les fils est uniquement du à l'épaisseur de la coquille), il est possible de réduire la surface occupée par la fonctionnalité d'un facteur 4,1. Concernant les performances en vitesse, on peut noter que le délai intrinsèque est amélioré d'un facteur 4,9 tandis que la charge voit son influence quasiment tripler. Le délai intrinsèque dépend du chemin de

données interne à la structure. Du fait des très faibles dimensions, il est possible de réduire de façon drastique les éléments parasites internes et donc de réduire le temps de propagation dans les étages internes. Néanmoins, cette réduction des dimensions impacte aussi la taille des transistors et tout particulièrement des transistors de sortie pilotant la charge. Ces faibles dimensions de l'étage de sortie (en fait, tous les transistors sont de dimensions identiques) rendent la charge plus influente sur le délai. Enfin, concernant la puissance consommée, on observe une augmentation de la consommation d'un facteur d'environ 2 par rapport à un circuit CMOS standard. Celle-ci est due aux fuites dans les étages internes.

**Tableau b-i. Evaluation globale des performances du MUX à base de CB-NWFET**

	Surface ( $\mu\text{m}^2$ )	Délai intrinsèque (ps)	$K_{\text{Load}}$ (ps.fF <sup>-1</sup> )	Puissance moyenne à 4 GHz ( $\mu\text{W}$ )
<i>MUX MOS</i>	1,191	79	5,7	3
<i>CB-NWFET</i>	0,289	17	15,1	6,7
<i>Gain</i>	x 4,1	x 4,6	x 0,38	x 0,45

### B.5.3 Proposition 3: Structures Entrecroisées aux Dimensions Lithographiques

La proposition précédente exploite l'augmentation de densité des éléments actifs, réalisés par des techniques de croissance sublithographiques. Néanmoins, ces techniques sont souvent basées sur des hypothèses hautement spéculatives. En effet, si l'on considère le procédé précédent, il faut aligner des nanofils sur un substrat. Ces techniques conduisent à une grande variation dans la position et dans l'espacement des fils. De plus, il s'agit de réaliser ce procédé par deux fois afin de réaliser les deux niveaux de fils perpendiculaires, augmentant encore plus la difficulté de la tâche. Dans cette proposition, on se propose d'étudier une structure entrecroisée, basée sur un procédé lithographique industriel.

#### B.5.3.1 Hypothèses Technologiques

En se basant sur un procédé de Silicium-Sur-Isolant Totalelement Déplété (FDSOI - *Fully-Depleted Silicon-on-Insulator*), il est possible de réaliser des fils avec une très grande régularité, comme montré dans [171]. Dans ce travail, nous complétons les procédés existants pour un simple niveau de fils parallèles par la mise en place d'un second niveau à base de poly-silicium. Les fils du bas sont définis par photolithographie à l'espacement minimal. Les dimensions peuvent ensuite être contrôlées plus finement par oxydation et gravure en dessous des limites lithographiques [172]. Le dopage de ces lignes est autorisé. Les lignes supérieures en poly-silicium sont également lithographiques. A chaque intersection peut être défini un transistor ou une simple interconnexion si l'intersection a été préalablement dopée pour éviter tout contrôle électrostatique. Le procédé de fabrication est illustré par la figure b-dd.

#### B.5.3.2 Méthode de Conception

A partir du procédé technologique proposé, il est possible d'envisager plusieurs stratégies de réalisation de circuits.

##### i) Régions à Dopage Séparé

Il est possible de grouper séparément les régions de type  $p$  et les régions de type  $n$ . Cette technique est proche de la représentation habituelle des circuits CMOS et se trouve très adaptée aux circuits dont la sortie est un point commun à toutes les branches. Cette situation se trouve dans la plupart de circuits logiques. Un exemple est présenté sur la figure b-ee, qui nous montre un multiplexeur 4 vers 1. La structure du multiplexeur est une structure complémentaire comme montré par le schéma associé. Sur la structure entrecroisée, chaque branche correspond à une branche du circuit entre une alimentation et le point commun.

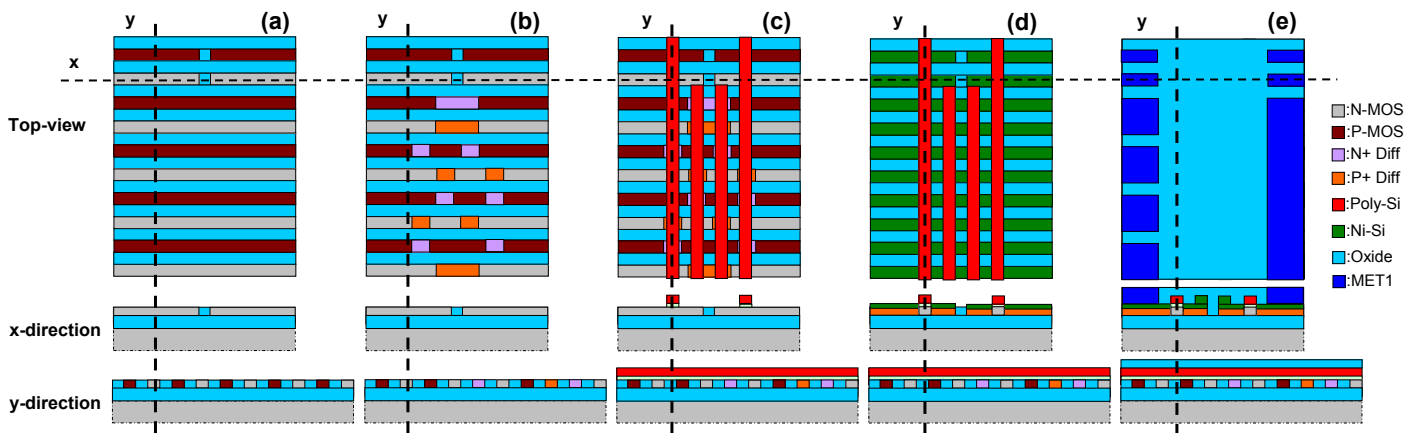


Figure b-dd. Procédé de fabrication d'une structure entrecroisée à base de FDSOI : a) Structuration en réseau et dopage des régions actives. b) Définition des régions passives. c) Dépôt de la grille. d) Finalisation des régions passives et silicidation. e) Métallisation des contacts associés aux régions passives.

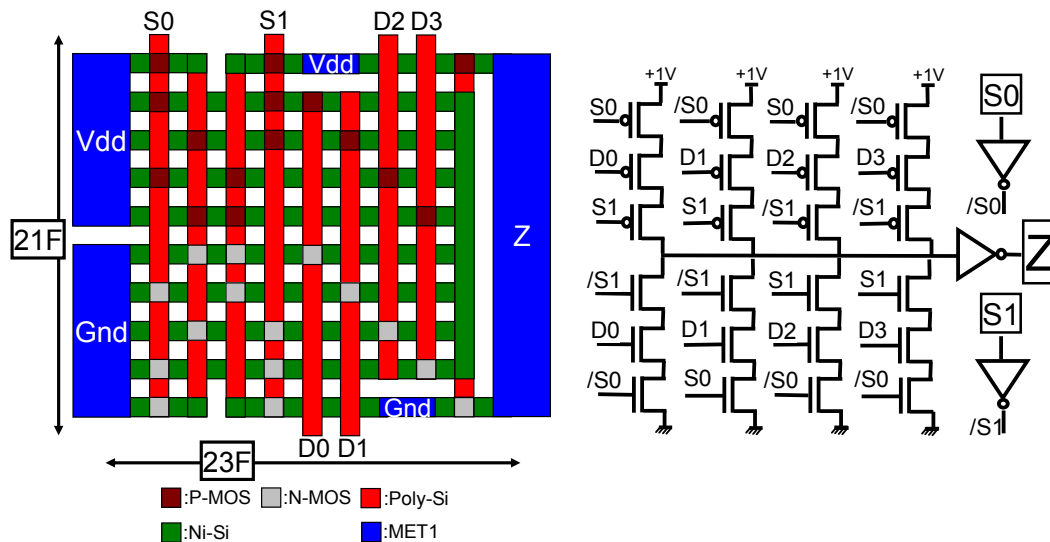


Figure b-ee. Implémentation d'un MUX 4:1 en structure complémentaire à base de régions à dopage séparé et schéma équivalent

## ii) Régions à Dopage Alterné

Une autre structuration possible consiste à alterner les régions de type  $p$  et de type  $n$ . Ce genre d'approche est tout particulièrement adapté à la réalisation de porte à transmission entre deux contacts métalliques. Cette technique est employée pour réaliser un multiplexeur 4 vers 1 dans la figure b-ff.

### B.5.3.3 Evaluation des Performances dans le Cas Idéal

Les performances d'une telle approche sont présentées dans le tableau b-j, au travers de la surface, du délai et de la puissance moyenne. La technologie étant compatible avec les procédés de fabrication actuels, la comparaison est effectuée avec une technologie 65-nm courante (et non pas une technologie 22-nm prospective). Le tableau présente les performances de la structure dans le cas idéal.

On constate la forte réduction de surface atteignable par le biais de la conception dense. En effet, la surface est réduite d'un facteur 6. Le délai intrinsèque et la puissance moyenne sont quant à eux améliorés de facteurs 1,3 et 1,9 respectivement. Ces bonnes propriétés sont dues aux performances intrinsèques du FDSOI. Enfin, on peut remarquer que la charge a une influence bien plus grande que sur l'équivalent bulk. Ceci est dû à la structure entrecroisée qui contraint les dimensions du transistor à de faibles dimensions avec  $W=L=F$  où  $F$  est

l'espacement lithographique minimal. Ainsi, les transistors ne sont pas capables de piloter des charges importantes de façon performante.

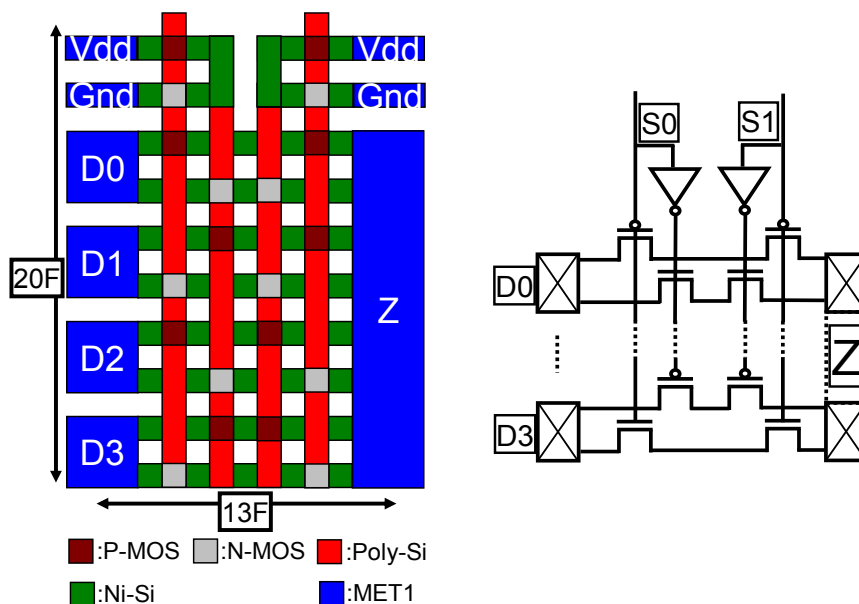


Figure b-ff. Implémentation d'un MUX 4:1 en structure à porte de transmission à base de régions à dopage alterné et schéma équivalent

Tableau b-j. Evaluation des performances d'un MUX 4-vers-1 utilisant une structure à fils entrecroisés lithographiques (parasites non pris en compte)

	Surface ( $\mu\text{m}^2$ )	Délai intrinsèque (ps)	$K_{\text{Load}}$ ( $\text{ps.fF}^{-1}$ )	Puissance moyenne à 650 MHz ( $\mu\text{W}$ )
MOS 65nm	10,4	155,2	11,2	1,47
Structure entrecroisée Dimensions compactes	1,74	121,8	97,5	0,77
Structure entrecroisée vs. CMOS	x 6	x 1,3	x 0,1	x 1,9

#### B.5.3.4 Evaluation des Performances dans le Cas Réel

Du fait de la très grande compacité atteignable, de nombreux composants parasites doivent être considérés dans les structures entrecroisées. La figure b-gg présente le modèle de parasites utilisé. On peut distinguer 2 types de croisement différents: les croisements actifs qui fonctionnent comme des transistors et les croisements passifs qui ont juste pour objectif d'assurer le routage. Dans notre procédé de fabrication, les intersections passives sont réalisées avec du silicium fortement dopé croisant des lignes de poly-silicium siliciurées. Ces lignes sont séparées par une fine couche d'oxyde de grille. Ainsi, les surfaces sous les lignes de poly-silicium forment de grandes capacités parasites (cf. l'élément A de la figure b-gg). De plus, nous avons modélisé l'impact des résistances parasites au travers des zones de routage faites en silicium fortement dopées (cf. l'élément B de la figure b-gg), ainsi que le couplage électrostatique entre les lignes (cf. l'élément C de la figure b-gg).

La figure b-hh illustre l'impact de l'oxyde de grille ( $T_{\text{OX}}$ ) sur le délai de propagation. En effet, les capacités aux intersections sont les parasites dominants. On peut voir qu'une épaisseur optimale peut être atteinte autour de 6-nm pour un nœud lithographique 65-nm. En effet, l'épaisseur d'oxyde a deux effets opposés. Lorsqu'il augmente, les capacités parasites A sont réduites, mais le contrôle électrostatique des transistors est dégradé. Il faut alors se positionner dans la zone optimale.

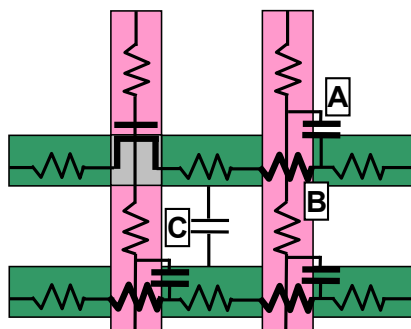


Figure b-gg. Modèle de parasites

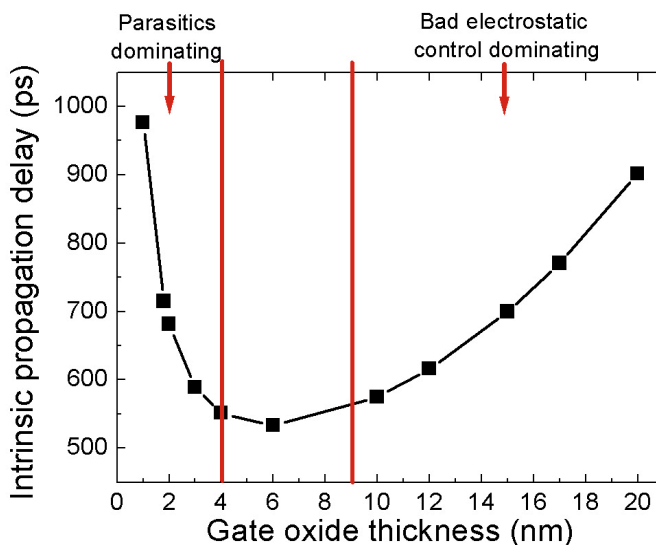


Figure b-hh. Influence de l'épaisseur de grille ( $T_{OX}$ ) sur le délai de propagation du MUX ( $F_{HP} = F_{HN} = F_V = F_I = F = 60 \text{ nm}$ )

Afin de résumer l'étude de la structure en présence de parasites, les métriques habituelles ont été évaluées et présentées dans le tableau b-k. Les évaluations ont été réalisées pour une épaisseur d'oxyde de 6-nm. On peut constater dans cette situation une forte dégradation des performances en termes de délai intrinsèque et une augmentation de l'influence de la charge. Il faut toutefois noter que ces résultats sont intimement liés au nœud technologique considéré. Effectivement, on peut s'attendre à dépasser les performances d'un circuit CMOS conventionnel pour des nœuds lithographiques autour de 16-nm comme montré par la figure b-ii.

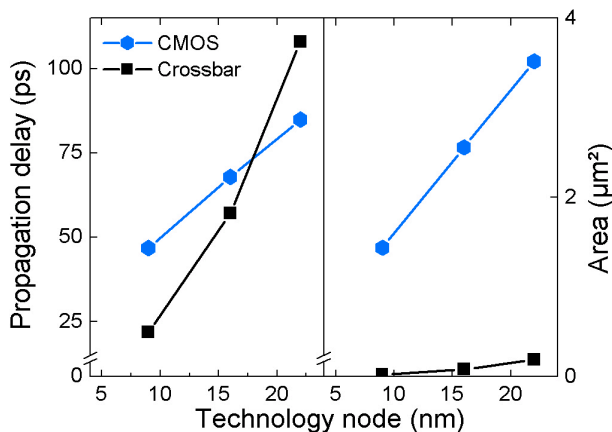


Figure b-ii. Impact de la réduction des dimensions sur la structure ( $F_{HP} = F_{HN} = F_V = F_I = F$ )

**Tableau b-k. Evaluation des performances d'un MUX 4-vers-1 utilisant une structure à fils entrecroisés lithographiques (parasites pris en compte)**

	Surface ( $\mu\text{m}^2$ )	Délai intrinsèque (ps)	$K_{\text{Load}}$ ( $\text{ps}\cdot\text{fF}^{-1}$ )	Puissance moyenne à 650 MHz ( $\mu\text{W}$ )
<i>MOS 65nm</i>	10.4	155.2	11.2	1.47
<i>Structure entrecroisée Dimensions compactes</i>	1.74	533	208	0.99
<i>Structure entrecroisée vs. CMOS</i>	x 6	x 0.29	x 0.05	x 1.48

### B.5.4 Conclusion

Dans ce chapitre, des technologies émergentes ont été utilisées afin de concevoir des blocs logiques innovants pour les applications reconfigurables. Le focus sur les éléments logique a pour objectif de définir de nouveaux circuits et d'aboutir à de nouveaux paradigmes architecturaux. L'étude a utilisé un procédé d'électronique carbone ambipolaire et deux procédés de fabrication de structures entrecroisées à base de nanofils.

L'électronique carbone présente des propriétés d'ambipolarité. Ceci consiste en la présence de comportements de type  $n$  et de type  $p$  accessibles au sein du même composant. Il est ainsi possible de réaliser des composants à la polarité contrôlable, grâce à l'ajout d'une seconde grille. Cette augmentation de la fonctionnalité du composant est employée pour réaliser une cellule logique reconfigurable et dans une nouvelle approche de conception de cellule logique dynamique. Une réduction de surface d'un facteur 3,1 est observée, grâce au faible nombre de transistors utilisés.

Nous avons ensuite étudié la réalisation de structures de nanofils entrecroisés pour réaliser une cellule logique reconfigurable compacte. On a pu constater que l'utilisation de l'organisation entrecroisée permet de réduire les dimensions jusqu'à un facteur 4 par rapport à l'équivalent CMOS. Néanmoins, le procédé technologique reposant sur un arrangement de fils sublithographiques reste incertain du fait de sa maturité. Il a ainsi été proposé de réaliser ces arrangements à partir d'un procédé industriel FDSOI éprouvé. Par ce biais, il est possible de construire des portes logiques avec une réduction de dimensions jusqu'à un facteur 6. Néanmoins, ce procédé génère de nombreux composants parasites, qui doivent être pris en compte dans l'étude des performances. Il est toutefois possible de les réduire par le biais d'une optimisation technologique, telle que le choix de l'épaisseur d'oxyde. Avec un choix correct, le délai peut être réduit d'un facteur 1,3 et il a finalement été montré que les performances de la structure à fils croisés doivent dépasser les performances de l'organisation CMOS traditionnel à partir du nœud 16-nm.

## B.6 Propositions Architecturales en Ruptures et Analyse de Performances

Dans le chapitre précédent, nous avons introduit un certain nombre de propositions permettant d'obtenir des cellules logiques compactes. Dans la section B.2.1.2, nous avons évoqué le fait que seule une faible portion d'un FPGA est utilisée pour réaliser les opérations logiques combinatoires et séquentielles (14% de la surface totale), par opposition aux circuits de routage et aux mémoires (86% de la surface totale). Ainsi, l'utilisation de cellules logique à grain ultra-fin apparaît étrange, car le rapport de surface ne peut qu'empirer si l'on considère que les circuits de routage ne réduisent pas leurs dimensions dans les mêmes proportions que la logique. Dans ce chapitre, nous allons décrire une nouvelle organisation architecturale permettant d'utiliser à bon escient les circuits logiques à grain ultra-fin. En sus de l'architecture, un outil d'évaluation ad-hoc sera présenté et utilisé pour évaluer les performances de l'organisation.

### B.6.1 Proposition Architecturale

L'architecture traditionnelle FPGA est organisée selon une répartition précise, comme décrite par la figure b-jj-b. Cette organisation selon 4 niveaux (LUT, BLE, CLB puis FPGA) suit le gabarit architectural propre à notre étude (figure b-jj-a)

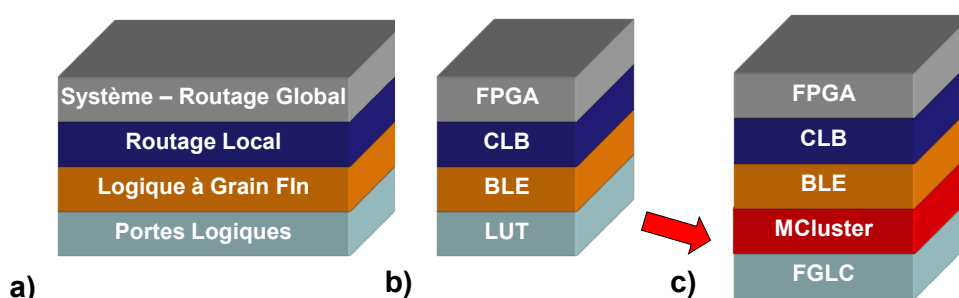


Figure b-jj. Gabarit d'architecture (a), modèle FGPA (b) et modèle adapté (c)

Dans la situation abordée dans ce chapitre, les traditionnelles LUTs sont remplacées par des Cellules Logiques à Grain Fin (FGLC - *Fine Grain Logic Cell*). Néanmoins, il apparaît difficile de réaliser un bloc logique élémentaire en utilisant directement une unique FGLC. En effet, la faible granularité de la cellule conduirait à un large déséquilibre entre les circuits logiques et les circuits de routage. Afin de pallier à ce problème, on se propose de modifier la hiérarchie FPGA comme montré par la figure b-jj-c. Pour réduire le déséquilibre, on regroupe les cellules au sein d'une structure intermédiaire, appelée *MCluster*, compatible en termes de taille avec les niveaux hiérarchiques standard des FPGAs. Il est alors possible de disposer de blocs logiques présentant des performances accrues par rapport aux LUTs. La convergence vers l'architecture FPGA se fait au niveau du BLE.

#### B.6.1.1 *MCluster*

Au niveau des *MClusters*, des cellules logiques élémentaires à grain ultra-fin sont utilisées. La terminologie « grain ultra-fin » recouvre deux aspects. Premièrement, elle correspond à la taille du circuit logique, qui se trouve être compactée en comparaison de son équivalent CMOS. Deuxièmement, elle couvre aussi la granularité de la fonction réalisée. Dans les cas étudiés précédemment, on s'est centré sur des opérateurs logiques à 2 entrées tandis que, dans un FPGA conventionnel, les blocs sont souvent à 4 entrées.

Afin d'améliorer la couverture logique, on propose de regrouper les cellules logiques en assemblages matriciels 2-D. Les cellules logiques sont alors arrangées par niveaux, où des connexions sont possibles uniquement entre les étages adjacents. Ceci permet d'éviter les longues interconnexions et de maximiser la connectivité locale. L'organisation est décrite par la figure b-kk. Les dimensions de la structure sont décrites par la largeur  $w$  (width) et la profondeur  $d$  (depth) et apparaissent ainsi:  $MCluster\_d\_w$ .

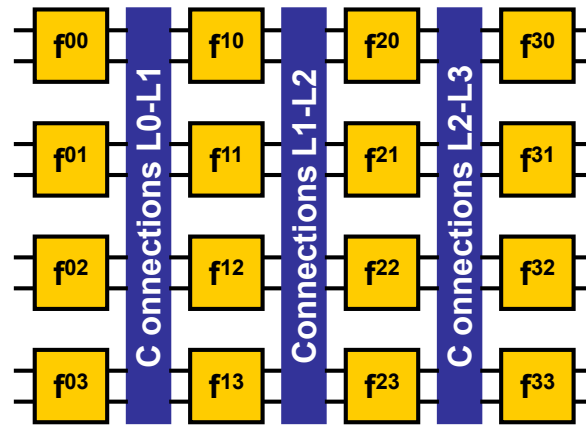


Figure b-kk. Approche MCluster pour les architectures configurables (MCluster\_4\_4)

Concernant les connexions au sein de la matrice, une topologie d'interconnexion totale n'est pas envisageable, car elle conduirait à une trop grande complexité et à une importante surface perdue. A la place, on propose d'utiliser une topologie d'interconnexion incomplète. Afin de maximiser la diversité de la topologie, on adapte les technologies utilisées dans les réseaux d'information. En effet, les Réseaux d'Interconnexions Multi-niveaux (MINs - *Multistage Interconnection Networks*) sont conçus pour interconnecter avec le meilleur compromis entre la performance et le nombre de liens nécessaires. Il existe de nombreuses topologies MIN. Néanmoins, on se concentre sur quatre permutations typiques [180, 181]: Banyan (Figure b-ll-a), Baseline (Figure b-ll-b), Flip (Figure b-ll-c) et Modified Omega (Figure b-ll-d), pour laquelle des modifications ont été apportées afin de maximiser le brassage. Etant donné que les topologies sont fixes et statiques, le choix de la topologie est réalisé par le concepteur durant la phase de développement. Les performances des topologies seront étudiées dans la suite du chapitre.

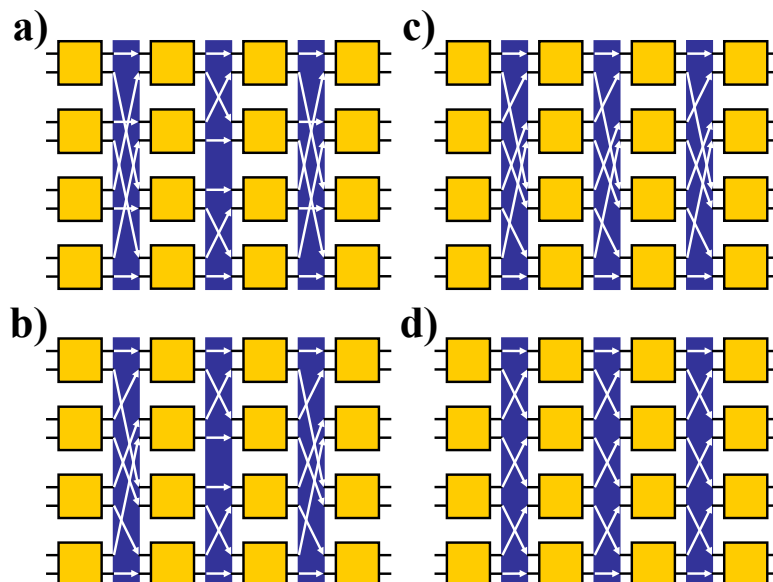


Figure b-ll. Matrices de 16 portes reconfigurables avec des topologies d'interconnexions fixes a) Banyan, b) Baseline, c) Flip, d) Modified Omega)

### B.6.1.2 Organisation de la Hiérarchie BLE/CLB/FPGA

La figure b-mm montre l'organisation des blocs logiques jusqu'au niveau CLB. Les BLEs sont composés de  $MClusters_{d_w}$  afin de réaliser la logique combinatoire. Dans l'approche traditionnelle, chaque sortie de LUT peut être verrouillée par une bascule. Le reste de l'organisation du CLB reste similaire à l'approche d'origine. Toutes les sorties des BLEs sont connectées à l'interconnexion globale mais peuvent également être connectées aux entrées du CLB selon une interconnexion complète. Les  $d*N$  signaux de retour ainsi que les  $I$  entrées du CLB peuvent alors être routées vers n'importe quel MCluster par le biais de multiplexeurs



reconfigurables. Enfin, un unique signal d'horloge est utilisé pour contrôler les éléments séquentiels. Les blocs logiques ainsi adaptés, l'architecture standard du FPGA en îlot est envisagée.

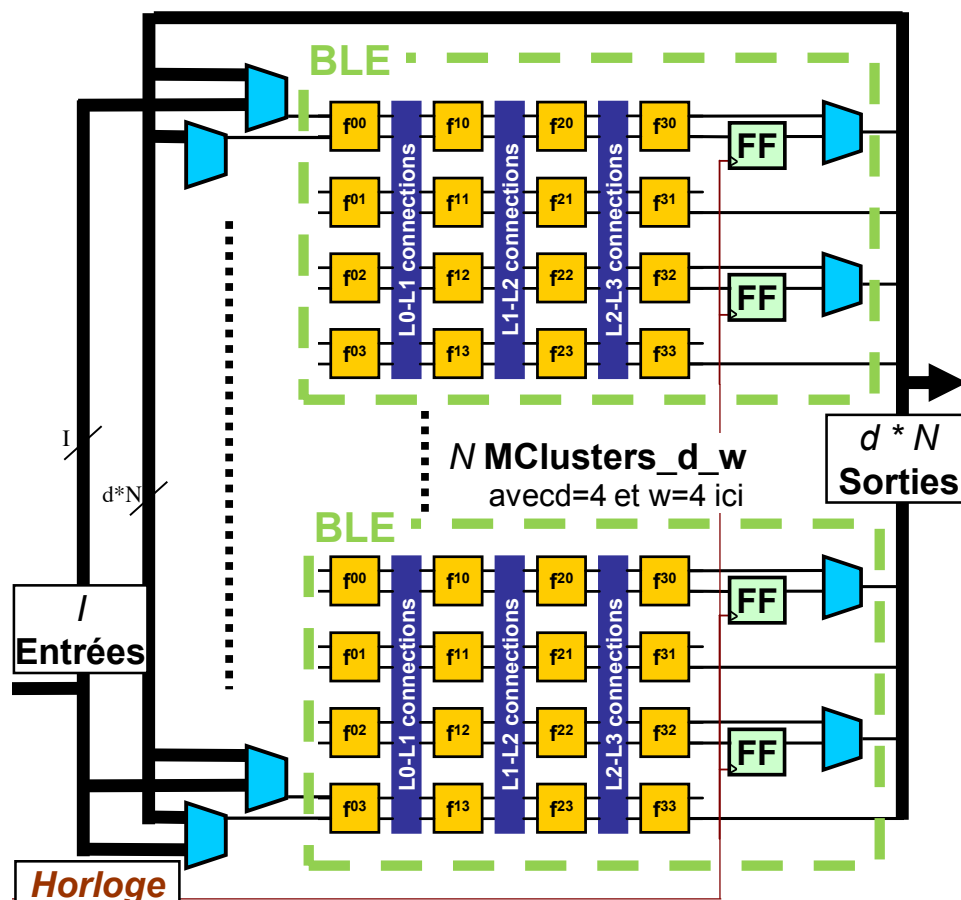


Figure b-mm. Proposition de CLB à base de MClusters

## B.6.2 Outils d'Evaluation

Afin de comparer de façon juste l'architecture proposée à un schéma FPGA existant, nous définissons un flot d'évaluation adapté aux spécificités de l'architecture.

### B.6.2.1 Flot

L'évaluation d'une architecture FPGA peut être menée par une suite d'outils adaptés, comme décrite dans la partie B.4.1. Ces outils permettent de décrire de façon détaillée l'architecture. La suite d'outils VTR (AAPACK et VPR [187]) utilise en effet un langage XML de description d'architecture afin de modéliser de nombreux blocs logiques complexes. Néanmoins, il est nécessaire de remarquer que certaines spécificités de notre architecture ne sont pas supportées par l'outil. En particulier, les topologies d'interconnexions fixes ne peuvent être décrites nativement par la suite VTR. Un outil spécifique a ainsi été conçu afin de permettre le groupement de cellules logiques sous forme de MClusters. L'outil dédié est appelé MPack. Le flot complet, présenté par la figure b-nn permet ainsi de définir de très nombreuses organisations architecturales.

### B.6.2.2 MPack

L'outil MPack (pour Matrix Packer) a pour objectif de grouper un ensemble de cellules logiques sous forme de MClusters. La figure b-oo décrit l'organisation interne de l'outil. Il est basé sur deux algorithmes principaux ayant en charge respectivement : l'affectation d'un groupe de cellules sur une matrice (Mapper) et le regroupement des cellules par affinité (Clusterer). Afin d'assurer la communication avec le reste du flot, les circuits sont importés et exportés selon une description au format BLIF. Les caractéristiques du MCluster sont fournies en termes de dimensions et de topologies, afin de permettre la construction du gabarit par le générateur d'architecture.

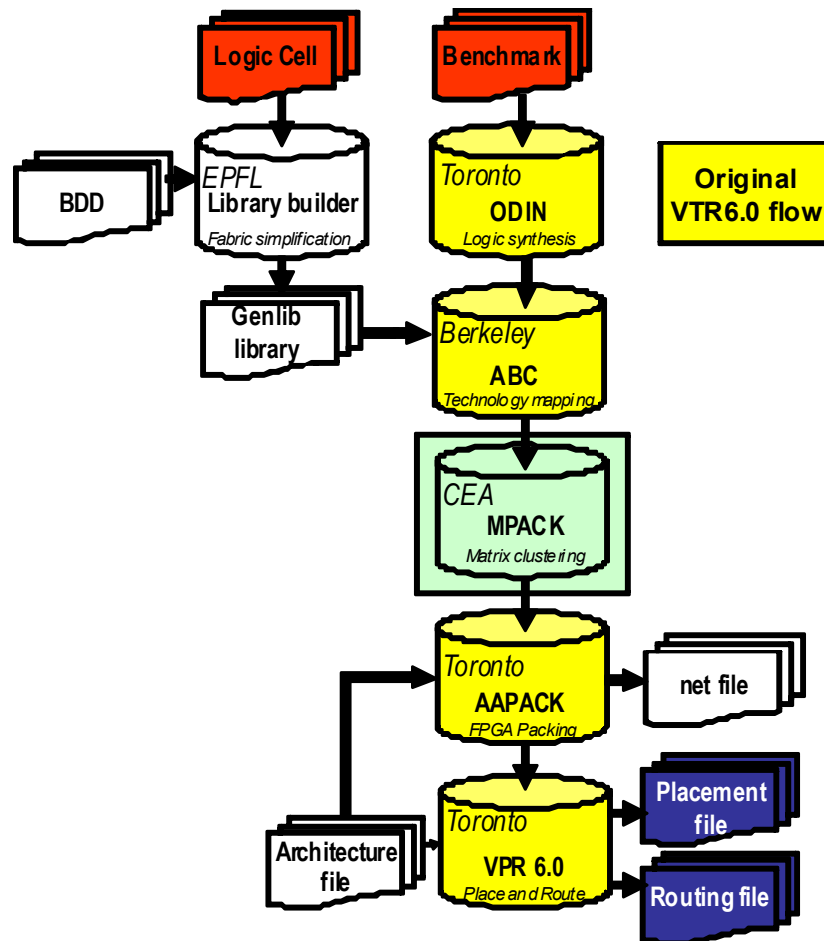


Figure b-nn. Flot d'évaluation compatible avec les technologies en ruptures

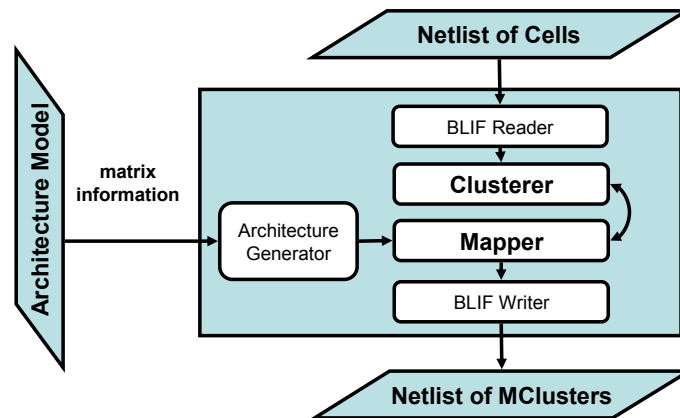


Figure b-oo. Flot de fonctionnement de MPack

i) Algorithme d'Affectation

L'algorithme d'affectation a pour objectif de faire correspondre un ensemble de cellules logiques sur une architecture de type MCluster. Cette opération est effectuée en deux étapes.

Tout d'abord, une optimisation du circuit à affecter est réalisée pour s'adapter aux contraintes physiques de l'architecture. En effet, du fait de l'organisation matricielle et de la topologie fixe, il est nécessaire d'ajouter des éléments de synchronisation. Cette synchronisation permet d'étendre les chemins d'entrée et de sortie, ainsi que d'adapter les connexions à leur réalité physique. A titre d'exemple, dans un circuit, une porte peut être connectée à de nombreuses autres portes, alors que physiquement, seul un nombre limité de portes est connecté à une sortie dans l'organisation MCluster. Il est alors nécessaire d'ajouter des opérateurs logiques pour augmenter le nombre de sortie disponible.

Par la suite, la fonction logique est affectée sur l'architecture par la force brute. Tout d'abord, le circuit à affecter est analysé en profondeur. En effet, pour chaque nœud de la structure, les branches filles sont identifiées et récursivement explorées. Ceci permet l'identification des connexions entre les nœuds et l'initialisation de la séquence d'affectation. Les nœuds sont ensuite affectés aux cellules, tout en respectant les connexions physiques. L'affectation est faite de manière exhaustive, c'est-à-dire que toutes les combinaisons possibles sont étudiées tant qu'une solution n'a pas été trouvée et que toutes les combinaisons n'ont pas été testées. L'algorithme d'affectation récursif est présenté par la figure b-pp.

```

FUNCTION Mapping_procedure(current node)
  IF (all children/parents of current node already placed) THEN
    FOR (empty cell connected to children or parents) THEN
      Store the potential positions
    END FOR
  ELSE Store all the positions on the considered level
  END IF

  IF (No positions have been found) THEN
    MAPPING FAILURE
  END IF

  FOR (Each potential position)
    Place the current node at the position
    Ret=Mapping_procedure(Next node)
    IF (Ret is MAPPING FAILURE) THEN
      Unplace the current node from the position
    END IF
    IF (Ret is MAPPING SUCCESS AND No more nodes) THEN
      MAPPING SUCCESS
    END IF
  END FOR

  IF (No potential positions have been correct) THEN
    MAPPING FAILURE
  END IF
END FUNCTION

```

Figure b-pp. Algorithme d'affectation de l'outil MPack (pseudo-code)

## ii) Algorithme de Groupement

L'algorithme de groupement a pour objectif de grouper les cellules logiques en paquets de MClusters. Ce type d'opération est très courant dans le monde des outils FPGA. Ainsi, l'algorithme utilisé est dérivé de l'outil VPACK [17]. L'outil construit chaque MCluster séquentiellement. Le pseudo-code de l'algorithme est donné par la figure b-qq

L'algorithme commence par choisir une graine pour le MCluster courant. Comme décrit dans [17], le meilleur choix de graine est fait en sélectionnant la cellule avec le plus grand nombre

d'entrées, étant donné que celles-ci sont des ressources rares des CLBs. Ensuite, l'outil choisi la cellule avec la plus grande « attraction » pour le MCluster courant, tout en s'assurant que le MCluster résultant reste viable. Le test de viabilité est fait par la méthode d'affectation précédemment décrite et consiste à vérifier que la cellule ajoutée au MCluster courant forme un MCluster qui peut être affecté sur l'architecture physique. Si tel est le cas, la cellule est définitivement ajoutée. L'attraction entre une cellule logique  $LC$  et le MCluster courant  $MC$  est donnée, comme décrit par [17], par le nombre d'entrées et de sorties communes:

$$Attraction(LC) = |Nets(LC) \cap Nets(MC)|$$

Cette procédure continue tant que (i) le MCluster n'est pas plein ou que (ii) l'ajout d'une cellule ne conduit pas à un assemblage non viable. Lorsque le MCluster est plein, une nouvelle graine est trouvée et le groupement recommence pour un nouveau MCluster.

```

Remaining_Cells: set of unclustered logic cells
C: set of cells packed in the current MCluster
MClusters: set of MClusters

MClusters = NULL
WHILE (Remaining_Cells != NULL) // More Cells to cluster
  C = Seed (Remaining_Cells) // Most Used Inputs and legal Cell
  WHILE (|C| < (MCluster_w.MCluster_d)) // Cluster is not full
    Best_Cell = MaxAttractionLegalCell(C, Remaining_Cells)

    IF (Best_Cell == NULL)
      BREAK
    END IF

    Remaining_Cells = Remaining_Cells - Best_Cell
    C = C U Best_Cell
  END WHILE
MClusters = MClusters U C
END WHILE

```

Figure b-qq. Algorithme de groupement de l'outil MPack (pseudo-code)

### B.6.3 Evaluation des Topologies d'Interconnexions

Grâce à l'outil décrit dans la section précédente, on se propose d'étudier l'impact des topologies d'interconnexions sur l'architecture de type MCluster. L'analyse porte sur une structure MCluster\_4\_4. Des jeux de 1000 graphes aléatoires de complexité allant de 6 à 16 nœuds sont affectés sur l'architecture.

Il est ainsi possible d'étudier la capacité du couple matrice/topologie à réaliser des fonctions complexes. En considérant le pourcentage de fonctions correctement affectées sur la matrice par rapport au nombre d'échantillons de test, on obtient le taux de succès (*success rate*). La figure b-rr montre les résultats obtenus pour les topologies Banyan, Baseline, Flip et Modified Omega. Pour les topologies Banyan, Flip et Baseline, le taux de succès est autour de 80% pour une fonction contenant 6 nœuds. A 12 nœuds, le taux est d'environ 25%. La différence entre ces trois topologies est relativement faible. En revanche, pour la topologie Modified Omega, le taux de succès avoisine les 90% pour un graphe de 6 nœuds et 40% pour un graphe

de 12 nœuds. Ceci montre très clairement la meilleure performance de la topologie Modified Omega pour l'application MCluster.

Ce résultat s'explique car cette topologie présente le meilleur taux de brassage. En effet, dans les couches d'interconnexions, on peut remarquer des paires de cellules possédant les mêmes entrées. Le jeu de fonctions réalisables par deux cellules possédant les mêmes entrées est  $N$ , tandis que si les entrées sont disjointes, il est possible de réaliser  $2.N$  fonctions différentes. Dans la topologie Banyan, il y a 6 paires de cellules possédant les mêmes entrées pour seulement 2 paires dans la topologie Modified Omega. Ce meilleur brassage induit donc une meilleure diversité logique au niveau du MCluster.

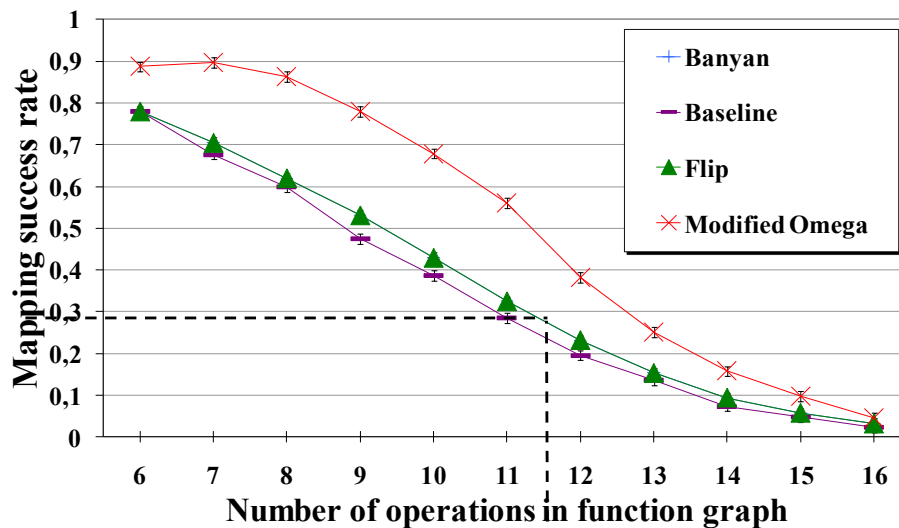


Figure b-rr. Taux de succès de l'étape d'affectation sur les topologies d'interconnexions Banyan, Modified Omega, Flip et Baseline avec des matrices de dimensions 4 par 4.

### B.6.4 Evaluation de l'Architecture Complète

On se propose à présent d'évaluer les performances de la structure complète par rapport à un équivalent CMOS. La suite d'outils d'évaluation va donc être utilisée dans son ensemble pour évaluer les circuits de tests MCNC et ISCAS'89 [204] sur l'architecture proposée. Le circuit CMOS de référence est formé par des CLB de 10 BLEs, eux-mêmes réalisés par les LUTs à 4 entrées. 22 entrées connectent les CLB avec le routage global. Les paramètres technologiques des différentes structures sont extraits d'une technologie 65-nm. Les métriques considérées sont la surface et le délai du chemin critique.

La figure b-ss représente l'estimation de surface pour un FPGA réalisé avec des MClusters\_3\_3; la granularité 3 par 3 correspondant à la meilleure compacité possible. La figure nous montre une réduction de la surface allant de 38% à 62%, avec une moyenne à 46%. Ces résultats sont clairement dus à la compacité des cellules logiques à grain fin et à l'efficacité du groupement MCluster. En effet, un MCluster\_3\_3 occupe une surface de  $63433\lambda^2$  (avec  $\lambda$  la moitié de la longueur de grille minimum) alors qu'une LUT à 4 entrées occupe  $54292\lambda^2$ . Il est alors pertinent de remarquer que pour une taille du même ordre de grandeur, la fonctionnalité du MCluster est bien plus grande que son équivalent LUT. En effet, à taille identique, un MClusters peut recevoir 50% de données d'entrée en plus et sortir 3x plus de signaux. En lien avec l'efficacité de l'outil de groupement, il est alors possible d'atteindre ces intéressantes performances.

La figure b-tt montre la distribution du délai du chemin critique. Tandis que le délai moyen ne varie que très peu (environ 10% de gain), il est intéressant de voir que la structure MCluster aboutit à une réduction de l'écart type dans la distribution du délai. Ce comportement peut être expliqué par l'impact de la granularité ultra-fine, qui favorise l'utilisation de routage local plutôt que les longues interconnexions. Ainsi, on observe une réduction du chemin critique.

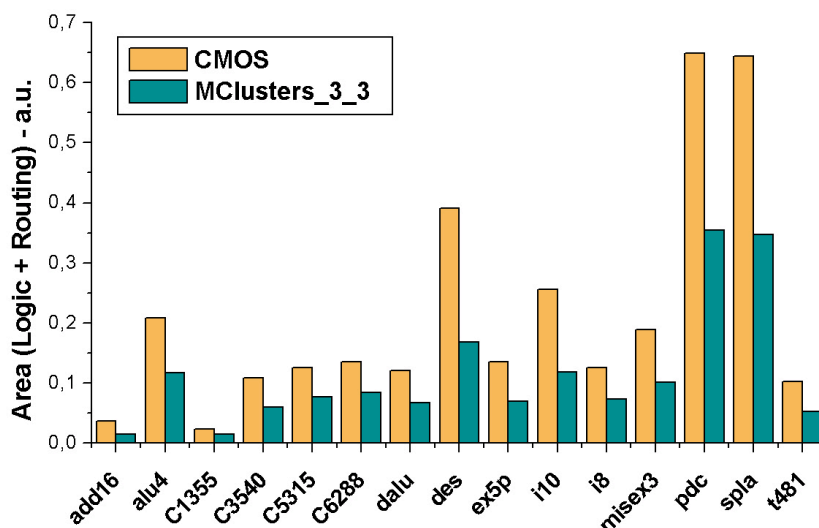


Figure b-ss. Estimation de surface pour un FPGA réalisé avec des MClusters\_3\_3 et des technologies classique CMOS

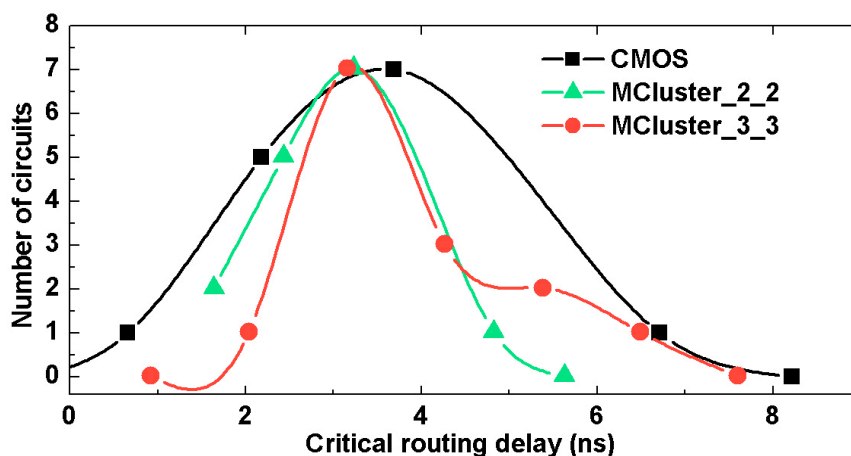


Figure b-tt. Répartition du délai de routage critique pour des FPGAs basés sur du CMOS, des MClusters\_2\_2 et des MClusters\_3\_3

## B.6.5 Conclusion

Dans ce chapitre, nous avons exploré l'impact sur l'architecture des cellules à grain ultra-fin réalisées dans le chapitre précédent. L'utilisation de telles cellules logiques nécessite d'adapter l'architecture. Effectivement, il n'est pas possible de remplacer directement la logique combinatoire d'un FPGA par ce type de cellules sans arriver à un large déséquilibre entre les structures de routage et la logique.

Nous avons ainsi proposé une architecture adaptée à cette logique à grain ultra-fin. Les cellules logiques sont regroupées selon une architecture matricielle et interconnectées par niveaux grâce à topologies d'interconnexions fixes, afin de réduire l'impact inhérent au routage. Ces blocs permettent de remplacer les LUTs dans l'architecture FPGA conventionnelle.

Afin de comparer de manière objective cette organisation avec l'équivalent LUT, il a été nécessaire de développer une suite d'outils d'évaluation, capable de décrire l'architecture conçue. Basé sur l'outil VTR, un outil spécifique gérant les topologies d'interconnexion fixes a été créé et a ainsi permis l'évaluation de l'architecture.

Dans une première approche, une étude locale de l'architecture est réalisée, afin d'étudier l'impact de la topologie d'interconnexion sur la capacité d'affectation. Nous avons vu que la topologie Modified Omega donne le meilleur taux d'affectation avec environ 90% de réussite pour des fonctions de 6 opérateurs. Dans une seconde approche, l'étude de l'architecture

complète a été menée et nous avons pu observer une réduction de la surface de 46% en moyenne par rapport à son équivalent MOS. De plus, nous avons aussi découvert que le délai de routage tend vers une plus faible distribution, rendant la structure plus contrôlable.

La structure proposée est extrêmement prometteuse pour les futurs systèmes FPGA. De plus, la suite d'outils générique mise en place pour l'évaluation ouvre la voie vers l'étude de structures toujours plus en rupture.

## **B.7 Conclusion**

Dans cette thèse, nous avons évalué le potentiel des technologies émergentes pour les futurs systèmes reconfigurables. Le sujet est tellement vaste, qu'il n'est pas envisageable de réaliser de coûteux développements technologiques pour chacune des technologies en lice. Ainsi, afin de réduire le temps de cycle et les coûts, on se propose de mettre en place, d'un point de vue architectural, une stratégie d'évaluation rapide pour les technologies en rupture. L'évaluation se fait selon un schéma architectural précis. L'architecture FPGA a été choisie car elle permet de tirer le meilleur parti des technologies émergentes tout en assurant un gabarit commun et équitable pour la comparaison. Ainsi, après avoir positionné le contexte et la méthodologie dans le chapitre B.1, un état de l'art relatif aux architectures reconfigurables a été présenté dans le chapitre B.2. Afin d'améliorer l'architecture FPGA traditionnelle, nous nous sommes concentrés dans le chapitre B.3 sur l'amélioration des structures de routage et de configuration. Alors que ce chapitre se focalise sur l'étude des circuits de base, le chapitre B.4 évalue l'impact de ces différentes améliorations du point de vue de l'architecture complète. Dans le chapitre B.5, un travail est effectué sur les blocs logiques de base, afin de concevoir de nouvelles graines architecturales. L'utilisation de ces nouveaux blocs dans l'architecture est évoquée dans le chapitre B.6, où une nouvelle approche est présentée et évaluée. Cette évaluation passe par la mise en place d'outils d'évaluation spécifiques. Ainsi, ce travail voit des contributions à de nombreux niveaux allant des circuits de base à l'architecture en passant par la méthodologie et les outils. A titre de conclusion globale sur ce travail, les contributions de différentes natures vont être détaillées.

### **B.7.1 Méthodes et les Outils**

Afin de proposer une méthode d'étude rapide et à faible coût des technologies émergentes, nous avons développé une nouvelle méthode d'évaluation basée sur l'évaluation architecturale (chapitre A.1). En effet, un gabarit d'architecture générique est défini. Chaque technologie émergente est alors employée pour améliorer les performances du gabarit. Il apparaît ainsi possible de définir un flot d'évaluation capable de s'adapter à de nombreuses applications. Plus particulièrement, nous proposons un flot d'évaluation complet capable d'évaluer les performances du gabarit d'architecture au travers de nombreux paramètres technologiques. (Chapitres B.4 et B.6). De plus, afin de s'adapter à des architectures non conventionnelles telles que l'utilisation d'arrangements matriciels avec des topologies d'interconnexions fixes, un outil spécifique a été développé (chapitre B.6). Cet outil appelé MPack est intégralement compatible avec les outils standards.

### **B.7.2 Mémoires et Circuits de Routage pour FPGA**

Dans un FPGA, les ressources de routage et la mémoire occupent plus de 80% de la surface complète du circuit. Ces circuits étant distribués tout au long du circuit, afin d'en assurer la programmation, il apparaît intéressant de chercher à les répartir dans une approche 3-D (chapitre B.3). Nous avons étudié l'intérêt de trois technologies différentes : les mémoires résistives à changement de phase, l'intégration 3-D monolithique et l'utilisation de transistors verticaux.

Les mémoires résistives sont des composants passifs non-volatiles qui peuvent être programmés entre deux états résistifs stables. Ces composants peuvent être intégrés dans les niveaux d'interconnexions, du fait de leur compatibilité avec les procédés métalliques. Un point mémoire de configuration a été proposé. Il est basé sur l'utilisation de mémoires résistives arrangées selon un diviseur de tension, afin de retenir intrinsèquement un niveau logique. Ce nœud améliore son équivalent non-volatile flash d'un facteur 1,5 en termes de surface et d'un facteur 16,6 concernant le temps d'écriture. De plus, du fait de la faible résistance à l'état passant, il est tentant d'utiliser les mémoires directement dans le chemin de donnée pour créer des interrupteurs à haute performances. Il devient alors envisageable de combiner la porte à transmission et la mémoire de configuration en un unique nœud résistif à



deux terminaux. Nous avons ainsi montré que ce nœud améliore la surface d'un facteur 3,4 par rapport à la mémoire flash.

L'intégration 3-D monolithique a pour objectif de superposer plusieurs couches de silicium actif, avec une grande densité de connexions. Il est alors possible de réaliser les circuits relatifs au chemin de données du FPGA sur une couche, tandis que les mémoires de configurations sont positionnées sur l'autre couche. Afin de valider le concept, nous avons réalisé une simple LUT à 2 entrées et un point de croisement programmable. Nous avons montré qu'une LUT réalisée dans cette technologie peut apporter un gain en surface d'un facteur 2 par rapport à son équivalent planaire. Il a également été montré que le délai intrinsèque pouvait être réduit d'un facteur 1,6, l'influence de la charge réduite d'un facteur 6,1 et la puissance moyenne réduite d'un facteur 2. Cette tendance est également observée dans les circuits de routage.

De façon plus prospective, une technologie d'intégration de transistors verticaux a été proposée afin de placer des transistors réellement dans la troisième dimension. Ce genre de circuits ouvre la voie vers des transistors de larges dimensions possédant un très faible impact sur la surface active planaire. Nous avons ainsi démontré une réduction d'un facteur 31 pour la surface, d'un facteur 2 pour le délai intrinsèque et d'un facteur 14 pour les fuites par rapport à une technologie CMOS équivalente.

### **B.7.3 Blocs Logiques de Traitements pour FPGA**

Tandis qu'une réduction de la taille des blocs logiques d'un FPGA peut sembler désuète de prime abord (seulement 14% de la surface du FPGA occupée par les blocs logiques), il faut en fait considérer les opportunités qui sont données de voir apparaître de nouveaux circuits de base. Ces circuits peuvent ensuite conduire à de nouvelles organisations architectures dont la finalité est de réduire le déséquilibre entre la logique et les circuits annexes. Deux principaux types de technologies ont été employés : l'utilisation d'une technologie améliorant la fonctionnalité du composant (soit un composant avec des fonctionnalités innovantes pour une dimension identique à l'existant) et l'utilisation de technologies améliorant la densité d'intégration (soit plus de composants dans le même espace).

Concernant l'amélioration de la fonctionnalité, nous avons employé la technologie carbone et le DG-CNFET. Un DG-CNFET est un transistor pouvant être configuré dans un de ses trois différents états (*n*, *p* et *off*) par simple changement de la tension de grille arrière. Ce composant permet de concevoir une cellule logique compacte capable de réaliser 14 fonctions Booléennes avec seulement 7 transistors. Nous avons montré qu'une telle cellule améliore la surface d'un facteur 3,1 par rapport à son équivalent CMOS.

Deux technologies basées sur l'augmentation de la densité d'intégration ont ensuite été explorées. Les technologies émergentes permettent l'utilisation en maillage dense de composants aux dimensions en dessous de la limite lithographique. Ce type de structures entrecroisées est utilisé pour réaliser une cellule logique à grain ultra-fin. Nous avons montré qu'une telle structure permet un gain d'un facteur 4,1 en surface et 4,6 en délai intrinsèque, grâce aux très petites dimensions atteintes par les composants. Néanmoins, il est nécessaire de remarquer qu'une telle technologie reste très controversée vis-à-vis de sa faisabilité technologique. En effet, les dimensions sublithographiques nécessitent des techniques non encore maîtrisées à ce jour, en particulier pour l'alignement des fils. En ce sens, nous avons proposé un procédé de fabrication de structures entrecroisées à partir d'un procédé FDSOI industriel, basé sur la lithographie. Ainsi, nous avons pu démontrer une amélioration de la surface d'un facteur 6 tandis que la consommation se trouve améliorée d'un facteur 1,48 en comparaison de son équivalent CMOS.

### **B.7.4 Architectures**

Grâce aux circuits de base proposés, il est possible d'améliorer les architectures de traitement selon deux axes distincts : une amélioration incrémentale des circuits de routage et de

configuration (chapitre B.4) et une amélioration en rupture des blocs logiques, ce qui conduit à de nouvelles hypothèses architecturales (chapitre B.6).

Concernant l'amélioration des structures de routage et de mémoire, on peut remarquer qu'il est possible de réduire la surface de 46% par rapport à un FPGA standard à base de SRAM CMOS. Ce gain est obtenu grâce à l'utilisation de la technologie la plus compacte qui se base sur la réalisation de transistors verticaux. Cette technologie est effectivement la meilleure pour réduire la surface des circuits de routage. Les autres technologies permettent de réduire la surface de 21% et de 13% pour l'intégration 3-D monolithique et pour les mémoires résistives à changement de phase respectivement. Concernant le délai critique, un gain maximum de 44% est observé pour les mémoires résistives. Ces dernières présentent la plus faible résistance à l'état passant parmi les technologies utilisées, et forment ainsi des interrupteurs de bonnes performances. L'intégration verticale de transistor donne également de bons résultats avec un gain de 42% et l'intégration 3-D monolithique améliore le délai de 22%.

Vis à vis de l'organisation d'architecture en rupture, nous avons décrit un nouveau schéma architectural adapté à la logique à grain ultra-fin (chapitre B.6). Alors qu'il n'est pas possible d'utiliser l'architecture d'un FPGA standard, du fait du trop grand déséquilibre entre la logique et le routage, on propose de mettre en place des matrices de cellules. Ces matrices ont pour but de réaliser les fonctions combinatoires. Afin de prévenir la complexité d'interconnexion, ces matrices sont connectées par étage en utilisant des topologies fixes. Ces structures sont appelées MClusters et sont utilisées en tant que remplaçants des LUTs. Grâce à l'outil spécifique que nous avons développé, nous avons étudié l'impact de la topologie sur les performances du MCluster. Il a ainsi été montré que la topologie Modified Omega donne les meilleures performances, grâce à ses bonnes propriétés de brassage. Nous avons ensuite réalisé l'évaluation d'une architecture FPGA complète. Utilisant le meilleur compromis de taille pour le cluster (MCluster\_3\_3), il a été montré que la surface totale d'un FPGA pouvait être réduite de 46% et que le délai critique pouvait être mieux contrôlé.





# *Abstract*

*Keywords:* Process-Design co-integration, PCM, Monolithic 3-D, NWFET, DG-CNFET, Crossbars, Nanoarchitectures, Benchmarking

---

For the last four decades, the semiconductor industry has experienced an exponential growth. According to the ITRS, as we advance into the era of nanotechnology, the traditional CMOS electronics is reaching its physical and economical limits. The main objective of this thesis is to explore novel design opportunities for reconfigurable architectures given by the emerging technologies. On the one hand, the thesis will focus on the traditional FPGA architecture scheme, and survey some structural improvements brought by disruptive technologies. While the memories and routing structures occupy the major part of the FPGAs total area and mainly limit the performances, 3-D integration appears as a good candidate to embed all this circuitry into the metal layers. Configuration and routing circuits based on back-end compatible resistive memories, a monolithic 3-D process flow and a prospective vertical FETs process flow are introduced and assessed within a complete architectural context. On the other hand, the thesis will present some novel architectural schemes for ultra-fine grain computing. The size of the logic elements can be reduced thanks to inherent properties of the technologies, such as the crossbar organization or the controllable polarity of carbon electronics. Considering the granularity of the logic elements, specific fixed and incomplete interconnection topologies are required to prevent the large overhead of a configurable interconnection pattern. To evaluate the potentiality of this new architectural scheme, a specific benchmarking flow will be presented in order to explore the ultra-fine grain architectural design space.

# *Résumé*

*Mots Clés:* Conception Proche-Techno, PCM, 3-D monolithique, Nanofils, DG-CNFET, Crossbars, Nanoarchitectures, Benchmarking

---

Durant les quatre dernières décennies, l'industrie des semi-conducteurs a connu une croissance exponentielle. En accord avec l'ITRS et à mesure de l'approche vers le nanomètre, les promesses sont énormes et les composants sont réduits à leurs limites physiques et économiques ultimes. L'objectif principal de cette thèse est d'explorer les opportunités offertes par les technologies émergentes pour la conception d'architectures reconfigurables. Tout d'abord, la thèse se centre sur l'architecture FPGA traditionnelle et étudie des améliorations structurelles apportées par des technologies en ruptures. Tandis que les structures de configuration et de routage occupent la majeure partie de la surface d'un FPGA et limitent ces performances, l'intégration 3-D apparait comme une bonne opportunité pour déplacer ces circuits dans les niveaux métalliques. Des circuits de configuration et de routage utilisant des mémoires résistives compatibles back-end, un procédé d'intégration 3-D ou encore un procédé de réalisation de transistors verticaux seront introduits et évalués dans un contexte architectural complet. Par la suite, la thèse présente de nouvelles propositions architecturales pour la logique à grain ultra-fin. La taille des éléments logiques peut être réduite grâce aux propriétés inhérentes de certaines technologies, telles que l'arrangement en structures entrecroisées de nanofils ou la polarité contrôlable des transistors carbonés. Considérant le changement de granularité des opérateurs logiques, des topologies d'interconnexions fixes sont nécessaires afin d'éviter l'important surcoût dû à l'interconnexion programmable. Afin d'étudier les possibilités de cette organisation, un flot d'évaluation est présenté et utilisé pour explorer l'espace de conception relatif aux architectures à grain ultra-fin.