



# Types for Detecting XML Query-Update Independence

Federico Ulliana

## ► To cite this version:

Federico Ulliana. Types for Detecting XML Query-Update Independence. Databases [cs.DB]. Université Paris Sud - Paris XI, 2012. English. NNT : . tel-00757597

**HAL Id: tel-00757597**  
**<https://theses.hal.science/tel-00757597>**

Submitted on 27 Nov 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.



Université Paris Sud

Laboratoire de Recherche en Informatique

## THÈSE DE DOCTORAT

*présentée par :* **Federico Ulliana**

# Types for Detecting XML Query-Update Independence

### THÈSE DIRIGÉE PAR

BIDOIT-TOLLU Nicole  
COLAZZO Dario

*Professeur, Université Paris Sud  
Maître de Conférence, HDR, Université Paris Sud*

### RAPPORTEURS

SENEILLART Pierre  
SCHMITT Alan

*Maître de Conférence, HDR, Telecom ParisTech  
CR, HDR, INRIA Rennes*

### EXAMINATEURS

MANOUSSAKIS Yannis  
TALBOT Jean-Marc

*Professeur, Université Paris Sud  
Professeur, Université de Provence*



# Abstract

In the last decade XML became one of the main standards for data storage and exchange on the Web. Detecting XML query-update independence is crucial to efficiently perform data management tasks, like those concerning view-maintenance, concurrency control, and security. This thesis presents a novel static analysis technique to detect XML query-update independence, in the presence of a schema. Rather than types, the presented system infers *chains* of types. Each chain represents a path that can be traversed on a valid document during query/update evaluation. The resulting independence analysis is precise, although it raises a challenging issue: recursive schemas may lead to infer infinitely many chains. This thesis presents a sound and complete approximation technique ensuring a finite analysis in any case, together with an efficient implementation performing the chain-based analysis in polynomial space and time.



# Résumé

Pendant la dernière décennie, le format de données XML est devenu l'un des principaux moyens de représentation et d'échange de données sur le Web. La détection de l'indépendance entre une requête et une mise à jour, qui a lieu en absence d'impact d'une mise à jour sur une requête, est un problème crucial pour la gestion efficace de tâches comme la maintenance des vues, le contrôle de concurrence et de sécurité. Cette thèse présente une nouvelle technique d'analyse statique pour détecter l'indépendance entre requête et mise à jour XML, dans le cas où les données sont typées par un schéma. La contribution de la thèse repose sur une notion de type plus riche que celle employée jusqu'ici dans la littérature. Au lieu de caractériser les éléments d'un document XML utiles ou touchés par une requête ou mise à jour en utilisant un ensemble d'étiquettes, ceux-ci sont caractérisés par un ensemble de chaînes d'étiquettes, correspondants aux chemins parcourus pendant l'évaluation de l'expression dans un document valide pour le schéma. L'analyse d'indépendance résulte du développement d'un système d'inférence de type pour les chaînes. Cette analyse précise soulève une question importante et difficile liés aux schémas récurifs: un ensemble infini de chaînes pouvant être inférées dans ce cas, est-il possible et comment se ramener à une analyse effective donc finie. Cette thèse présente donc une technique d'approximation correcte et complète assurant une analyse finie. L'analyse de cette technique a conduit à développer des algorithmes pour une implantation efficace de l'analyse, et de mener une large série de tests validant à la fois la qualité de l'approche et son efficacité.



# Contents

<b>1</b>	<b>Introduction</b>	<b>ix</b>
<b>2</b>	<b>Preliminaries</b>	<b>1</b>
2.1	XML . . . . .	1
2.2	Schema Languages . . . . .	3
2.3	Query and Update Languages . . . . .	8
2.4	Formal Semantics . . . . .	13
2.4.1	Semantics of XPath Steps . . . . .	13
2.4.2	XQuery Semantics . . . . .	15
2.4.3	XQuery Update Semantics . . . . .	19
2.5	Query-Update Independence . . . . .	25
<b>3</b>	<b>Query Update Independence: State of the Art</b>	<b>27</b>
3.1	Decidability and Complexity of Exact Static Independence Analysis . . .	27
3.2	Approximate Algorithms . . . . .	31
3.2.1	Schema-based independence analysis for XML updates [BC09a] . .	32
3.2.2	Commutativity analysis for XML updates [GRS08] . . . . .	36
3.2.3	Destabilizers and independence of XML updates [BC10] . . . . .	41
<b>4</b>	<b>Independence Analysis based on Schema-Chains</b>	<b>45</b>
4.1	Introduction . . . . .	45
4.2	Chain Inference . . . . .	48
4.2.1	Step Chain Inference . . . . .	48
4.2.2	Query Chain Inference . . . . .	51
4.2.3	Update Chain Inference . . . . .	56
4.3	Correctness and Precision of Chain Inference . . . . .	61
4.3.1	Soundness of Query Chain Inference . . . . .	62
4.3.2	Soundness of Update Chain Inference . . . . .	65
4.4	Independence Analysis . . . . .	67
4.4.1	Infinite Analysis . . . . .	67
4.4.2	Finite Analysis . . . . .	70
4.4.3	Soundness of the Finite Chain Analysis . . . . .	79



<b>5</b>	<b>Implementing the Chain-based Independence Analysis</b>	<b>81</b>
5.1	Introduction . . . . .	81
5.2	Storing Chains on DAGs . . . . .	86
5.3	Chain Inference with CDAGs . . . . .	93
5.4	Complexity of Chain Inference with CDAGs . . . . .	114
<b>6</b>	<b>Experiments</b>	<b>117</b>
6.1	Benchmarks and Experimental Settings . . . . .	117
6.2	Experimental Results . . . . .	119
<b>7</b>	<b>Extensions</b>	<b>127</b>
7.1	Queries and Updates . . . . .	127
7.1.1	Query rewriting . . . . .	127
7.1.2	New Inference Rules . . . . .	128
7.2	Schemas . . . . .	133
7.2.1	Attributes . . . . .	133
7.2.2	Keys and foreign-keys integrity constraints . . . . .	133
7.2.3	Extended DTDs . . . . .	135
<b>8</b>	<b>Conclusions and Future Perspectives</b>	<b>137</b>
8.1	Future Perspectives . . . . .	138
<b>9</b>	<b>Proofs</b>	<b>141</b>
9.1	Soundness of Chain Inference in the Possibly Infinite Case . . . . .	141
9.1.1	Soundness and Completeness of Schema Type Relations . . . . .	142
9.1.2	Axis chain inference . . . . .	144
9.1.3	Completeness of step chain inference . . . . .	147
9.1.4	Soundness of Query Chain Inference . . . . .	148
9.1.5	Soundness of Update chain Inference . . . . .	157
9.1.6	Soundness of Chain-based Independence . . . . .	161
9.2	Soundness of the Chain Analysis for the Finite Case . . . . .	162
9.2.1	Folding Lemma for a Single Expression . . . . .	162
9.2.2	Folding Lemma with conflict preservation . . . . .	175

# Chapter 1

## Introduction

In the last decade, XML has established itself as one of the main standards for data storage and exchange on the Web. Along with XML, a family of languages for querying and manipulating XML data have been devised, such as W3C standards XPath, XQuery and XSLT languages [SCF<sup>+</sup>07]. The recent W3C standardization of an update facility for XML [RCD<sup>+</sup>11] attracted a big deal of attention by the database community. Many old questions related to update optimization, already recognized as of crucial importance for relational databases, are open again for XML. Three central issues related to XML updates are view-maintenance, concurrency, and security.

**View-maintenance** Materialized views are precomputed queries from the database, that are used to expedite query answering. Materialized views have been shown to improve query evaluation performance by up to several orders of magnitude. View-maintenance is the problem of propagating into a materialized view the changes in the base data made by updates. View maintenance is a non-trivial process. Incremental re-computations of the view can take time proportional to the database size just to determine that no operation has to be done on the view and, in general, the problem is non-trivial even for XPath based queries and updates [BGMM09].

**Concurrency** It is well known that, for ensuring high performances in the execution of concurrent query and update expressions, it is fundamental to evaluate them *in parallel*. Parallel evaluations, with a proper level of safety, can be ensured once one is able to define concurrency control mechanisms, locking schemes and schedulers which guarantee serializability of transactions on XML documents. All of this could be achieved, for instance, by means of a *commutativity* analysis, able to check if two expressions can be safely executed on a database independently of their order. Such an analysis would permit also to optimize logical plans for query languages with side effects, as shown in [GORS08].

**Security** Security views are a well studied framework to provide controlled access to data XML [FCG04]. The idea is that each user-group is provided with an XML view

consisting of all and only the information that the users of that group are authorized to access or to update. It is of crucial importance to verify that updates coming from a user-group are executed respecting access control policies.

A way to entail an advance of technology in all such contexts, is to provide an efficient solution to the problem of detecting *XML query-update independence*. Query-update independence holds when a query and an update over a database do not interact, and therefore the update does not alter query result.

How can independence detection help tasks of view-maintenance, concurrency and security control? Intuitively, when a view is specified as a query, being able to determine that there is no interaction between a query and an update, allows to skip the whole view maintenance phase. Concerning concurrency control, if a query and an update do not interact, then they commute, meaning that they can be safely executed in any order, or even in parallel. Concerning security views, when these specify the part of data that cannot be accessed or modified, independence can be used to ensure that updates can be executed without violating the security policy expressed by the views.

A query and an update are independent when the query result is not affected by update execution, on any possible input database. In all contexts where detecting query-update independence is of crucial importance, benefits are amplified when query-update independence can be checked *statically*. In order to be useful, every static analysis technique must be sound: if query-update independence is statically detected, then independence does hold. The inverse implication (completeness) cannot be ensured in the general case, since static independence detection is undecidable (see [BC09a]). This means that if a static analyzer is used, for instance, in a view maintenance system, sometimes views are re-materialized after updates even if not needed, because the analysis has not been smart enough to statically detect a view-update independence. Useless view re-materialization frequently occurs if a static analyzer with low precision is adopted. This can lead to great waste of time, since view materialization cost can be proportional to the database size.

High precision of static independence analysis can be ensured by taking into account schema information. In many contexts, schemas are defined by users, mainly by means of the DTD or XML Schema languages, while in other contexts quite precise schemas, in the form of a DTD, can be automatically inferred, by using accurate and efficient existing techniques like the one proposed by Bex et al. in [BNSV10]. Schema-based detection of XML query-update independence has been recently investigated. The state of the art technique has been presented by Benedikt and Cheney in [BC09a]. This technique infers from the schema the set of node types traversed by the query, and the set of node types impacted by the update. The query and the update are then deemed as independent if the two sets do not overlap. This technique is effective since the static analysis *i)* is able to manage a wide class of XQuery queries and updates, *ii)* can be performed in a negligible time, and *iii)* as a consequence, even on small documents, can avoid expensive query re-computation when independence wrt an update is detected. However, the technique has

some weaknesses. As illustrated in [BC09a], in some cases, independence is not detected by the static analysis, due to some over-approximation made by the type inference rules.

## Contributions

In this dissertation we propose a novel schema-based approach for detecting XML query-update independence. Differently from [BC09a, Che08, CGMS06], our system infers sequences of labels, hereafter called *chains*, to perform a static analysis of the data accessed by the queries and the updates. Intuitively, for each node that can be selected by a query/update path in a schema instance, the system infers a chain recording *i*) all labels that are encountered from the root to the node and *ii*) the order of traversal. This information is at the basis of a precise static independence analysis.

The main contribution of this work is a precise algorithm to detect independence for a query-update pair  $q-u$  knowing that documents are valid wrt a DTD  $d$ . It strongly relies on the following developments.

- Chain-based independence for queries and updates, a static notion, is the foundation of our algorithm: starting from the set of all possible chains associated with a DTD, our inference system extracts subsets of chains for the query and the update, capturing the navigation through valid documents made by the evaluation of the query and the update, respectively. Our inference system is cautiously specified for dealing with all XPath axes. Chain-based independence is the result of the absence of overlapping pair of chains for the query and the update. Our inference system is formally proved to be *sound*.
- A major step of our work concerns recursive schemas, for which chain-based independence analysis may cripplingly involve to deal with an infinite number of chains. Our technique enabling the restriction of the analysis to *finite* subsets of query and update chain is a key contribution, and the core of our algorithm is the resulting finite analysis. It is proved to be *equivalent* to the infinite analysis
- Aiming at designing a tractable analysis, we show that by using a DAG-based representation of inferred chains, the finite analysis can run in polynomial space and time. Our technique has been carefully implemented, and extensive tests have been performed to validate our claim of precision and efficiency. Concerning precision, our results show that our technique outperforms [BC09a] to a large extent.

This work has been published in the *International Conference of Very Large Databases 2012* [BTCU12], while preliminary versions have been presented in the *26ème Journée des Bases de Données Avancées 2010* [BTCU10a] and in the *International Formal Methods Workshop 2010* [BTCU10b].

The thesis is organized as follows. Chapter 2 introduces basic definitions about XML, schemas, queries and updates. Chapter 3 revises the state-of-the-art for detecting XML

query-update independence, and provides motivations for the work here presented. Chapter 4 presents our static independence analysis based on schemas and chains. Chapter 5 describes how the chain analysis can be efficiently implemented. Chapter 6 presents experiments validating our method. Chapter 7 discusses extensions, while Chapter 8 draws conclusions and discusses future perspectives.

## Chapter 2

# Preliminaries

In this chapter we define the formal framework on top of which we build our work. We begin by presenting XML and DTDs, as the data model and the type language of semi-structured databases. Then, we present the XML query and update languages that we consider, together with their formal semantics, which is used to prove the correctness of our analysis. We conclude by defining the XML query-update independence problem.

### 2.1 XML

XML [BPSM<sup>+</sup>06] is a W3C-recommended language for representing data on the Web. XML has been designed to facilitate information exchange across different systems, thus overcoming the limits of its predecessor HTML, that was purposely conceived for displaying. XML is capable of wrapping many different kinds of data, because the data is represented in a textual form, organized hierarchically in a tree, and the document structure is self-describing by means of markup tags.

An XML Document is written as a unicode text with markup tags and other meta-information representing all document nodes. The most important ones we consider are element and textual nodes. As an example, consider the XML document  $t_1$  in Figure 2.1. It is a fragment of a bibliography, containing a single book, with its own title and author. Textual nodes here are written as the text they represent, for instance “B. Pierce”. Element nodes are denoted by pairs of open and closing markup-tags, for instance `<book></book>`, and are possibly nested inside other elements.

An XML document in its textual form must be *well-formed*. This essentially means that it defines a tree structure, the labels of open and closing tags match, and are also properly nested, so that the closing tag of a node appears always before the close-tag of its parent. The document  $t_1$  is well-formed.

The abstract model underlying XML documents is tree-based. In its textual form, an XML document looks like the linearization of a tree. In this work, we represent an

$$\begin{array}{l}
\text{<bib>} \\
\text{  <book>} \\
\text{    <title>} \\
\text{      Types and Programming Languages} \\
\text{  </title>} \\
\text{  <author>} \\
\text{    <name>B. Pierce</name>} \\
\text{  </author>} \\
\text{</book>} \\
\text{</bib>}
\end{array}
\quad
\sigma = \begin{cases} l_{t_1} \leftarrow \text{bib}[l_1] \\ l_1 \leftarrow \text{book}[l_2, l_3] \\ l_2 \leftarrow \text{title}[l'_2] \\ l'_2 \leftarrow \text{"Types and ..."} \\ l_3 \leftarrow \text{author}[l'_3] \\ l'_3 \leftarrow \text{name}[l''_3] \\ l''_3 \leftarrow \text{"B. Pierce"} \end{cases}$$

Figure 2.1: Document  $t_1$  and store  $(\sigma, l_{t_1})$ 

XML document as a *store*  $\sigma$ , which is an environment associating each *node location*  $l$  with either an element node  $\mathbf{a}[L]$  or a text node “**txt**”. In  $\mathbf{a}[L]$ ,  $\mathbf{a}$  is the element *tag* (or label), while  $L = (l_1, \dots, l_n)$  is the ordered sequence of *children* locations in  $\sigma$ . We denote by  $\text{dom}(\sigma)$  the set of locations of  $\sigma$ . A tree is a pair

$$t = (\sigma, l_t)$$

where  $l_t$  is the root location. With a little abuse of notation, we denote by  $\text{dom}(t)$  the domain of  $t$ , which consists of  $l_t$  together with all locations connected to  $l_t$ . We denote by  $\sigma @ l$  the *tree*  $t$  in  $\sigma$  rooted at  $l$ . We denote by  $t \sim t'$  two isomorphic trees  $t$  and  $t'$ ; they may differ only in terms of names of node locations. The empty sequence of locations is denoted by  $()$ . The concatenation of two sequences of locations  $L$  and  $L'$  is denoted by  $L \cdot L'$ , and defined such that  $L \cdot () = () \cdot L = L$ .

In Figure 2.1 we draw the store associated to the XML document  $t_1$ ; this store is a driving example for most of the definitions provided in the chapter. As already mentioned, each element tag of  $t_1$  corresponds to a single node in the store, and vice versa. As a subtlety, notice that each node is referenced by a location. A location can be simply seen as the identifier of the node.

The XML data model has many features, but in this work we consider only element and textual nodes, since they are the most relevant ones. Most of the missing features, e.g., string attribute nodes, can be straightforwardly encoded in our model. We consider tree instances without document-node [BPSM<sup>+</sup>06] (an unlabeled node which is the parent of the root). Other mechanisms such as key referencing (ID/IDREF) require to extend our model, by adding further constraints. We discuss extensions in Chapter 7.

## 2.2 Schema Languages

The concept of *schema* is standard in databases: it is a set of constraints defining a collection of data. XML documents are self-describing, hence a schema is an optional feature of XML databases. Nevertheless, any data collection is likely to follow a structure, and schemas are a very powerful means to express structural constraints. In all the context where a schema is available, it can be exploited to ensure *safe* and *efficient* XML processing [CGS11].

Many formalisms for writing XML schemas have been proposed, with the notable distinction of Document Type Definition (DTD) and XML Schema (XSD) as W3C standards, and RelaxNG as OASIS standard [CM01]. In this work we will focus on the widely used DTDs. DTDs embody all of the main features of XML schema languages, whilst still allowing to keep the formalization concise. The extension of DTDs to more powerful languages such as XSD, by means of *extended DTDs* [PV00] will be also discussed.

A DTD essentially consists of a list of declarations of elements. As an example, consider the schema `bibliography.dtd` reported in Figure 2.2. The declaration of elements follows a specific syntax like `<!ELEMENT bib (book*) >` stating that an element tag is defined, that the name of the element is `bib`, and that the content of the element is modeled by the regular expression `book*`. The reserved word `#PCDATA` denotes the string type.

DTDs can be formally defined along the line of [MLMK05].

**Definition 2.2.1** (DTD). *A DTD is a 3-tuple  $(d, s, \Sigma)$  where*

- $\Sigma$  is a finite alphabet for element tags, denoted by  $a, a_2, a'$ ;
- $s \in \Sigma$  is a non-recursive root-type;
- $d$  is a function from  $\Sigma$  to regular expressions over  $\Sigma \cup \{\text{String}\}$ , where **String** denotes the string base-type.

For convenience, in this work we assume that the root-type of a DTD is non-recursive. This implies that the root-element tag is never re-used inside a valid document, as happens in most practical scenarios. As we will discuss in Chapter 4, Section 4.4.2, this assumption is made for the sake of formalizations of the finite analysis, and is not restrictive. From now on we will use only the  $d$  component to specify a DTD, and we will write  $\Sigma_S$  to denote the set  $\Sigma \cup \{\text{String}\}$ .

In Figure 2.2 we show how a DTD is encoded in our formalism. The verbose declarations are denoted in our formalism in a simpler way, by associating each element tag with a regular expression modeling the element tag content.

An XML database compliant with all of the structural constraints defined by a schema is said to be *valid* against the schema.



<ELEMENT bib (book*)>	bib ← book*
<ELEMENT book (title,author+)>	book ← title,author+
<ELEMENT title (#PCDATA)>	title ← String
<ELEMENT author (name)>	author ← name
<ELEMENT name (#PCDATA)>	name ← String

Figure 2.2: DTD bibliography.dtd and the schema model  $d_1$ 

**Definition 2.2.2** (Validity). A tree  $t = (\sigma, l_t)$  is valid wrt a DTD  $(d, s, \Sigma)$ , denoted by  $t \in d$ , if and only if there exists a mapping  $\lambda: \text{dom}(t) \mapsto \Sigma_S$  such that

- $\lambda(l_t) = s_d$
- $\lambda(l) = \text{String}$  implies that  $\sigma(l)$  is a text node
- $\lambda(l) = a$  implies that  $\sigma(l)$  is an element-node  $a[L]$ , and that the word  $\lambda(L)$  is generated by the regular expression  $d(a)$ , i.e.,  $\lambda(L) \in \text{Lang}(d(a))$

where  $\text{Lang}(r)$  denotes the language generated by the regular expression  $r$ .

As an example, the document  $t_1$  in Figure 2.1 is valid wrt the DTD  $d_1$  in Figure 2.2, because there exists a mapping  $\lambda$  such that

- for the root  $l_{t_1}$ , we have  $\lambda(l_{t_1}) = \text{bib}$
- for all textual nodes, namely  $l'_2$  and  $l''_3$ , we have  $\lambda(l'_2) = \lambda(l''_3) = \text{String} \in \text{Lang}(\text{String})$
- for all the sequences of children locations, namely  $(l_1)$ ,  $(l_2, l_3)$  and  $(l'_3)$  we have

$$\begin{aligned}
 \lambda(l_1) &= \text{book} && \in \text{Lang}(\text{book}*) \\
 \lambda(l_2, l_3) &= \text{title author} && \in \text{Lang}(\text{title,author}+) \\
 \lambda(l'_3) &= \text{name} && \in \text{Lang}(\text{name})
 \end{aligned}$$

A recursive DTD  $d$  may not have any XML tree  $t$  such that  $t \in d$ . This is because the content description of some element  $a$  in  $d$  is *non-terminating*, i.e., there exists no finite subtree rooted at  $a$  that satisfies all constraints of  $d$ . A DTD  $d$  may also feature the content description of some elements that are *unreachable* from the root. No subtree of a valid tree is typed with an unreachable type definition, because there is a unique root type, and each element type definition is unique.

One can determine whether an element type  $a$  in  $d$  is non-terminating or unreachable in linear time (in the size of  $d$ ). These problems can be reduced to the emptiness problem for context-free grammars, which can be solved in linear time [HU00]. To simplify the discussion, in the sequel we assume that all element types in a DTD are terminating, and that all types in  $\Sigma_S$  are reachable from the root.

DTDs capture all of the main features of XML schema languages. The principal restriction a DTD imposes is that each element tag definition must be *unique*. This is partially relaxed in XSDs, where a tag can have multiple type definitions, but at a specific condition: two type definitions for the same element tag, are never used in the same element content description. This is the *single type* restriction.

To illustrate, a legal fragment of XSD containing two element tag definitions that cannot be expressed by DTDs, is the following one.

```
<xs:element name = "author">
  <xs:complexType>
    <xs:element name = "name" type = "AuthorName"/>
  </xs:complexType>
</xs:element>

<xs:element name = "publisher">
  <xs:complexType>
    <xs:element name = "name" type = "PublisherName"/>
  </xs:complexType>
</xs:element>
```

Notice that XSDs and DTDs have different syntax (XSDs are written as XML documents). The command `xs:element` defines a new element tag, while `xs:complexType` is used here to define the type structure of the element tag. In the example, two element tags are defined, namely `author` and `publisher`. Both have a common child element, `name`, specifying the name of the author and that of the publisher.

In DTDs, the definition of `name` must be unique while, in XSDs, we can use multiple type definitions. In this case, the schema refers to distinct type definitions `AuthorName` and `PublisherName` (defined elsewhere in the schema) depending on whether `name` is nested in `author` and `publisher`, respectively. What is powerful here is that the two types can define completely different structures of the `name` element.

A fragment of XSD that does not respect the single type restriction is the following.

```
<xs:element name = "person">
  <xs:complexType>
    <xs:choice>
      <xs:element name = "name" type = "AuthorName"/>
      <xs:element name = "name" type = "PublisherName"/>
    </xs:choice>
  </xs:complexType>
</xs:element>
```

Here `xs:choice` says that only one between the elements in the list can be chosen. In particular, one between the two `name` element definitions can be used. This definition is

illegal in XML Schema, since it is illegal to have two element types with the same name (i.e., `name`) but different types (i.e., `AuthorName`, `PublisherName`) in the same regular expression. Indeed, this breaks the single type restriction, and it is not expressible in DTDs nor in XSDs. Nevertheless it is expressible in RelaxNG.

A formalism that goes beyond that of DTDs and captures the expressivity of XML Schemas is that of *Extended DTDs* (EDTDs).

**Definition 2.2.3** (Extended DTD and Validation [PV00]). *An Extended DTD is a triple  $(\Sigma, d', \mu)$  where*

- $d'$  is a DTD defined over the tag alphabet  $\Sigma'$
- $\mu : \Sigma'_S \mapsto \Sigma_S$  is a function, that is the identity on input `String`

*A tree  $\mathcal{A}$  over  $\Sigma$  is valid wrt the extended DTD  $(\Sigma, d', \mu)$  if and only if there exists a tree  $t'$  over  $\Sigma'$  such that  $\mu(t') = t$  and  $t'$  is valid wrt the DTD  $d'$ .*

The crux of the definition is to interpret the two tag alphabets  $\Sigma$  and  $\Sigma'$  in a way such that  $\Sigma$  represents all element tags found in the tree instance, while  $\Sigma'$  contains a different tag for each type definition in the schema. In case of a tag with multiple definitions, such as `name` in the last example, we would have  $\Sigma \supseteq \{\text{name}\}$  and  $\Sigma' \supseteq \{\text{AuthorName}, \text{PublisherName}\}$  with  $\mu(\text{PublisherName}) = \text{name}$  and  $\mu(\text{AuthorName}) = \text{name}$ . The validation of the tree  $t$  over  $\Sigma$  wrt the extended DTD is then defined in terms of validation of a tree  $t'$  over  $\Sigma'$  wrt a DTD over a “type” alphabet  $\Sigma'$ .

It is worth noticing that EDTDs do not make any assumption on the use of types in the regular expressions, and hence they are strictly more expressive than XML Schemas, and roughly correspond to RelaxNG (without interleaving and counting operators). In this work, we will develop a static analysis relying on DTDs, and in Chapter 7 we will show that it can be lifted to EDTDs very easily.

XML *types* are founded on regular tree grammars [HVP05]. Indeed, a schema language (e.g., DTD) roughly corresponds to a class of regular tree grammars [MLMK05]. This is needed for designing typed XML query and update language, where types of input and output expressions can be checked and type safety possibly ensured.

However, our ultimate goal is to develop an *access analysis* based on schema informations, which is a problem radically different from that of type-checking. To avoid any confusion, in this work we adopt the following definition of XML type.

**Definition 2.2.4** (XML Type). *Let  $d$  be a DTD, a type  $\tau \in \Sigma_S$  stands for the set of nodes  $\mathcal{A}$  such that*

- *if  $\tau$  is the label  $a$  then  $\mathcal{A}$  is the set of all nodes labeled with  $a$  belonging to some tree  $t \in d$*
- *if  $\tau$  is `String` then  $\mathcal{A}$  is the set of all textual nodes belonging to some tree  $t \in d$*

From this definition, it follows that the *tag* of each element definition in a DTD is a type. As an example, consider the bibliographic DTD  $d_1$ . The type **author** stands for all nodes labeled as **author** belonging to a tree valid wrt the DTD  $d_1$ .

A DTD induces a type dependency graph, defined as follows.

**Definition 2.2.5** (Dependency graph). *Let  $d$  be a DTD, the dependency graph of  $d$ , denoted by  $\mathcal{G}_d$ , is a directed graph defined as  $\mathcal{G}_d = (\Sigma_S, \Rightarrow_d)$ .*

The vertex-set of  $\mathcal{G}_d$  coincides with  $\Sigma_S$ , which is the tag-alphabet of the schema together with the **String** type. The edge-set of  $\mathcal{G}_d$  is  $\Rightarrow_d$ , therefore two (immediately) reachable types in  $\Sigma_S$  are two adjacent vertex in  $\mathcal{G}_d$ . The root of the graph is  $s_d$ . The graph is connected because all types in  $\Sigma_S$  are reachable from  $s_d$ . A DTD  $d$  is recursive if and only if  $\mathcal{G}_d$  has a cycle.

**Proposition 2.2.6.** *Let  $\mathcal{G}_d$  be a type dependency graph with  $n$  vertices and  $m$  edges.*

$$i) \quad n-1 \leq m \leq n(n-1)$$

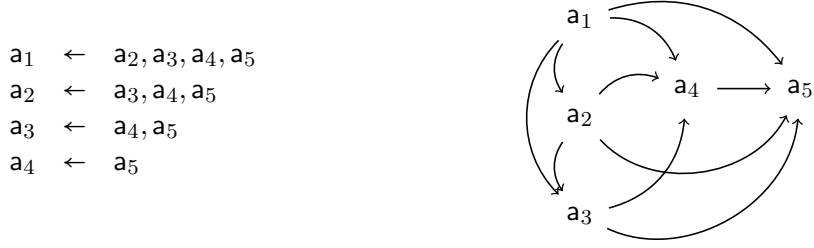
$$ii) \quad \text{if } m > n(n-1)/2 \text{ then } \mathcal{G}_d \text{ is cyclic}$$

*Proof.* *i)* If  $\mathcal{G}_d$  has  $n$  vertices, the number of edges in  $\mathcal{G}_d$  varies from  $n-1$ , because the graph is connected, to  $n(n-1)$ , because all types can be recursive, except for the root that we assumed to be non-recursive. *ii)* An *acyclic* dependency graph of  $n$  vertices with the maximum number of edges, is such that only one type (the root) is connected with at most  $n-1$  types, only one type is connected with at most with  $n-2$  types, and so on until a single type with no children (e.g., **String**). This is  $\sum_{i=1}^n (i-1) = n(n-1)/2$ .  $\square$

We call *saturated* a DTD  $d$  that maximizes the number of edges in the dependency graph  $\mathcal{G}_d$ . Saturated DTDs are defined as follows.

**Definition 2.2.7.** (*Saturated DTDs*) *A recursive DTD  $d$  is said to be saturated if  $\mathcal{G}_d$  has  $n$  vertexes and  $n(n-1)$  edges. A non-recursive DTD  $d$  is said to be saturated if  $\mathcal{G}_d$  has  $n$  vertexes and  $n(n-1)/2$  edges.*

An example of a saturated and non-recursive DTD and relative dependency graph is depicted in Figure 2.3. The DTD  $d_2$  features 5 elements. The root type  $a_1$  is defined in terms of all other types  $a_2, \dots, a_5$ , the type  $a_2$  is defined in terms of all other types except the root and itself, namely  $a_3, \dots, a_5$ , and so on until the singleton type  $a_5$  that has empty content. Note that no edge can be added in  $\mathcal{G}_{d_2}$  without introducing recursion.

Figure 2.3: DTD  $d_2$  and saturated type-dependency graph  $\mathcal{G}_{d_2}$ 

## Inference and Derivations

Throughout this dissertation we will describe several inference systems for checking properties of XML queries and updates. Inference systems have been nicely formalized in [Gra03].

**Definition 2.2.8** (Inference System). *Let  $\mathbf{J}$  be a set of predicates called judgments. An inference system  $\mathbf{I}$  is a set of rules defined over judgments. A rule  $R$  is of the form*

$$\frac{\Upsilon}{\beta}$$

where  $\Upsilon \subseteq \mathbf{J}$  is a set of premises and  $\beta \in \mathbf{J}$  is a judgment called the goal of the rule.

An inference system proves properties (ground judgments) by repeated backward application of inference rules, in Prolog style. The whole proof is called derivation tree, that is a tree whose nodes are judgements and edges denote the application of inference rules. In particular, the root of the derivation tree is the property one aims to prove and the leafs are axioms. A set of rules is said to be deterministic if for each judgement created, at most one rule can be applied. This property is necessary to ease the derivation of an algorithm from the rules. In this dissertation, we will present only deterministic inference systems.

## 2.3 Query and Update Languages

XQuery is the XML query language standardized by the W3C [SCF<sup>+</sup>07]. XQuery stands to XML databases as SQL stands to relational databases. Nevertheless, the two languages are deeply different, because they are conceived for different data-models (trees vs. relations), and different deployments (the web vs. mainframes applications).

XQuery and SQL are both declarative languages. This is a crucial requirement for any query language in order to ensure logical and physical data independence. XQuery is a Turing-complete language, while SQL (without advanced features) is not. This just means that XQuery provides sufficient - perhaps, superabundant - expressive query power.

XQuery and SQL are both typed languages, this means that a type representing the result of an expression can be inferred. XQuery is a functional language, while SQL is not. This means that XQuery expressions are functions that can be nested or composed, and that the evaluation of each expression always returns a value with no side effect.

To illustrate, consider the following XQuery expression, where  $y$  is a free variable bound to the root of a bibliographic XML document and  $x$  is a bound variable of the query.

```
for x in y//book return (x/title, x/author)
```

The expression iterates over book nodes in the document. At each iteration step, it binds the current node to the variable  $x$ , and then it returns a *sequence* of title and author of the book.

Being declarative, XQuery allows to retrieve all books in the document, by invoking the XPath navigation  $y//book$ , independently from how they are effectively stored and retrieved.

Being functional, expressions can be composed in a sequence, to form the return clause, or nested to form the whole iterative construct.

Being typed, it is possible to infer the regular expression type  $(title*, author*)^*$ , to capture the fact that the result of the query is a possibly empty sequence of title and author lists.

XQuery is built on top of XPath. XPath is a language designed for navigating tree structures, along the axes of ancestors, descendants and siblings of a node. XQuery wraps-up XPath expressions, by introducing programming-language style constructs such as iteration, let binding, conditional expression, sequence construction, element construction, type-switching, etc.

The W3C working group has published a rich documentation on the XML query and update languages [SCF<sup>+</sup>07, RCD<sup>+</sup>11]. There is a large amount of details to consider in a rich language such as XQuery, hence it is a common practice to focus on a subset of the language that capture the most interesting constructs, so as to make the formal development intelligible yet effective.

In this work we deal with a large fragment of XQuery considered also in related approaches [BC09a, BC10], which is defined by the grammar of Table 2.1.

We comment on the main constructs of the language.

(EMPTY) The empty query is an expression performing no operation.

(CONCATENATION) A concatenation of queries is used to build new sequences of trees.

$q ::= ()$	(EMPTY)
$  q, q$	(CONCATENATION)
$  \langle a \rangle q \langle /a \rangle$	(ELEMENTNODE)
$  \text{"txt"}$	(TEXTNODE)
$  x/\text{axis}::\phi$	(STEP)
$  \text{for } x \text{ in } q \text{ return } q$	(FOR)
$  \text{let } x := q \text{ return } q$	(LET)
$  \text{if } (q) \text{ then } q \text{ else } q$	(IF)
$\text{axis} ::= \text{self}$	(AXES)
$  \text{child}$	
$  \text{descendant}$	
$  \text{descendant-or-self}$	
$  \text{parent}$	
$  \text{ancestor}$	
$  \text{ancestor-or-self}$	
$  \text{preceding-sibling}$	
$  \text{following-sibling}$	
$\phi ::= a$	(NODE TESTS)
$  \text{text}()$	
$  \text{node}()$	

Table 2.1: Query Language

(NODE ELEMENT AND NODE TEXT CONSTRUCTION) The importance of node element and text constructions stems from the fact that they are the only query expression that extend the input store with new data. Their use is twofold. They can be used in order to re-structure nodes of the base data, and then present it to the user (e.g. output a node changing the order of its children). They can be used by update expressions, in order to create new forests to be added to the store (e.g., insert a fresh textual node “**txt**” in a given position).

(XPATH STEPS) Navigational capabilities are at the heart of all query languages for semi-structured and, in general, for hierarchical data. Our language grammar includes downward, upward and sibling XPath navigational axis. This is a large fragment, considered also by related studies such as [BC09a, BC10].

Preceding and following axes, are not included. They can be easily encoded by standard rewriting. As we will show later, this can be done in a both correct and complete way. Therefore, the language we consider simulates all XPath navigational axes.

(ITERATION, LET BINDING AND CONDITIONAL) Iteration (FOR), let binding (LET) and conditional expression (IF) represent the most powerful constructs of the language, in terms of expressivity. This set of constructs subsumes the standard FLWR XQuery fragment [SCF<sup>+</sup>07], that features the where clause instead of the conditional if. Indeed, while the where clause nested in a let or for expression can be simulated by using an if-then-else expression, the converse is not always true since the where clause can express only an if-then condition. Therefore the language is strictly more expressive than its FLWR version.

For the sake of presentation, we adopt some minor syntactic deviation from the W3C specification [SCF<sup>+</sup>07]. In XQuery variables are always prefixed by \$, in order to avoid ambiguity during the parsing phase. However, in this work variables will be always clear from the context and thus, we avoid to prefix them with a reserved symbol. Also, element construction is denoted by  $\langle a \rangle \mathbf{q} \langle /a \rangle$ , without enclosing the inner query  $\mathbf{q}$  into brackets, as defined in [SCF<sup>+</sup>07].

XPath node tests are denoted by  $\phi$ . In the grammar,  $a$  stands for a tag name. We will use  $/\phi$  as a shortcut for  $/\mathbf{child}::\phi$  while  $//\phi$  is a (strict) descendant navigation that, according to the W3C, stands for

$$/\mathbf{descendant-or-self}::\mathbf{node}()/\mathbf{child}::\phi$$

In the following we will use **step** to indicate an XPath step. The node filter **node()** is often denoted by the wildcard  $*$ . XPath expressions of the form

$$\mathbf{x}/\mathbf{step}_1/\mathbf{step}_2/\dots/\mathbf{step}_n$$



<code>u ::= ()</code>	(EMPTY)
<code>  u, u</code>	(CONCATENATION)
<code>  for x in q return u</code>	(FOR)
<code>  let x := q return u</code>	(LET)
<code>  if (q) then u<sub>1</sub> else u<sub>2</sub></code>	(IF)
<code>  delete q<sub>0</sub></code>	(DELETE)
<code>  rename q<sub>0</sub> as a</code>	(RENAME)
<code>  insert q pos q<sub>0</sub></code>	(INSERT)
<code>  replace q<sub>0</sub> with q</code>	(REPLACE)
<code>pos ::= into (as first   as last)?</code>	(POSITIONS)
<code>  after</code>	
<code>  before</code>	

Table 2.2: Update Language

are supported by the grammar by mean of iteration, as follows.

```

for x1 in x/step1 return
  for x2 in x1/step2 return
    ⋮
    for xn in xn-1/stepn return xn/self :: node()

```

Recently the W3C standardized an update language for XML called XQuery Update Facility [RCD<sup>+</sup>11]. This includes a set of commands that can be used to modify instances of the XML data model, thus introducing *effects* to the query language.

As for queries, update commands are declarative and typed. This means that it is possible to modify an XML document, independently from how this is stored or retrieved, and also to compute a type describing the modified document. However, XML updates are not written in a functional way. They are not true functions with a return value that can be arbitrarily composed or nested.

An example of XML update is the following.

```
insert <book/> into x/bib
```

which inserts a fresh node labeled as *book* below the root labeled with *bib*.

XQuery Update Facility is built on XQuery and XPath, that are indeed used to fill the arguments of update commands.

The subset of XQuery Update Facility [RCD<sup>+</sup>11] which we consider is defined in Table 2.2. It is essentially composed of two parts, one of programming constructs shared with the query language, and one of commands modifying the input store.

(BASIC PROGRAMMING) Like for queries, updates can be composed sequentially or can be composed by means of let/for statements, where only the return part can contain update operations. The conditional expression allows to apply different updates, depending on the boolean value of a conditional query.

(SIDE-EFFECTS) In this work we consider all main constructs of the XQuery Update Facility, namely node deletion, renaming, insertion and replacement. We consider only replace update affecting *node* locations. The extension to the replace node-value command is straightforward <sup>1</sup>, as well as that to the remaining update commands. Other extensions to the update language are discussed in Chapter 7.

## 2.4 Formal Semantics

In this section we give the formal semantics of XQuery and XQuery Update Facility. This will be used in the proofs to show the correctness of our independence analysis. We begin by defining the semantics of XPath steps as they are the building blocks of both languages.

### 2.4.1 Semantics of XPath Steps

**Tree relations** The semantics of axis navigation and text filtering is defined in [DFF<sup>+</sup>10]. If  $\sigma$  is a store, the binary parent-child relation among store locations is denoted by  $Child_\sigma \subseteq dom(\sigma) \times dom(\sigma)$ , and  $(l, l') \in Child_\sigma$  when  $l = a[L \cdot l' \cdot L'] \in \sigma$ . We denote by  $Parent_\sigma$  the inverse of  $Child_\sigma$ , by  $Descendant_\sigma$  the transitive closure of  $Child_\sigma$ , and by  $Ancestor_\sigma$  the inverse of  $Descendant_\sigma$ . Finally,  $Self_\sigma$  denotes the identity function on  $dom(\sigma)$ . The binary relation that holds between a location and its immediate-right-sibling is denoted by  $NextSibling_\sigma \subseteq dom(\sigma) \times dom(\sigma)$ . We have that  $(l, l') \in NextSibling_\sigma$  when  $l_a \leftarrow a[L \cdot l \cdot l' \cdot L'] \in \sigma$ , for some  $a$ . We denote by  $FollowingSibling_\sigma$  the transitive closure of  $NextSibling_\sigma$ , and by  $PrecedingSibling_\sigma$  the inverse of  $FollowingSibling_\sigma$ . If  $R$  is one of the binary relations defined above, provided a location  $l$ , we write  $R(l)$  for the set of locations  $l'$  such that  $(l, l') \in R$ . For instance  $Child_\sigma(l)$  denotes the set of children of  $l$  and  $FollowingSibling_\sigma(l)$  the set of following sibling nodes of  $l$ . Of course, it holds  $Self_\sigma(l) = \{ l \}$ . Finally  $docOrder_\sigma \subseteq dom(\sigma) \times dom(\sigma)$  denotes the partial document-order among store locations, that is total for locations belonging to the same tree, but is undefined for locations belonging to distinct trees of the store. If  $R_\sigma$  is a relation over  $\sigma$  defined above, we denote by  $R_t \subseteq R_\sigma$  its subset concerning all and only the locations belonging to the tree  $t = (\sigma, l_t)$ .

---

<sup>1</sup>replace value of  $q_0$  with  $q$  becomes `let x :=  $q_0$  return replace x/text() with  $q$`

**Axis Semantics** The semantics of an XPath step is given by the composition of the semantics of a navigational axis **axis** and that of a test filtering  $\phi$ . The semantics of a navigation axis **axis** wrt a store  $\sigma$  and an input location  $l$ , is denoted by  $\llbracket \mathbf{axis} \rrbracket_{\sigma}^l$  and defined as follows.

$$\begin{aligned}
\llbracket \mathbf{self} \rrbracket_{\sigma}^l &\stackrel{def}{=} Self_{\sigma}(l) \\
\llbracket \mathbf{child} \rrbracket_{\sigma}^l &\stackrel{def}{=} Child_{\sigma}(l) \\
\llbracket \mathbf{descendant} \rrbracket_{\sigma}^l &\stackrel{def}{=} Descendant_{\sigma}(l) \\
\llbracket \mathbf{descendant-or-self} \rrbracket_{\sigma}^l &\stackrel{def}{=} Descendant_{\sigma}(l) \cup Self_{\sigma}(l) \\
\llbracket \mathbf{parent} \rrbracket_{\sigma}^l &\stackrel{def}{=} Parent_{\sigma}(l) \\
\llbracket \mathbf{ancestor} \rrbracket_{\sigma}^l &\stackrel{def}{=} Ancestor_{\sigma}(l) \\
\llbracket \mathbf{ancestor-or-self} \rrbracket_{\sigma}^l &\stackrel{def}{=} Ancestor_{\sigma}(l) \cup Self_{\sigma}(l) \\
\llbracket \mathbf{following-sibling} \rrbracket_{\sigma}^l &\stackrel{def}{=} FollowingSibling_{\sigma}(l) \\
\llbracket \mathbf{preceding-sibling} \rrbracket_{\sigma}^l &\stackrel{def}{=} PrecedingSibling_{\sigma}(l)
\end{aligned}$$

The semantics of a test filtering wrt a store  $\sigma$  and an input location  $l$ , is denoted by  $\llbracket \phi \rrbracket_{\sigma}^l$  and defined by the following equations.

$$\begin{aligned}
\llbracket \mathbf{a} \rrbracket_{\sigma}^l &\stackrel{def}{=} \{ l \mid l \leftarrow \mathbf{a}[L] \in \sigma \} \\
\llbracket \mathbf{text}() \rrbracket_{\sigma}^l &\stackrel{def}{=} \{ l \mid l \leftarrow \mathbf{“txt”} \in \sigma \} \\
\llbracket \mathbf{node}() \rrbracket_{\sigma}^l &\stackrel{def}{=} \{ l \mid l \in \sigma \}
\end{aligned}$$

Notice that this rules output either input location, when it satisfies the filtering condition, or the empty set.

The following examples illustrate the above definitions. Consider the store  $\sigma$  in Figure 2.1.

$$\begin{aligned}
\llbracket \mathbf{descendant-or-self} \rrbracket_{\sigma}^{l_{t1}} &= dom(\sigma) \\
\llbracket \mathbf{following-sibling} \rrbracket_{\sigma}^{l_2} &= \{ l_3 \} \\
\llbracket \mathbf{text}() \rrbracket_{\sigma}^{l'_2} &= \{ l'_2 \} \\
\llbracket \mathbf{author} \rrbracket_{\sigma}^{l'_2} &= \emptyset
\end{aligned}$$

An XPath axis **axis** is said to be downward if it is either **child**, **descendant** or **descendant-or-self**, backward if it is **parent**, **ancestor** or **ancestor-or-self**, and horizontal if it is **preceding-sibling** or **following-sibling**.

**Step Semantics** An XPath step is evaluated by composing axis navigation and step filtering. An axis navigation results in a set of locations  $\mathcal{L}$ . This set is explicitly casted into a sequence  $L$  that respects document order. We model this by using an abstract predicate  $order_\sigma(\mathcal{L}) = L$  that holds iff (i)  $l \in \mathcal{L} \iff l \in L$  and (ii)  $L = (L_1 \cdot l \cdot l' \cdot L_2) \implies (l, l') \in docOrder_\sigma$ . Finally,  $L$  is filtered according to  $\phi$ . Provided that  $\llbracket \mathbf{axis} \rrbracket_\sigma^l = \mathcal{L}$  and  $order_\sigma(\mathcal{L}) = (l_1, \dots, l_n)$  step semantics is defined as follows.

$$\llbracket \mathbf{axis} :: \phi \rrbracket_\sigma^l \stackrel{def}{=} (\llbracket \phi \rrbracket_\sigma^{l_1}, \llbracket \phi \rrbracket_\sigma^{l_2}, \dots, \llbracket \phi \rrbracket_\sigma^{l_n})$$

When  $\llbracket \phi \rrbracket_\sigma^l = \emptyset$  it is treated as the empty sequence.

Notice that this definition slightly differs from the one in [SCF<sup>+</sup>07] where sequences obtained by the evaluation of backward and preceding-sibling axes are returned in inverse document order. Of course, this detail has no impact on our static analysis.

The following examples illustrate the above definition. Consider the store  $\sigma$  in Figure 2.1.

$$\begin{aligned} \llbracket \mathbf{descendant-or-self} :: \mathbf{node}() \rrbracket_\sigma^{l_{t_1}} &= dom(\sigma) \\ \llbracket \mathbf{descendant-or-self} :: \mathbf{text}() \rrbracket_\sigma^{l_{t_1}} &= (l'_2, l''_3) \\ \llbracket \mathbf{following-sibling} :: \mathbf{author} \rrbracket_\sigma^{l_2} &= (l_3) \end{aligned}$$

Notice that  $\llbracket \mathbf{axis} :: \phi \rrbracket_\sigma^l$  produces a sequence of nodes without duplicates.

### 2.4.2 XQuery Semantics

Query semantics is specified in [DFF<sup>+</sup>10], while we are highly inspired by the succinct and elegant formalization made in [BC09b], from which we borrow some notions that are needed for our own presentation.

Query semantics is denoted by the following judgment

$$\sigma, \gamma \models \mathbf{q} \Rightarrow \sigma', L$$

meaning that the execution of the query  $\mathbf{q}$  over  $\sigma$  outputs a sequence of locations  $L$  and a new store  $\sigma'$ , including  $\sigma$  plus new elements built by  $\mathbf{q}$ . The dynamic environment  $\gamma$  binds each free variable of  $\mathbf{q}$  to a sequence of locations  $L$  in  $\sigma$ , this is denoted by  $\gamma[\mathbf{x} \mapsto L]$  or simply by  $\gamma[\mathbf{x} \mapsto l]$  when  $L = (l)$ .

Query semantics is reported in Table 2.3. In the rules we also make use of two auxiliary judgements  $\models^{\text{copy}}$  and  $\models^*$  that model the copy of a sequence of elements and the iteration over a sequence of elements. These are defined by rule (SQ-COPY) and (SQ-ITER), respectively.

$$\begin{array}{c}
\frac{}{\sigma, \gamma \models () \Rightarrow \sigma, ()} \text{ (SQ-EMPTY)} \\
\\
\frac{\sigma, \gamma \models q_1 \Rightarrow \sigma_2, L_1 \quad \sigma_2, \gamma \models q_2 \Rightarrow \sigma_3, L_2}{\sigma, \gamma \models q_1, q_2 \Rightarrow \sigma_3, L_1 \cdot L_2} \text{ (SQ-CONCAT)} \\
\\
\frac{\sigma, \gamma \models q \xRightarrow{\text{copy}} \sigma_2, L \quad l \notin \text{dom}(\sigma_2)}{\sigma, \gamma \models \langle a \rangle q \langle /a \rangle \Rightarrow \sigma_2[l := a[L]], l} \text{ (SQ-ELT)} \quad \frac{l \notin \text{dom}(\sigma)}{\sigma, \gamma \models \text{"txt"} \Rightarrow \sigma_2[l := \text{"txt"}], l} \text{ (SQ-TEXT)} \\
\\
\frac{\sigma, \gamma \models q \Rightarrow \sigma_0, L_0 \quad (L', \sigma_1) = \text{copy}(L_0, \sigma_0) \quad \sigma' = \sigma_0 \cup \sigma_1}{\sigma, \gamma \models q \xRightarrow{\text{copy}} \sigma', L'} \text{ (SQ-COPY)} \\
\\
\frac{\sigma, \gamma \models q_1 \Rightarrow \sigma_2, L \quad \sigma_2, \gamma[x \mapsto L] \models^* q_2 \Rightarrow \sigma_3, L'}{\sigma, \gamma \models \text{for } x \text{ in } q_1 \text{ return } q_2 \Rightarrow \sigma_3, L'} \text{ (SQ-FOR)} \\
\\
\frac{}{\sigma, \gamma[x \mapsto ()] \models^* q \Rightarrow \sigma, ()} \text{ (SQ-ITERBASE)} \quad \frac{\sigma, \gamma[x \mapsto l] \models q \Rightarrow \sigma_2, L_1 \quad \sigma_2, \gamma[x \mapsto L] \models^* q \Rightarrow \sigma_3, L_2}{\sigma, \gamma[x \mapsto l \cdot L] \models^* q \Rightarrow \sigma_3, L_1 \cdot L_2} \text{ (SQ-ITER)} \\
\\
\frac{\sigma, \gamma \models q_1 \Rightarrow \sigma_2, L \quad \sigma_2, \gamma[x \mapsto L] \models q_2 \Rightarrow \sigma_3, L'}{\sigma, \gamma \models \text{let } x := q_1 \text{ return } q_2 \Rightarrow \sigma_3, L'} \text{ (SQ-LET)} \\
\\
\frac{\sigma, \gamma \models q \Rightarrow \sigma_2, l \cdot L \quad \sigma_2, \gamma \models q_1 \Rightarrow \sigma_3, L_1}{\sigma, \gamma \models \text{if } (q) \text{ then } q_1 \text{ else } q_2 \Rightarrow \sigma_3, L_1} \text{ (SQ-IF1)} \\
\\
\frac{\sigma, \gamma \models q \Rightarrow \sigma_2, () \quad \sigma_2, \gamma \models q_2 \Rightarrow \sigma_3, L_2}{\sigma, \gamma \models \text{if } (q) \text{ then } q_1 \text{ else } q_2 \Rightarrow \sigma_3, L_2} \text{ (SQ-IF2)} \\
\\
\frac{\gamma(x) = l \cdot () \quad \llbracket [\text{axis} :: \phi] \rrbracket_\sigma^l = L}{\sigma, \gamma \models x/\text{axis} :: \phi \Rightarrow \sigma, L} \text{ (SQ-STEP)}
\end{array}$$

Table 2.3: Query evaluation rules

We illustrate the semantics of the language through some running examples. Consider for instance the following iteration, that returns all parents of textual nodes, evaluated over the tree  $t_1$  in Figure 2.1, wrt the environment  $\gamma = \{ \mathbf{x} \mapsto l_{t_1} \}$ .

$$\frac{\textcircled{1} \quad \sigma, \gamma \models \mathbf{x}/\text{descendant} :: \text{text}() \Rightarrow \sigma, L}{\sigma, \gamma \models \text{for } \mathbf{y} \text{ in } \mathbf{x}/\text{descendant} :: \text{text}() \text{ return } \mathbf{y}/\text{parent} :: \text{node}() \Rightarrow \sigma, L'} \quad \frac{\textcircled{2} \quad \sigma, \gamma[\mathbf{y} \mapsto L] \models^* \mathbf{y}/\text{parent} :: \text{node}() \Rightarrow \sigma, L'}{\sigma, \gamma \models \text{for } \mathbf{y} \text{ in } \mathbf{x}/\text{descendant} :: \text{text}() \text{ return } \mathbf{y}/\text{parent} :: \text{node}() \Rightarrow \sigma, L'}$$

The semantics of iteration is standard and modeled by rule (SQ-FOR). The resulting sequence of the left subexpression  $L$  is bound with variable  $\mathbf{y}$  and then used for iterating over the right subexpression. Notice that the evaluation does not change the store  $\sigma$ .

Let us explore more in detail the derivation.

- ① The left branch is derived by applying rule (SQ-STEP) as follows

$$\frac{\llbracket \text{descendant} :: \text{text}() \rrbracket_{\sigma}^{l_{t_1}} = (l'_2, l''_3)}{\sigma, \gamma \models \mathbf{x}/\text{descendant} :: \text{text}() \Rightarrow \sigma, (l'_2, l''_3)}$$

The result sequence is composed of all textual nodes that are descendants of the root.

- ② The right branch is derived by applying rule (SQ-ITER) as follows

$$\frac{\frac{\llbracket \text{parent} :: \text{node}() \rrbracket_{\sigma}^{l'_2} = (l_2)}{\sigma, \gamma[\mathbf{y} \mapsto l'_2] \models \mathbf{y}/\text{parent} :: \text{node}() \Rightarrow \sigma, (l_2)} \quad \frac{\llbracket \text{parent} :: \text{node}() \rrbracket_{\sigma}^{l'_3} = (l'_3)}{\sigma, \gamma[\mathbf{y} \mapsto l'_3] \models \mathbf{y}/\text{parent} :: \text{node}() \Rightarrow \sigma, (l'_3)}}{\sigma, \gamma[\mathbf{y} \mapsto (l'_2, l''_3)] \models^* \mathbf{y}/\text{parent} :: \text{node}() \Rightarrow \sigma, (l_2, l'_3)}$$

Here step evaluation is performed for each location bound to  $\mathbf{y}$ . The results of the evaluations are concatenated and then returned.

The construction of new elements is the only command that extends the input store during query evaluation.

To illustrate, consider the following query listing all books of tree  $t_1$  in Figure 2.1, wrt the environment  $\gamma = \{ \mathbf{x} \mapsto l_{t_1} \}$ .

$$\frac{\frac{\sigma, \gamma \models \mathbf{x}/\text{descendant} :: \text{book} \Rightarrow \sigma, (l_1)}{(l_c, \sigma_c) = \text{copy}(l_1, \sigma) \quad \sigma' = \sigma \cup \sigma_c}}{\sigma, \gamma \models \mathbf{x}/\text{descendant} :: \text{book} \xrightarrow{\text{copy}} \sigma', (l_c)} \quad l \notin \text{dom}(\sigma') \\ \sigma, \gamma \models \langle \text{list} \rangle \mathbf{x}/\text{descendant} :: \text{book} \langle / \text{list} \rangle \Rightarrow \sigma'[l := \text{list}[l_c]], l_c$$

In this case rule (SQ-ELT) is applied. First, the inner query selecting all books in the store is evaluated, resulting in the only location  $l_1$ . Second, the whole subtree rooted at  $l_1$  is copied to a fresh isomorphic tree  $\sigma_c@l_c$ , by means of rule (SQ-COPY). This operation is modeled here by using an abstract predicate  $copy(l_1, \sigma)$ . The output store  $\sigma'$  is the extension of  $\sigma$  with  $\sigma_c$ . Finally,  $l_c$  becomes the child of a fresh node location  $l$  representing the element node labeled with `list`.

**Rewriting of preceding and following axes** As said before, following and preceding axes are not directly dealt with by our system, but they can be expressed by means of rewriting. We can show that the rewriting is complete, in the sense that it gives the same result set as the standard semantics.

We show the case of the **following** axis, since **preceding** is analogous. Below  $\sigma$  is a fixed store, while  $l, l_a$  and  $l_f$  are locations belonging to  $\sigma$ . The W3C semantics for the axis is as follows.

$$\llbracket \text{following} :: \phi \rrbracket_{\sigma}^l \stackrel{\text{def}}{=} \text{order}_{\sigma}(\mathcal{L}) \quad (2.1)$$

where  $\mathcal{L}$  is

$$\begin{aligned} & \bigcup \{ l' \} \\ & l_a \in \llbracket \text{ancestor-or-self} :: \text{node}() \rrbracket_{\sigma}^l \\ & l_f \in \llbracket \text{following-sibling} :: \text{node}() \rrbracket_{\sigma}^{l_a} \\ & l' \in \llbracket \text{descendant-or-self} :: \phi \rrbracket_{\sigma}^{l_f} \end{aligned}$$

Step  $x/\text{following} :: \phi$  is rewritten in our language in the following way

$$x/\text{following} :: \phi \approx \text{q}_f(\phi)$$

where

$$\begin{aligned} \text{q}_f(\phi) = & \text{for } x_a \text{ in } x/\text{ancestor-or-self} :: \text{node}() \\ & \text{for } x_f \text{ in } x_a/\text{following-sibling} :: \text{node}() \\ & \text{return } x_f/\text{descendant} :: \phi \end{aligned}$$

Hence, provided that the judgment  $\sigma, x \mapsto l \models \text{q}_f(\phi) \Rightarrow \sigma, L$  holds, we have that  $L$  is

$$\begin{aligned} & \prod L_{(i,j)} \\ & (l_{a_1}, \dots, l_{a_n}) = \text{order}_{\sigma}(\llbracket \text{ancestor} :: \text{node}() \rrbracket_{\sigma}^l) \\ & (l_{f_{i,1}}, \dots, l_{f_{i,m}}) = \text{order}_{\sigma}(\llbracket \text{following-sibling} :: \text{node}() \rrbracket_{\sigma}^{l_{a_i}}) \\ & L_{(i,j)} = \text{order}_{\sigma}(\llbracket \text{descendant-or-self} :: \phi \rrbracket_{\sigma}^{l_{f_{i,j}}}) \end{aligned} \quad (2.2)$$

where  $\prod$  denotes sequence concatenation, with  $i = 1..n$  and  $j = 1..m$ .

We claim that a location is in the result of (2.1) if and only if it is in the result of (2.2). This holds since the resulting locations of all intermediate steps are the same. Furthermore, if one also wants to enforce that resulting sequences of both (2.1) and (2.2) are output in the same order (that is, document order), it is sufficient to impose that locations in the result of step  $\llbracket \text{ancestor} :: \text{node}() \rrbracket_\sigma^l$  in (2.2) are iterated in inverted document order, for instance by imposing  $i = n..1$ .

### 2.4.3 XQuery Update Semantics

In this section we formally define the semantics of update expressions.

As a convention in update expressions, we write  $q_0$  for the *target* expression that is retrieving the nodes in the input document that are target of the update (e.g. deleted nodes are target nodes). We simply write  $q$  for the *source* update subexpression, that is retrieving or producing elements to insert in the database by means of insert or replace commands. According to the W3C semantics [RCD<sup>+</sup>11] the target expression  $q_0$  is required to output a single node otherwise a run time error occurs.

In the W3C specification update evaluation is split into three phases.

1. Creation of an *update pending list* (UPL) of atomic update commands;
2. Execution of a sanity check on this list, and reorder of atomic updates;
3. Application of the UPL on the input store so as to update the document.

An update pending list  $\omega$  is a sequence of *atomic update commands*. An atomic update commands  $\iota$  in of the following form.

$$\begin{aligned} \iota \quad ::= \quad & \text{del}(l) \\ & | \text{ren}(l, a) \\ & | \text{ins}(L, \text{pos}, l) \\ & | \text{repl}(l, L) \end{aligned}$$

where  $l$  is the *target location*,  $L$  the sequence of roots of source elements to be inserted, and  $\text{pos}$  is an insertion position defined as before.

**Creation of the UPL** The creation of the UPL from an update  $u$  is denoted as follows

$$\sigma, \gamma \models u \Rightarrow \sigma_\omega, \omega$$

As usual,  $\gamma$  binds  $u$  free variables to locations in  $\sigma$  and the store  $\sigma_\omega$  that extends  $\sigma$  contains newly created locations potentially used in the UPL  $\omega$ .



$\frac{}{\sigma, \gamma \models () \Rightarrow \sigma, \epsilon}$	(SU-EMPTY)
$\frac{\sigma_1, \gamma \models u_1 \Rightarrow \sigma_2, \omega_1 \quad \sigma_2, \gamma \models u_2 \Rightarrow \sigma_3, \omega_2}{\sigma_1, \gamma \models u_1, u_2 \Rightarrow \sigma_3, \omega_1; \omega_2}$	(SU-CONCAT)
$\frac{\sigma_1, \gamma \models q \Rightarrow L, \sigma_2 \quad \sigma_2, \gamma[x \mapsto L] \models^* u \Rightarrow \sigma_3, \omega}{\sigma_1, \gamma \models \text{for } x \text{ in } q \text{ return } u \Rightarrow \sigma_3, \omega}$	(SU-FOR)
$\frac{\sigma_1, \gamma \models q \Rightarrow L, \sigma_2 \quad \sigma_2, \gamma[x \mapsto L] \models u \Rightarrow \sigma_3, \omega}{\sigma_1, \gamma \models \text{let } x := q \text{ return } u \Rightarrow \sigma_3, \omega}$	(SU-LET)
$\frac{\sigma_1, \gamma \models q \Rightarrow \sigma_2, l \cdot L \quad \sigma_2, \gamma \models u_1 \Rightarrow \sigma_3, \omega_1}{\sigma_1, \gamma \models \text{if } (q) \text{ then } u_1 \text{ else } u_2 \Rightarrow \sigma_3, \omega_1}$	(SU-IF1)
$\frac{\sigma_1, \gamma \models q \Rightarrow \sigma_2, () \quad \sigma_2, \gamma \models u_2 \Rightarrow \sigma_3, \omega_2}{\sigma_1, \gamma \models \text{if } (q) \text{ then } u_1 \text{ else } u_2 \Rightarrow \sigma_3, \omega_2}$	(SU-IF2)
$\frac{\sigma_1, \gamma \models q \xRightarrow{\text{copy}} \sigma_2, L \quad \sigma_2, \gamma \models q_0 \Rightarrow \sigma_3, l_0}{\sigma_1, \gamma \models \text{insert } q \text{ pos } q_0 \Rightarrow \sigma_3, \text{ins}(L, \text{pos}, l_0)}$	(SU-INSERT)
$\frac{\sigma_1, \gamma \models q_0 \Rightarrow \sigma_2, l_0}{\sigma_1, \gamma \models \text{delete } q_0 \Rightarrow \sigma_2, \text{del}(l_0)}$	(SU-DELETE)
$\frac{\sigma_1, \gamma \models q_0 \Rightarrow \sigma_2, l_0 \quad \sigma_2, \gamma \models q \xRightarrow{\text{copy}} \sigma_3, L}{\sigma_1, \gamma \models \text{replace } q_0 \text{ with } q \Rightarrow \sigma_3, \text{repl}(l_0, L)}$	(SU-REPLACE)
$\frac{\sigma_1, \gamma \models q_0 \Rightarrow \sigma_2, l}{\sigma_1, \gamma \models \text{rename } q_0 \text{ as } b \Rightarrow \sigma_2, \text{ren}(l_0, b)}$	(SU-RENAME)
$\frac{}{\sigma, \gamma[x \mapsto ()] \models^* u \Rightarrow \sigma, \epsilon}$	(SU-ITERBASE)
$\frac{\sigma_1, \gamma[x \mapsto l] \models u \Rightarrow \sigma_2, \omega_1 \quad \sigma_2, \gamma[x \mapsto L] \models^* u \Rightarrow \sigma_3, \omega_2}{\sigma_1, \gamma[x \mapsto l \cdot L] \models^* u \Rightarrow \sigma_3, \omega_1; \omega_2}$	(SU-ITER)

Table 2.4: Rules for evaluating update expressions to pending update lists

$$\begin{array}{c}
\frac{\sigma(l') = a[L_1 \cdot l \cdot L_2]}{\sigma \models \text{del}(l) \rightsquigarrow \sigma[l' := a[L_1 \cdot L_2]]} \text{ (SAU-DELETE)} \\
\\
\frac{\sigma(l) = a[L]}{\sigma \models \text{ren}(l, b) \rightsquigarrow \sigma[l := b[L]]} \text{ (SAU-RENAME)} \\
\\
\frac{\sigma(l') = a[L_1 \cdot l \cdot L_2]}{\sigma \models \text{ins}(L, \text{before}, l) \rightsquigarrow \sigma[l' := a[L_1 \cdot L \cdot l \cdot L_2]]} \text{ (SAU-INSERTBEFORE)} \\
\\
\frac{\sigma(l') = a[L_1 \cdot l \cdot L_2]}{\sigma \models \text{ins}(L, \text{after}, l) \rightsquigarrow \sigma[l' := a[L_1 \cdot l \cdot L \cdot L_2]]} \text{ (SAU-INSERTAFTER)} \\
\\
\frac{\sigma(l) = a[L_1 \cdot L_2]}{\sigma \models \text{ins}(L, \text{into}, l) \rightsquigarrow \sigma[l := a[L_1 \cdot L \cdot L_2]]} \text{ (SAU-INSERTINTO)} \\
\\
\frac{\sigma(l) = a[L_1]}{\sigma \models \text{ins}(L, \text{into as first}, l) \rightsquigarrow \sigma[l := a[L \cdot L_1]]} \text{ (SAU-INSERTFIRST)} \\
\\
\frac{\sigma(l) = a[L_1]}{\sigma \models \text{ins}(L, \text{into as last}, l) \rightsquigarrow \sigma[l := a[L_1 \cdot L]]} \text{ (SAU-INSERTLAST)} \\
\\
\frac{\sigma(l') = a[L_1 \cdot l \cdot L_2]}{\sigma \models \text{repl}(l, L) \rightsquigarrow \sigma[l' := a[L_1 \cdot L \cdot L_2]]} \text{ (SAU-REPLACE)} \\
\\
\frac{}{\sigma \models \epsilon \rightsquigarrow \sigma} \text{ (SAU-EMPTY)} \quad \frac{\sigma \models \omega_j \rightsquigarrow \sigma' \quad \sigma' \models \omega_k \rightsquigarrow \sigma'' \quad \{j, k\} = \{1, 2\}}{\sigma \models \omega_1, \omega_2 \rightsquigarrow \sigma''} \text{ (SAU-ITER)}
\end{array}$$

Table 2.5: Update pending list application

**Sanity Check and Update Reordering** In this phase, the sanity checks verify that UPL is consistent. As in [BC09a], we do not model this phase here, and we defer the work to an abstract predicate *sanitycheck*( $\omega$ ). The check enforces for instance that the same node is not a target of multiple rename or replace operators, that text nodes are not the targets of inserts, and that the elements being inserted are rooted at element nodes, as specified in [RCD<sup>+</sup>11].

**Application of the UPL** Applying the UPL  $\omega$  to the input store  $\sigma_\omega$  produces the updated store  $\sigma_u$ . This is denoted by

$$\sigma_\omega \models \omega \rightsquigarrow \sigma_u$$

Of course, all locations used in  $\omega$  are contained in  $\sigma_\omega$ .

The composition of all phases of update semantics is denoted by

$$\sigma, \gamma \models u : \sigma_u$$

and it is defined as follows

$$\frac{\sigma, \gamma \models u \Rightarrow \sigma_\omega, \omega \quad \text{sanitycheck}(\omega) \quad \sigma_\omega \models \omega \rightsquigarrow \sigma_u}{\sigma, \gamma \models u : \sigma_u}$$

As for queries, we explain the semantics of updates through some examples. Let us consider the execution of the update below that deletes all textual nodes in the tree  $t_1$  of Figure 2.1, wrt the environment  $\gamma = \{ \mathbf{x} \mapsto l_{t_1} \}$ .

for y in x/descendant :: text() return delete y

The end-to-end semantics of the update consists of generating the UPL of nodes to delete, and then doing the sanity check on the UPL and of applying the update pending list so as to obtain the updated store  $\sigma'$ .

$$\frac{\frac{\frac{\textcircled{1}}{\sigma, \gamma \models \text{for y in x/descendant :: text() return delete y} \Rightarrow \sigma, (\text{del}(l'_2), \text{del}(l''_3))} \quad \text{sanitycheck}(\text{del}(l'_2), \text{del}(l''_3))}{\frac{\textcircled{2}}{\sigma \models \text{del}(l'_2), \text{del}(l''_3) \rightsquigarrow \sigma'}}}{\sigma, \gamma \models \text{for y in x/descendant :: text() return delete y} : \sigma'}$$

- ① The derivation of the branch generating the UPL with locations to be deleted is the following.

$$\begin{array}{c}
\sigma, \gamma \models \mathbf{x}/\text{descendant} :: \text{text}() \Rightarrow \sigma, (l'_2, l''_3) \\
\frac{\sigma, \gamma[y \mapsto l'_2] \models \text{delete } y \Rightarrow \sigma, (\text{del}(l'_2)) \quad \sigma, \gamma[y \mapsto l''_3] \models \text{delete } y \Rightarrow \sigma, (\text{del}(l''_3))}{\sigma, \gamma[y \mapsto (l'_2, l''_3)] \models^* \text{delete } y \Rightarrow \sigma, (\text{del}(l'_2), \text{del}(l''_3))} \\
\hline
\sigma, \gamma \models \text{for } y \text{ in } \mathbf{x}/\text{descendant} :: \text{text}() \text{ return delete } y \Rightarrow \sigma, (\text{del}(l'_2), \text{del}(l''_3))
\end{array}$$

Here the UPL is computed by first evaluating the navigational expression selecting all textual nodes, and then by iterating over the resulting sequence of locations.

- ② The derivation of the branch applying the UPL is the following.

$$\begin{array}{c}
\frac{\sigma(l'_3) = \text{name}[l''_3] \quad \sigma'' = \sigma[l'_3 := \text{name}[]]}{\sigma \models \text{del}(l'_3) \rightsquigarrow \sigma''} \quad \frac{\sigma''(l_2) = \text{title}[l'_2] \quad \sigma' = \sigma[l_2 := \text{title}[]]}{\sigma'' \models \text{del}(l'_2) \rightsquigarrow \sigma'} \\
\hline
\sigma \models \text{del}(l'_2), \text{del}(l'_3) \rightsquigarrow \sigma'
\end{array}$$

Here the application of the UPL consists of removing the use of the target locations by means of rule (SAU-ITER) and (SAU-DELETE). As mentioned before, in this semantics we do not explicitly take into account order of update application, in fact rule (SAU-ITER) specifies that two update pending *lists* can be executed in any order. In order to stress this, in the example, we chose to remove first location  $l'_3$  and then  $l_2$ , while the input UPL presents them in inverse order. For instance for all valid sequences of updates targeting the same locations like

```

for x in y//a
return (delete x, rename x as b)

```

the semantics allows to apply non-deterministically the atomic delete and rename commands in any order. Whatever order of application is chosen, the result of the update is to delete all nodes labeled as  $a$ . Notice that once atomic updates are collected, we may have that a location is both target of an insert and of a rename. In fact, the sanity check just enforces that no location is target of multiple rename or replace.

Another interesting point to stress is that for deleted and replace updates, no node is destroyed in the store. In fact, if  $\sigma, \gamma \models \mathbf{u} \Rightarrow \sigma_\omega, \omega$  and  $\sigma_\omega \models \omega \rightsquigarrow \sigma_u$  then  $\text{dom}(\sigma) \subseteq \text{dom}(\sigma_u)$ . For a tree  $t = (\sigma, l_t)$ ,  $\mathbf{u}(t)$  denotes the tree  $(\sigma_u @ l_t, l_t)$ . Note that  $\text{dom}(\sigma) \subseteq \text{dom}(\sigma_u @ l_t)$  may not hold anymore, since  $\sigma_u @ l_t$  does not keep locations disconnected to the root  $l_t$  after the update.

To see how the input store is extended by updates such as insert and replace, we consider the following expression on the tree  $t_1$  in Figure 2.1, wrt the environment  $\gamma = \{x \mapsto l_{t_1}\}$ . The update inserts a new tree listing all books as a child of the root.

$$\begin{array}{c}
 \textcircled{1} \\
 \hline
 \sigma, \gamma \models \text{insert } \langle \text{list} \rangle x / \text{descendant} :: \text{book} \langle / \text{list} \rangle \text{ into } x / \text{self} :: \text{bib} \Rightarrow \text{ins}(L, \text{into}, l_{t_1}) \\
 \text{sanitycheck}(\text{ins}(L, \text{into}, l_{t_1})) \\
 \textcircled{2} \\
 \hline
 \sigma \models \text{ins}(L, \text{into}, l_0) \rightsquigarrow \sigma' \\
 \hline
 \sigma, \gamma \models \text{insert } \langle \text{list} \rangle x / \text{descendant} :: \text{book} \langle / \text{list} \rangle \text{ into } x / \text{self} :: \text{bib} : \sigma'
 \end{array}$$

We explore a possible derivation tree of the update semantics, by means of an example.

- ① The derivation for the branch generating the UPL is the following.

$$\begin{array}{c}
 \sigma, \gamma \models x / \text{self} :: \text{bib} \Rightarrow \sigma, (l_{t_1}) \\
 \\
 \sigma, \gamma \models x / \text{descendant} :: \text{book} \Rightarrow \sigma_1, l_1 \\
 (l_c, \sigma_c) = \text{copy}(l_1, \sigma_1) \quad \sigma' = \sigma_1 \cup \sigma_c \\
 \hline
 \sigma, \gamma \models \langle \text{list} \rangle x / \text{descendant} :: \text{book} \langle / \text{list} \rangle \xrightarrow{\text{copy}} \sigma', l_c \\
 \hline
 \sigma, \gamma \models \text{insert } \langle \text{list} \rangle x / \text{descendant} :: \text{book} \langle / \text{list} \rangle \text{ into } x / \text{self} :: \text{bib} \Rightarrow \text{ins}(l_c, \text{into}, l_{t_1})
 \end{array}$$

Here the UPL is composed by a single atomic insert update command. The arguments of the atomic update are taken from the outputs of the *target* and the *source* query. From one hand, the target query simply returns the root of the document. From the other hand, the source query returns the copy of a fresh tree listing all books.

- ② The derivation concerning the application of the UPL is the following.

$$\begin{array}{c}
 \sigma(l_{t_1}) = \text{bib}[l_1] \quad \sigma' = \sigma[l_{t_1} := \text{bib}[l_1, l_c]] \\
 \hline
 \sigma \models \text{ins}(l_c, \text{into}, l_{t_1}) \rightsquigarrow \sigma'
 \end{array}$$

Here the root location receives one more child. In this example we chose to insert the new location in last position. It is worth noticing that when position *into* is specified, our semantics, as well as the W3C one, allows to insert the node non-deterministically, in any position.

**Assumptions on the Query Language** In order to simplify the formal treatment, we do some assumptions on the query language.

(Node Element Construction) We assume that node element construction  $\langle a \rangle q \langle /a \rangle$  is not used in the left-hand side expression of a for/let-expression. This restriction is met by a very large class of queries used in practice, while queries like

$$\text{let } x := \langle a \rangle q \langle /a \rangle \text{ return } \langle b \rangle x \langle /b \rangle$$

can be rewritten by simple variable substitution into

$$\langle b \rangle \langle a \rangle q \langle /a \rangle \langle /b \rangle$$

(Text Element Construction) Given that a constant query is independent of any update, we assume that “**txt**” is either nested in a for expression.

**Updates Preserving the Schema** We assume also that the updates are preserving the schema, in the sense that they do not introduce new chains. Delete updates that break schema constraints are however allowed. As we show later in Chapter 3, schema evolution is an issue that goes beyond the scope of detecting independence.

## 2.5 Query-Update Independence

We now formalize query-update independence. To this end we need a notion of equivalence over XML stores.

**Definition 2.5.1** (Equivalence). *Given two stores  $\sigma$  and  $\sigma'$ , two locations  $l \in \sigma$  and  $l' \in \sigma'$  are said to be value equivalent, written  $(\sigma, l) \cong (\sigma', l')$ , iff the two trees  $\sigma @ l$  and  $\sigma' @ l'$  are isomorphic (they possibly differ only in terms of locations). We write  $(\sigma, L) \cong (\sigma', L')$  to denote the value equivalence on location sequences  $L = (l_1, \dots, l_n)$  and  $L' = (l'_1, \dots, l'_n)$ , with  $l_i \in \sigma$  and  $l'_i \in \sigma'$ , and holding iff  $(\sigma, l_i) \cong (\sigma', l'_i)$  for  $i = 1..n$ .*

**Definition 2.5.2** (Independence:  $q \perp u$ ,  $q \perp_d u$ ). *Let  $\sigma$  be a store and  $\gamma$  a variable environment over  $\sigma$ .*

*A query  $q$  and an update  $u$  are said to be independent wrt  $(\sigma, \gamma)$  if*

$$\sigma, \gamma \models q \Rightarrow \sigma_q, L_q \quad \sigma, \gamma \models u : \sigma_u \quad \sigma_u, \gamma \models q \Rightarrow \sigma'_q, L'_q$$

*implies  $(\sigma_q, L_q) \cong (\sigma'_q, L'_q)$ .*

*Also,  $q$  and  $u$  are independent, written  $q \perp u$ , iff they are independent for any pair  $(\sigma, \gamma)$ .*

*Finally,  $q$  and  $u$  are independent wrt the DTD  $d$ , written  $q \perp_d u$ , iff for every tree  $t = (\sigma, l_t) \in d$  and  $\gamma$ , they are independent for  $(\sigma, \gamma)$ .*

As a natural consequence of the fact that XML data are typed by a schema, we assume that our independence analysis is run in a context where all data remain consistent wrt the schema after each update. In case an update triggers schema evolution, then a larger task of schema maintenance has to be carried on. This task may imply existing views (queries) to be reformulated in order to be correct wrt the new schema, and thus it is likely to exclude any other kind of schema-based analysis until its completion.

## Conclusions

In this chapter we presented the formal framework on top of which we build our work. We presented XML databases, schemas, query and update languages, and defined the problem we study: XML query update independence. In the next chapter we will revise the state of the art for detecting independence via static analysis.

## Chapter 3

# Query Update Independence: State of the Art

In this chapter we present the state of the art techniques for checking XML query-update independence. In the first part, we report on the theoretical background that is relevant to detecting independence via static analysis, and outline the intractability of the problem for practical cases. In the second part, we present a family of approximate algorithms developed for checking independence that, being sound but not complete, try to ensure a tractable analysis.

### 3.1 Decidability and Complexity of Exact Static Independence Analysis

Former studies by Benedikt and Cheney [BC09a] and by Raghavachari and Shmueli [RS06] provided the theoretical background of checking independence via static analysis. These works show that independence essentially inherits its computational difficulty from more fundamental problems such as query *equivalence* and *containment*. Because of this, the problem is challenging.

The following results give a flavor on how checking independence ranges from undecidability to tractability. Below, boolean atomic queries and updates are expressions without aggregates and tests on data values returning either the empty sequence or a constant value [BK06].

**Theorem 3.1.1** ([BC09a]). *Independence is undecidable for XQuery queries and updates.*

**Theorem 3.1.2** ([BC09a]). *For boolean atomic XQuery queries and updates, independence is decidable but non-elementary.*

**Theorem 3.1.3** ([RS06]). *For queries and updates defined by using XPath expressions featuring descendant navigation, \* wildcard operator and predicates, independence is NP-Complete.*



**Theorem 3.1.4** ([RS06]). *Independence is polynomial for queries and updates defined by downward XPath expressions using at most two of the following features: descendant axis, \* wildcard and predicates.*

In the following, we discuss the main steps of the proofs of the above theorems, with the aim of giving more details on the links between independence, equivalence and containment.

The main observation for relating independence with equivalence suggested in [BC09a] is that we can solve independence if we can solve the equivalence problem for *compositions* of queries and updates. Equivalence means that two programs always agree on their output, once evaluated on the same input database. That is, a query  $q$  is independent from an update  $u$  if the equivalence  $q \equiv q \circ u$  holds, provided that the composition  $q \circ u$  and the notion of equivalence  $\equiv$  are properly defined. Results on the decidability of independence follow from this idea, as we now illustrate.

**General Undecidability of Independence** Benedikt and Cheney showed that independence is undecidable with a reduction to the satisfiability problem for queries expressed as first order logic formulas over data trees (corresponding to XQuery expressions) [BC09a]. It is well known that satisfiability of first order logic is undecidable, and this is also true in the signature of trees. With signature of trees, we mean the binary and unary relations that allow to define trees (see [BK09] for more details). The idea is that we could solve independence if we had solved satisfiability of first order logic over data trees. Fix an XQuery query  $q$ , and define a query update pair in the following way

$$q' = \text{if } (/root/a) \text{ then } q \text{ else } () \quad u = \text{insert } \langle a \rangle \text{ into } /root$$

Then,

$$q \text{ is unsatisfiable for all } \sigma, \gamma \quad \text{if and only if} \quad q' \text{ is independent of } u \text{ for all } \sigma, \gamma$$

In fact, query  $q'$  outputs constantly the empty sequence (before and after  $u$ ) on any  $\sigma, \gamma$ , if and only if  $q$  is not satisfiable on any  $\sigma, \gamma$ . Since the satisfiability problem for XQuery queries is undecidable, we conclude that independence for XQuery queries and updates is generally undecidable. This is stated by Theorem 3.1.1.

**Decidability of Boolean Independence for Core Atomic XQuery** Benedikt and Cheney showed also that independence is decidable for a relevant subset of XML queries and updates [BC09a], as stated by Theorem 3.1.2. Again, the result is obtained from decidability of satisfiability of first order logic over trees. This result is interesting because it sheds light on the links between independence and query equivalence.

Each step of the proof is quite dense, and will be detailed next.

1. Define independence using a weaker notion of equivalence.

2. Restrict the query and update languages.
3. Show that, with such restrictions, queries and updates correspond to *first order interpretations*.
4. Show that first order interpretations can be composed and tested to be equivalent, thus conclude decidability of independence.

1) Define independence using a weaker notion of equivalence.

Several notions of equivalence can be used to define independence (i.e.,  $q \equiv q \circ u$ ), notably, boolean equivalence, node equivalence or subtree equivalence. They naturally differ on the test they do. *Boolean* equivalence checks if the results of two expressions always agrees on non-emptiness. *Node* equivalence checks if node ids in result sequences are isomorphic. *Subtree* equivalence checks if forests of the two expressions are isomorphic. Decidability of boolean query equivalence has been showed in [BK06] for the class of atomic queries. Differently, decidability of node-equivalence and subtree-equivalence are still open problems. Consequently, in [BC09a] the authors studied the analogous problem of *boolean independence*. Boolean independence just checks if the result of a query changes from empty to non empty or vice-versa, after update. Now, for testing boolean independence it is sufficient to translate a query  $q$  in a boolean-query of the form `if (q) then 1 else 0`.

2) Restrict the query and update language.

Benedikt and Cheney [BC09a] considered XML queries and updates without aggregates and tests on data values, which is originally referred as Core Atomic XQuery in [BK06]. This is roughly equivalent to the language defined in Section 2.3. Because we want to check boolean independence, we will consider Boolean Core Atomic XQuery queries of the form `if (q) then 1 else 0`, and Core Atomic XQuery updates built on such expressions

3) Show that queries and updates correspond to first order interpretations.

Benedikt and Koch [BK06] showed that Boolean Core Atomic XQuery queries of XQuery are given by first-order queries in the signature of trees. First order interpretations are formulas that define the output structure in function of a given input. First order interpretations model queries. The crucial point is the following. It can be shown (see [BC09a]) that the addition of updates to Core Atomic XQuery still yields first order interpretations. This means that, while in general it is true that updates are more expressive than queries, this is not the case for Core Atomic XQuery. For this fragment queries and updates have the same expressive power, and they precisely correspond to first order interpretations.

4) Conclude decidability of independence.

Since first-order interpretations are closed under composition (Theorem 5.1 of [BK09]), the composition of a query and an update is again a first order interpretation. So the

goal now is to test equivalence of first order interpretations. In the special case of a boolean query, a first-order interpretation is just a first-order sentence (a formula with no variables). First-order sentences can be translated into tree automata. Equivalence of automata obtained by first order sentences can be decided. Therefore, equivalence of first order interpretations for boolean queries can be decided. Hence, equivalence of a query against the composition of a query and an update can be decided. Hence, independence can be effectively decided. It is worth noticing that it is critical for decidability that data value comparisons are disallowed.

Benedikt and Cheney also give the complexity of the problem when the tag alphabet is fixed. This captures the case where a schema is available. In this case, because the satisfiability problem for first order logic over labeled trees (with fixed alphabet) is non-elementary [Vor96] we obtain that the independence problem in the presence of the schema is non-elementary. As outlined in [Ben10], decidability of equivalence, and hence independence, when considering powerful notions of equivalence such as node and subtree equivalence is still an open question.

**NP-Completeness for XPath( $//, *, []$ ) fragment** Another contribution to the study of XML independence comes from the work of Raghavachari and Shmueli [RS06]. This work investigates the complexity of detecting *conflicting* query and update expressions. A query and an update conflict when a node accessed by the query is also impacted by the update. This is a dual notion of independence, that indeed holds in the absence of conflicts. Detecting conflicting expressions has been proved NP-complete for XML queries and updates expressed in the downward fragment of XPath featuring  $*$  wildcard and branching, and a notion of delete/insert node-conflict reminiscent of node-equivalence (are the ids of resulting sequences isomorphic?). NP-completeness follows from a reduction to the complement of the containment problem for XPath. Actually, it follows from the dual non-containment problem. The non-containment problem for XPath asks whether two expressions disagree on the emptiness of the result, for some trees. In brief, it turns out that detecting this property is enough to detect if the result of a query may change after update. For instance, by taking an instance for which query result is empty and adding some data so as to get a non-empty query result. This operation is possible if the expressions are not in the containment relation. Updates are defined in terms of queries, hence solving non-containment between query-paths and update-paths would solve the problem of determining if the expressions conflict. Miklau and Suciu previously showed that non-containment is NP-hard [MS04], that in turn determines the hardness of independence.

**Polynomial fragments of XPath( $//, *, []$ )** Along the line of the previous result, Ravagachari and Shmueli also showed that by further restricting the query and update language polynomial time algorithms for independence can be given [RS06]. This is stated by Theorem 3.1.3. This holds for queries and updates written with path expressions belonging to a subset of XPath( $//, *, []$ ) which is employing at most two features among

descendant navigation “//”, \* wildcard and branching “[ ]”. The authors in [RS06] showed the relationship with query containment of XPath expressions for such languages [MS04, AYCLS01]. This shows that tractable algorithms for independence can be given only by substantially restricting the query and update language. These subsets of XPath are quite restrained and may not support realistic XQuery workloads.

The results we discussed show that independence is not tractable in practical cases. Its general undecidability and hardness suggest that in order to do static analysis one is obliged to proceed with *approximate* algorithms, trading completeness with tractability, and then devise analysis that can be useful in practice. This is what motivated recent research efforts towards approximate techniques discussed next.

## 3.2 Approximate Algorithms

The hardness of independence leads to attack the problem with approximate solutions. Approximate algorithms are sound but not complete methods, that is, they either correctly detect an independence or conservatively raise a conflict. As we will see throughout this dissertation, approximate solutions can be extremely effective by exploiting structural properties of queries and updates.

The state-of-the-art techniques for detecting independence reformulate the problem in terms of an intersection problem. Precisely, testing independence is made by checking whether the data needed by a query intersects with the data affected by an update. The techniques now decompose the problem in two parts: 1) providing an abstract representation of the data that is accessed by the query and impacted by the update, and 2) defining the intersection test. Before illustrating the state-of-the-art, we discuss the problem of providing an abstract representation of the data and show a link with XML projection.

The problem of computing a subset of the input database that is accessed by a query, and that is sufficient for its execution, is known as XML data-projection [MS03, BCCN06]. This problem has been studied also for updates in [BBC<sup>+</sup>11]. The set of nodes involved by query/update execution is usually captured via static analysis. XML projection has developed the ideas of using either *XPath navigational paths* or *XML types* (when a schema is available) in order to abstract over the collection of data. These ideas gave birth to a series of works showing how projection can be used to optimize query execution [MS03, BCCN06] or for locally updating a database [BBC<sup>+</sup>11]. Indeed, all approaches in the current literature dealing with independence [GRS08, BC09a, BC10, BTCU12] are likely to follow these ideas. However, substantial technical differences arise between methods for projection and independence, as we will see next.

In the following, we present the state-of-the-art techniques for independence in detail. We begin by presenting the analysis based on XML types [BC09a] which are the closest to our work since they use schema information. Then we present the analysis based on XPath paths [GRS08] and then the approach with stabilizers [BC10].

### 3.2.1 Schema-based independence analysis for XML updates [BC09a]

This work is the first to propose a static analysis for detecting independence, in the presence of a schema. The idea at the basis of the schema-based analysis of [BC09a] is to infer *type-names* of nodes that are accessed by a query and impacted by an update, and then check that these do not overlap. The main contribution of [BC09a] is a type system that features very low complexity and which is able to detect a consistent number of independencies statically.

The technique targets a fragment of XQuery that is essentially the one presented in Section 2.3. Mild differences comes from the fact that queries can navigate on constructed data, and updates do not need to preserve the schema.

The independence analysis relies on three inference systems.

- (**RES**) An inference system for typing the *result* (**RES**) of a query expression. The judgment  $d, \Gamma \vdash_{(\mathbf{RES})} q : \tau$  stands for a relation between the schema  $d$ , the static variable environment  $\Gamma$ , the query  $q$  and the set of types  $\tau$  where  $\tau$  contains at least all types of nodes, defined in  $d$ , that can be in the result sequences of  $q$ .
- (**SAC**) An inference system for typing the data *accessed* by a query, that is called static access cover (**SAC**). The judgment  $d, \Gamma \vdash_{(\mathbf{SAC})} q : \tau$  is a relation between the schema  $d$ , the static variable environment  $\Gamma$ , the query  $q$  and the set of types  $\tau$  where  $\tau$  contains at least all types of nodes, defined in  $d$ , that are needed to evaluate  $q$ .
- (**IMP**) An inference system for typing the data *impacted* by the update (**IMP**). The judgment  $d, \Gamma \vdash_{(\mathbf{IMP})} q : \tau$  is a relation between the schema  $d$ , the static variable environment  $\Gamma$ , the query  $u$  and the set of types  $\tau$  where  $\tau$  contains at least all types of nodes, defined in  $d$ , that are updated by  $u$ .

We now illustrate through some examples how these three systems work, and how type inference is done.

Consider the bibliographic DTD  $d_3$  and the following query and update pair, for which we want to verify independence.

$q = \text{for } x \text{ in doc/descendant::book return } x/\text{author} \quad u = \text{delete doc//title}$

We always assume that the context of the first navigational step of a query or update is the schema root type, thus in all of the examples of this chapter we assume that  $\Gamma$  is a static environment such that  $\Gamma(\text{doc}) = \{\text{bib}\}$ . The derivation tree for the inference of types accessed by  $q$  is the following.

```

      bib  ←  book*
      book ←  title, author+, publisher+
author, publisher ←  name
      title, name ←  String

```

Figure 3.1: DTD  $d_3$ 

$$\begin{array}{c}
\textcircled{1} \qquad \textcircled{2} \qquad \textcircled{3} \\
\hline
d_3, \Gamma \vdash_{(\mathbf{SAC})} \text{doc/descendant} :: \text{book} : \tau_1 \quad d_3, \Gamma \vdash_{(\mathbf{RES})} \text{doc/descendant} :: \text{book} : \tau_2 \quad d_3, \Gamma[x \mapsto \tau_2] \vdash_{(\mathbf{SAC})} x/\text{author} : \tau_3 \\
\hline
d_3, \Gamma \vdash_{(\mathbf{SAC})} \text{for } x \text{ in doc/descendant} :: \text{book} \text{ return } x/\text{author} : \tau_1 \cup \tau_3
\end{array}$$

The set of types accessed by  $q$  is the union of types accessed by its subexpressions, i.e.  $\tau_1 \cup \tau_3$ . Types are inferred by calling system **SAC** as shown by derivations  $\textcircled{1}$  and  $\textcircled{3}$ . What is inferred for the query in  $\tau_1$  and  $\tau_3$ ? Because **doc/descendant** :: *book* selects all book nodes in a store, in  $\tau_1$  we expect to find *book* and all of its ancestor types defined in  $d_1$ . Because  $x/\text{author}$  selects all author children of nodes bound to  $x$ , and  $x$  is only bound to book nodes, in  $\tau_3$  we expect to find types *book* and *author*. Derivation  $\textcircled{2}$  shows also the role of the inference system **RES** in the process. It is used in order to provide types for variable  $x$  in the subsequent **SAC** derivation.

By exploring the next level of derivation we can understand more peculiarities of type inference.

$\textcircled{1}$  The judgment  $d_3, \Gamma \vdash_{(\mathbf{SAC})} \text{doc/descendant} :: \text{book} : \tau_1$  is developed as follows

$$\frac{\Gamma(\text{doc}) = \{ \text{bib} \} \quad d_3, \Gamma \vdash_{(\mathbf{RES})} \text{doc/descendant} :: * : \Sigma_S^{d_3} \setminus \{ \text{bib} \}}{d_3, \Gamma \vdash_{(\mathbf{SAC})} \text{doc/descendant} :: \text{book} : \Sigma_S^{d_3}}$$

The judgment infers the static access cover for **doc/descendant** :: *book*. The inference result is the whole set of types of  $d_3$ . This is obtained as a union of the root type (that is the context type for the step **descendant** ::  $*$ ), together with all descendant types of the root (that type the result of the step **descendant** ::  $*$ ). This last one is obtained by calling the **RES** system.

The first peculiarity of the system **SAC** we can observe is how the static access cover is computed. As said before, a sound static access cover would contain types for *book* nodes and their ancestors. To accomplish this, the navigational step **descendant** :: *book* is approximated by replacing the *book* label with the  $*$  wildcard (denoting any node). This of course makes the analysis sound but imprecise,

because all schema types are inferred, and thus also all types unrelated with the step. As a consequence of this, the query is likely to be deemed dependent wrt almost all updates. We stress that this is always the case, for all queries starting with a descendant navigation.

- ② The judgment  $d_3, \Gamma \vdash_{(\mathbf{RES})} \text{doc}/\text{descendant} :: \text{book} : \tau_2$  is developed as follows.

$$\frac{\Gamma(\text{doc}) = \{\text{bib}\} \quad \tau_2 = \{ \text{book} \mid \text{book is a descendant type of bib in } d_3 \}}{d_3, \Gamma \vdash_{(\mathbf{RES})} \text{doc}/\text{descendant} :: \text{book} : \{ \text{book} \}}$$

It infers types for the result of  $\text{doc}/\text{descendant} :: \text{book}$ , this time for providing input to the next inference step. The derivation stresses that schema information is used in order to check if type **book** is a descendant of type **bib**. If this is not the case then  $\tau_2 = \emptyset$ .

- ③ The judgment  $d_3, \Gamma[x \mapsto \tau_2] \vdash_{(\mathbf{SAC})} x/\text{author} : \tau_3$  is developed in the following way

$$\frac{d_3, \Gamma[x \mapsto \{ \text{book} \}] \vdash_{(\mathbf{RES})} x/* : \{ \text{title}, \text{author}, \text{publisher} \}}{d_3, \Gamma[x \mapsto \{ \text{book} \}] \vdash_{(\mathbf{SAC})} x/\text{author} : \{ \text{book}, \text{title}, \text{author}, \text{publisher} \}}$$

This infers the static access cover for  $x/\text{author}$ . The result is constituted by the union of types bound to variable  $x$  with the types of nodes in the result of step  $x/*$ .

Once again we observe that the child step is approximated, by replacing the label *author* with  $*$ . This means that all children types of **book** are inferred namely **title**, **author** and **publisher**. Once again, the typing is sound but is not precise since while type **author** is needed by step  $x/\text{author}$ , types **title** and **publisher** are not. Analogously to the previous case, the query is likely to be deemed as dependent wrt all updates impacting a node related to **book**.

Once types are inferred for the query, the update is analyzed.

The types of a node impacted by  $u$  are inferred by the system **IMP** as follows.

$$\frac{d_3, \Gamma \vdash_{(\mathbf{RES})} \text{doc}/\text{title} : \{ \text{title} \} \quad d_3, \Gamma[x \mapsto \{ \text{title} \}] \vdash_{(\mathbf{RES})} x/\text{parent} :: * : \{ \text{book} \}}{d_3, \Gamma \vdash_{(\mathbf{IMP})} \text{delete doc}/\text{title} : \{ \text{book} \}}$$

The set of types deleted by the update is computed in two steps. First, system **RES** infer types for nodes effectively deleted, that are those returned by the target expression  $\text{doc}/\text{price}$ . Second, a parent navigation for such types is performed. This is the second peculiarity of the system, and it is meant to state that a node whose children are modified by the update is considered as impacted by the update as well. For the DTD  $d_3$ , this means that the data type impacted by the update is **book**. As a consequence of this, the

update is deemed dependent with all queries accessing `book` or one of its children. The same approximation is done for `rename`, `insert` and `replace` updates.

After type inference is done, for the analysis to be sound, all descendant types of nodes in the result of the query are also added to the set of types accessed by the query. For instance, because `author` types the root of a subtree output by `q`, all descendant types of `author`, namely `name` and `String`, are included in the result of the analysis.

Once type inference is done for `q` and `u`, providing sets of accessed and updated types  $\tau_q$  and  $\tau_u$ , static independence holds iff

$$\tau_q \cap \tau_u = \emptyset$$

Another peculiarity of this method is that by basing an analysis on type-names, even if type inference is precise, independence cannot be detected in many cases. In fact, when a type is *reused* in several definitions, it may be difficult, or even impossible, to disambiguate it.

To illustrate, consider again the DTD `d3`, the environment  $\Gamma$ , and the following pair of query and update expressions.

$$q = \text{doc} // \text{author} / \text{name} \quad \text{and} \quad u = \text{delete doc} // \text{publisher} / \text{name}$$

Assume that the type `name` is inferred for both expressions. This type captures both nodes returned by the query and nodes impacted by the update, that are labeled with `name`. However, the sets of accessed and updated nodes here are *disjoint*, because the query accessed nodes labeled with `name` nested into nodes labeled with `author` while the update deletes nodes labeled with `name` nested into nodes labeled with `publisher`.<sup>1</sup> By using sets of type names, without recording any information on ancestor types, independence cannot be detected in many cases such as this one.

So far, we have seen a running example of the independence analysis based on types, and the use of inference system **RES**, **SAC** and **IMP**. Now we discuss how this technique supports updates that do not preserve the schema.

As said before, not assuming schema preservation is a nice property of the type-based analysis. However, when flows of updates (possibly violating the initial schema) are evaluated this property alone does not ensure the success of the analysis: one needs to maintain the schema after each update which does not preserve the schema<sup>2</sup>. Indeed, a flow of updates may change completely the structure of the schema thus making the static analysis unsound.

---

<sup>1</sup>Recall that if more than one node labeled with `publisher` is selected by the update path then the update has to be encoded by means of iteration, otherwise a run-time error is raised.

<sup>2</sup>A sound method for determining if an update preserves a schema and for maintaining the schema can be found in [BC09b]



To see this, consider the following update renaming all labels of the schema, provided the DTD  $d_3$  and the environment  $\Gamma$ .

```

u1 =  for x in doc//book
        rename x as "livre"
        rename x/title as "titre"
        rename x/author as "auteur"
        rename x//name as "nom"

```

How does the analysis detect that a query like  $q = \text{doc//livre}$ , whose result is empty over any document valid wrt the DTD, is now drawing some data from this renaming? This case is covered since, while computing the static access cover of the query, system **SAC** approximates  $\text{doc//livre}$  with  $\text{doc//}^*$ . This makes the analysis sound, despite the fact that  $u_1$  does not preserve the schema. Indeed, in the intersection of type-names inferred from  $q$  and  $u_1$  is non-empty.

At this point, if the DTD  $d_3$  is not maintained, the analysis may turn to be unsound. Consider an update  $u_2$  possibly conflicting with  $q$  over the (invalid) document instance obtained after  $u_1$  such as

```

u2 = delete doc//livre

```

Because no definition of *livre* is found in  $d_3$ , the set of types inferred from  $u_2$  over  $d_3$  is empty, and thus independence of  $q$  and  $u_2$  is assumed, while it should not. Therefore, after an update breaking schema constraints the schema needs to be maintained if one aims at a sound analysis.

This short discussion shows that preservation of schema constraints is a problem orthogonal to that of independence. Schema evolution has been studied in [BC09b, CGM11] and it is out of the scope of our work.

We conclude this section by discussing simplicity and complexity of the type analysis of [BC09a] which are the most important features of the system. The approximations made by the type system discussed are motivated by avoiding a complex analysis and providing a simple inference system. The advantages of a simple type system are that it can be easily understood, implemented and proved correct. Proofs of the type systems have been provided and worst case complexity studied. The independence analysis runs in  $O((|d|^2 + |q|)^2 + |u|)$  and it turns out to be very fast in practice.

### 3.2.2 Commutativity analysis for XML updates [GRS08]

Another approach for detecting independence can be drawn from the commutativity analysis of XML updates proposed by Ghelli, Rose and Siméon in [GRS08]. This technique relies on an extension of the path-based analysis proposed in [MS03] that overapproximates the nodes accessed and modified by a given expression, by means of simple XPath expressions for which disjointness is decidable.

This work on commutativity distinguishes itself by two main aspects. First, it studies commutativity of updates, a different problem from that of independence. Second, and more importantly, this study considers XML languages with an iterative semantics, thus different from languages with a snapshot semantics like the one introduced in Section 2.3. As we discuss next, these two differences can be actually levelled out since commutativity allows to detect independence and the path-based approach suits to both kinds of semantics.

We discuss first the links between independence and commutativity. Roughly speaking, independence asks whether the result of a query is preserved after updating, while commutativity asks if the application order of two updates is irrelevant. Deciding commutativity allows to decide independence. Conceptually, this is true because updates are written in terms of queries. This implies that, an instance of the independence problem for  $q, u$  can be reduced to a commutativity one by rewriting  $q$  as a new update  $u' = \text{replace } q \text{ with } q$ , and then check that  $u'$  commutes with  $u$ . This reduction also shows that checking commutativity is at least as hard as checking independence.

The commutativity analysis of [GRS08] has been developed for query languages with side-effects (i.e., updates) such as of XQuery! [GRS06] and XQueryP [CCF<sup>+</sup>06]. These languages have a so called *iterative* semantics, they support immediate update application and explicit left-to-right evaluation order. This is a main difference with the languages considered in [BC09a, BC10, BTCU12], that have a *snapshot* semantics. Iterative languages have a strict left to right evaluation order, while languages with snapshot semantics reorder their updates before application: insert and replace are executed first, then rename and finally delete. Iterative languages allow to see the effects of updates during evaluation.

While iterative languages seem more intuitive, they have a more involved semantics. This is because any expression can read data that it previously inserted. To illustrate, a (possibly) non terminating update such as the following can easily be written in an iterative language.

```
for x in doc//name return insert <name/> after x
```

Assuming that the base data contains at least one node labeled with *name* (that is also a descendant of the node bound to *doc*), the expression keeps on extending the iterating sequence of *name* nodes forever.

Much more machinery effort has to be done in order to prove the soundness of the commutativity analysis for a language with iterative semantics, with respect to languages with snapshot semantics. Many expressions that do not commute in iterative languages do commute in snapshot language (simply because they do not have the same effect), so an analysis for an iterative language is likely to work also for a language with snapshot semantics. Of course, it is possible to define a more general (and less precise) static analysis that works for both kinds of languages.

We proceed now with the description of the path-based analysis, thus illustrating its characteristics.

The path-based analysis relies on a system of *simple path* inference. Simple paths are expressions belonging to the fragment of XPath composed by downward and upward navigation, plus the wildcard \*. Predicates, as well as data values and sibling axis, are not allowed.

The analysis relies on an inference system (**PATH**) for inferring the simple paths of an expression. The judgment  $\Gamma \vdash_{(\mathbf{PATH})} q : (\mathbf{rt}, \mathbf{ac}, \mathbf{up})$  is a relation between the static variable environment  $\Gamma$ , the query  $q$  and the sets of paths  $(\mathbf{rt}, \mathbf{ac}, \mathbf{up})$  where  $\mathbf{rt}$  is the set of path leading to returned subtrees of  $q$ ,  $\mathbf{ac}$  is the set of path representing nodes accessed but possibly not returned by  $q$ , and  $\mathbf{up}$  is the set of path representing nodes that are updated by  $q$ .

Distinguishing among these three kinds of paths allows to capture two main aspects of the language. First, side-effects, by distinguishing between update paths from the rest. Second, query semantics, by distinguishing between accessed and return paths.

Note that here the analysis is supported by a single inference system, and most importantly, the analysis does not depend on an input schema. In the following examples the steps navigating from the root of a tree are written without an explicit variable.

To illustrate the (**PATH**) inference system, we consider again the following query-update pair.

$q = \text{for } x \text{ in } /descendant :: book \text{ return } x/author \quad u = \text{delete } //price$

Inference of return, used and update paths for  $q$  is derived in the following way.

$$\frac{\frac{\textcircled{1}}{\Gamma \vdash_{(\mathbf{PATH})} /descendant :: book : (\mathbf{rt}, \mathbf{ac}, \mathbf{up})} \quad \frac{\textcircled{2}}{\Gamma[x \mapsto \mathbf{rt}] \vdash_{(\mathbf{PATH})} x/author : (\mathbf{rt}', \mathbf{ac}', \mathbf{up}')}}{\Gamma \vdash_{(\mathbf{PATH})} \text{for } x \text{ in } /descendant :: book \text{ return } x/author : (\mathbf{rt}', \mathbf{ac} \cup \mathbf{ac}', \mathbf{up} \cup \mathbf{up}')}$$

Accessed and updated paths inferred for the whole expression are those inferred from each subexpression only. Return paths for the whole expression are naturally those on the right subexpression. Next we will see how the static environment  $\Gamma$  helps in reconstructing navigational paths by following variable bindings.

Notice that the rule for path extraction introduces a bit of imprecision since all accessed and updated paths of the left subexpressions are kept in the final result, independently of the right subexpression. This means that, for instance, when the right branch of the iteration query performs no operation, and thus the query does nothing,

the set of inferred paths for the whole expression may *not* be empty, while it should. This is the case for the following query

`for x in //book return ()`

for which path `//book` is inferred, despite the fact that inferring an empty set of returned and accessed path would still be sound.

Let us explore further derivations and comment on other peculiarities of the analysis.

- ① The judgment  $\Gamma \vdash_{(\text{PATH})} \text{/descendant} :: \text{book} : (\mathbf{rt}, \mathbf{ac}, \mathbf{up})$  is developed as follows

$$\frac{\mathbf{rt} = \{\text{/descendant} :: \text{book}\}}{\Gamma \vdash_{(\text{PATH})} \text{/descendant} :: \text{book} : (\mathbf{rt}, \mathbf{rt}, \emptyset)}$$

In this case path inference is a simple path extraction. Notice that return and accessed paths coincide, and of course no update path is inferred.

- ② The judgment  $\Gamma[x \mapsto \mathbf{rt}] \vdash_{(\text{PATH})} x/\text{author} : (\mathbf{rt}', \mathbf{ac}', \mathbf{up}')$  is developed as follows

$$\frac{\Gamma[x \mapsto \mathbf{rt}] \vdash_{(\text{PATH})} x : (\mathbf{rt}, \emptyset, \emptyset) \quad \mathbf{rt}' = \mathbf{rt} \times \{\text{/author}\} \quad \mathbf{ac}' = \text{prefixes}(\mathbf{rt}')}{\Gamma[x \mapsto \mathbf{rt}] \vdash_{(\text{PATH})} x/\text{author} : (\mathbf{rt}', \mathbf{ac}', \emptyset)}$$

Here the system first reconstructs simple paths by following variable bindings. Indeed, the path `/descendant :: book` that is bound to `x` is combined in a cartesian product with step `/author` in order to infer  $\mathbf{rt}' = \{\text{/descendant} :: \text{book}/\text{author}\}$ . Accessed paths are defined as prefixes of inferred paths. Of course, no update path is inferred since the operation has no side-effect.

The application of the rule may lead to infer exponentially many paths for some rare expressions such as the following one.

`for x1 in (x0/a, x0/b) ... xn in (xn-1/a, xn-1/b) return xn`

this expression would produce exactly  $2^n$  return paths and  $O(n2^n)$  accessed paths. This kind of expressions are however very rare in practice.

Now that paths are inferred from the query, we focus on the update. Path inference for the update expression is derived as follows

$$\frac{\begin{array}{l} \Gamma \vdash_{(\text{PATH})} \text{//price} : (\{\text{//price}\}, \emptyset, \emptyset) \\ \mathbf{ac} = \{\text{//price}\} \quad \mathbf{up} = \{\text{//price}\} \times \{\text{/descendant-or-self} :: *\} \end{array}}{\Gamma \vdash_{(\text{PATH})} \text{delete //price} : (\emptyset, \mathbf{ac}, \mathbf{up})}$$

Here nodes deleted by the expression are captured by combining return paths of *//price* with a **descendant-or-self :: \*** step, denoting that whole subtrees rooted at prices are deleted. Return paths for the target query become also accessed paths for the whole expression. No return path is inferred for the whole update.

Notice that the path analysis precisely capture only the deleted subtrees without including their parent, as is done by the type analysis.

Another characteristics of the analysis is how updates paths are extracted from insert and replace commands. To illustrate, consider the following update.

**insert** *//author into /site/book*

Paths for the update are extracted in the following way.

$$\frac{\begin{array}{l} \Gamma \vdash_{(\text{PATH})} //author : (\mathbf{rt}_1, \mathbf{ac}_1, \mathbf{up}_1) \quad \Gamma \vdash_{(\text{PATH})} /site/book : (\mathbf{rt}_2, \mathbf{ac}_2, \mathbf{up}_2) \\ \mathbf{ac} = \mathbf{rt}_1 \times \{ /descendant\text{-or-self} :: * \} \quad \mathbf{up} = \mathbf{rt}_2 / \times \{ descendant :: * \} \end{array}}{\Gamma \vdash_{(\text{PATH})} \text{insert } //author \text{ into } /site/book : (\emptyset, \mathbf{ac}_1 \cup \mathbf{ac}_2 \cup \mathbf{ac}, \mathbf{up}_1 \cup \mathbf{up}_2 \cup \mathbf{up})}$$

Accessed paths extracted from the two subexpressions become accessed paths for the whole update. Moreover, return paths for the source expression *//author* are combined with a self or descendants axis, and become accessed paths for the whole update, mimicking the fact that an entire subtree is copied.

Update paths extracted from the subexpressions (if any) become update paths for the whole update. Notice that the set of *inserted* paths is inferred by combining path */site/book* with a descendant step, mimicking that the whole subtree rooted at book is updated. This of course is sound, but not complete. For instance, the update is deemed as non-commuting with query */site/book/title*, while it does.

Once return, accessed and update paths are inferred from the expressions, an intersection test is run. The subset of XPath considered for inference includes downward and backward navigational axes, and \* wildcard. Branching is excluded. For this language intersection is decidable. While results in [Hid03] show that the intersection problem for downward XPath expressions is NP-complete, [HKL05] approach shows that intersection is solvable in quadratic time for the downward fragment only.

One of the most important features of the path analysis, is that it works independently without requiring a schema. However, in all the cases where schema is available, schema constraints cannot be easily plugged into the path-based analysis. Soundness of the path-based commutativity analysis has been proved in [GRS08], and also its benefits have been showed when used for optimizing logical plans of queries with side-effects [GORS08].

### 3.2.3 Destabilizers and independence of XML updates [BC10]

Destabilizers [BC10] are a generalization of the idea of statically approximating the set of accessed nodes found in some prior work [BCCN06, MS03, GRS08] and then testing intersection, combined with the concrete use of the data. The destabilizer framework can work in principle for any language and data-model, in [BC10] they are formally implemented for XML queries and updates.

The formal definition of destabilizers is quite technical, thus here we provide just a gentle introduction to the approach, discuss the problematics in computing destabilizers and outline how they are concretely used.

Given a query  $q$ , the destabilizer of  $q$ , denoted by  $\Delta(q)$ , is defined as the set of updates that may change the answer of  $q$ . From this follows that an update  $u$  is independent wrt a query  $q$  when  $u$  does not overlap with  $\Delta(q)$ . In the following we will distinguish between *dynamic* or *static*, and *exact* or *approximate* destabilizers. Dynamic destabilizers are composed by atomic update commands, while static ones are update expressions. Exact destabilizers are *minimal*, in the sense that they produce no false negative (updates that do not change the answer of  $q$ ), while approximated ones may produce false negatives.

An exact dynamic destabilizer is a minimal set of atomic update commands that can affect the result of a query  $q$  evaluated over the tree  $t$ . We recall that an atomic update command is an operation of

$$\text{del}(\_) \quad \text{ins}(\_, \_, \_) \quad \text{ren}(\_, \_) \quad \text{repl}(\_, \_)$$

As an example, consider the document  $t = (\sigma, l_t)$  with two nodes

$$\sigma = \begin{cases} l_t & \leftarrow \text{bib}[l_1] \\ l_1 & \leftarrow \text{book}[] \end{cases}$$

Let  $q$  be the query */bib/book*. The exact dynamic destabilizer of  $q$  wrt  $t$ , denoted by  $\Delta^d(q, t)$ , is composed by all atomic update commands targeting either  $l_t$  or  $l_1$  in  $t$ . The atomic updates that can alter  $l_t$  are the followings

$$\text{del}(l_t) \quad \text{ren}(l_t, "a") \quad \text{repl}(l_t, l')$$

where  $l'$  is a location not isomorphic to  $l_t$ , and  $a$  is a string different from *bib*. The atomic updates that can alter  $l_1$  are the followings.

$$\text{del}(l_1) \quad \text{ins}(L_i, \text{pos}, l_t) \quad \text{ins}(L_i, \text{pos}', l_1) \quad \text{ren}(l_1, "b") \quad \text{repl}(l_1, L_i)$$

where  $L_i$  is a non empty sequence of locations containing a node not labeled with *book*,  $\text{pos}$  is one of  $\{\text{into}, \text{into as first}, \text{into as last}\}$ ,  $\text{pos}'$  is one of  $\{\text{before}, \text{after}\}$ , and  $b$  is a string different from *book*.

An update  $u$  is independent of a query  $q$  if the UPL produced by  $u$  over  $t$  does not overlap with  $\Delta^d(q, t)$ . Consider the update

$$u_1 = \text{insert } \langle \text{new}/\rangle \text{ into } /bib$$

$u_1$  over  $t$  generates the atomic update  $\text{ins}(l_{\text{new}}, \text{into}, l_t)$ . Because  $l_{\text{new}}$  is not labeled with *book*,  $q$  is independent of  $u_1$ . Differently, the update

$$u_2 = \text{delete } /bib/ \star$$

is not independent with  $q$  because it generates the atomic update  $\text{del}(l_1)$  which also belongs to  $\Delta^d(q, t)$ .

Benedikt and Cheney [BC10] illustrate that it is not feasible to deal with dynamic destabilizers as concrete sets of atomic update sequences. Therefore, destabilizers are treated at a symbolic level. The set of nodes involved in atomic updates in  $\Delta^d(q, t)$  is approximated by any update expression targeting those locations. Since updates are defined in terms of queries, it turns out that the destabilizer set is defined as a set of *query* expression targeting all nodes that are arguments of the atomic updates. For  $q = /bib/book$ , the *exact static* destabilizer  $\Delta^s(q, t)$  approximating the dynamic one, is simply

$$\Delta^s(q, t) = \{ /bib, /bib/book \}$$

In developing a static destabilizer, there is a trade-off between the precision (how close it is to the dynamic analog) and the complexity of analyzing the resulting query. The generic static destabilizer defined as above is frequently a significant over-approximation of the runtime destabilizer. For example, a query may be sensitive to update sequences that can only affect a node via a particular update operation. In its simplest form, the static destabilizer cannot capture this aspect. The authors presented a query-rewriting technique that provides a useful, sound approximation. Destabilizers sensitive to update operations (precise for each update), as well as destabilizers sensitive to different equalities (precise for boolean, node and subtree equivalence) are also defined in [BC10].

Benedikt and Cheney showed that calculating an *exact* static destabilizers is not feasible even for a core XML query language, hence the analysis needs to abstract away from some features in order to be effectively executed. This is stated by the following theorem.

**Theorem 3.2.1** ([BC10]). *There is no elementary time algorithm for constructing an operation-sensitive minimal static destabilizer.*

The formal definition of a destabilizer is reminiscent of a path analysis. The main difference with the path-based approach is that destabilizers also consider the source document for the overlapping test. To illustrate, consider the following query  $q$  and update  $u$  over the previous document  $t$ .

$$q = \text{for } x \text{ in } /bib \text{ return } \langle \text{list} \rangle x /book \langle /list \rangle \quad u = \text{delete } //book/title$$

The destabilizer for the query is

$$\Delta^s(q, t) = \{ /bib, /bib/book/descendant-or-self :: * \}$$

Now, while running the intersection test between  $u$  and  $\Delta^s(q, t)$ , we consider also the actual data of  $t$ .

The intersection tests of the destabilizers approach is more accurate than the one of previous techniques. While the path analysis (and also the type analysis) would exclude independence, the destabilizer analysis correctly deems the independence, because there is no data in the store corresponding to *title* nodes. In fact, the path analysis (and also the type analysis) does not take into account the document. This is because the path analysis would see that the path targeting delete nodes `//book/title/descendant-or-self :: *` intersects with path `/bib/book/descendant-or-self :: *` (inferred for  $q$ ).

The source document is taken into account by translating queries into logical formulas, and then transforming the intersection test of query and update paths to a satisfiability problem, solved with SAT-solvers. Several exact and approximate translations are proposed. These allow to trade between precision and efficiency of the analysis. In the following, we limit ourselves to discuss translations, without reporting the formal development, that can be found in [BC10].

We discuss now exact solutions. Queries without element construction can be translated to first-order formulas over the child, descendant and sibling relations, and then disjointness analysis for these queries reduces to satisfiability for first-order formulas over trees  $FO(Trees)$ . However it follows from results of [Vor96] that there is a non-elementary lower bound on the complexity of satisfiability. The solver used for resolution of this formula is Yices [DM06]. Another strategy is to convert expression into monadic second order logic (MSO). It follows from results in [SM73] that there is a non-elementary lower bound on the complexity of satisfiability of MSO over trees. The solver used for MSO resolution is MONA [KM11].

We discuss approximate solutions. To check intersection of query expressions on database instances, the authors propose first a reduction to existential second order logic (ESO) which uses as support an arithmetic encoding of the input document. This solution has an overhead that depends on the document to analyze, and this could not be reasonable for large documents. We recall that path intersection is NP-hard in general, even if for downward-only paths, overlapping is exactly decidable in quadratic time.

In [JGL12] Junedi, Genevès and Layida showed that the efficiency of destabilizer can be enhanced wrt the work of [BC10] by checking satisfiability of logical formulas using their  $\mu$ -solver [GLS07, GL08]. Their approach is based on translating XPath expressions into  $\mu$ -calculus formulas. As outlined above, the theoretical complexity of satisfiability on MSO formulas is known to be non-elementary, whereas it is exponential for simple



based linear arithmetic satisfiability. The complexity of satisfiability on  $\mu$ -formulas is simple exponential. The resulting approach resulted to be far more competitive wrt solver such as MONA and Yices. An evidence from experiments in [JGL12] is that the dominant inference cost does not necessarily reside in the resolution procedure, but mostly in constructing destabilizers, parsing formulas and initializing the auxiliary libraries for resolution. The time spent in the resolution procedure better reflects the intrinsic difficulty of each problem instance.

In conclusion, the merit of destabilizers is twofold. It is a modular framework that generalizes concepts of all previous approaches to independence. It promotes the practical use of SAT-solvers for intersection analysis. The approach does not require a schema, even if they can be handled by adding schema constraints in the satisfiability test. While former experiments in [BC10] showed some limits of the approach concerning the cost of the destabilizer analysis, recent works of [JGL12] showed that destabilizers can be implemented efficiently.

## Conclusions

In this chapter we presented the theoretical background and the state-of-the-art techniques for checking XML query update independence. We discussed theoretical results that show the computational difficulty of the problem. Then we discussed three approximated analysis that, being sound but not complete, try to make the analysis tractable. All approaches distinguish for simplicity, precision and complexity. The type-based analysis [BC09a] requires a schema in order to be executed, it is simple to implement and it runs in polynomial time, but it may be not precise. The path-based analysis [GRS08] works without requiring a schema but when structural constraints are available they cannot be easily integrated. It is also simple to implement but for some expression the path extraction may occur in blowups, and also exact checking of path disjunction turns out to be NP-hard. The destabilizer analysis [BC10] is a modular framework that generalizes all previous ideas, and allow to perform an independence analysis varying precision and complexity of the approach. The method works without schemas, but it can easily integrate structural constraints.

In the next chapters we present our technique for detecting independence based on schemas, showing how it can ensure a precise analysis while still keeping inference polynomial in time and space.

## Chapter 4

# Independence Analysis based on Schema-Chains

This chapter presents the main contribution of this thesis: a novel method for detecting XML query update independence, in the presence of a schema. We present an inference system for computing the set of navigational paths accessed by queries and updates, called *chains*. We prove the chain-based inference system to be sound and more precise than the related approach based on types. We present a static notion of independence, based on chains, that we show to be decidable, even for recursive schemas which can generate infinite sets of chains.

The chapter is organized in four parts. In the first part, we present a system for chain inference. In the second part, we show the correctness of the type system. In the third part, we present a static notion of independence based on chains. In the fourth part, we show the decidability of the notion of static independence.

### 4.1 Introduction

As explained in the previous chapter, the type analysis for checking independence defined in [BC09a] is designed so as to infer information about accessed data, by abstracting away from some of the structural properties expressed by queries and updates. This is done in order to keep the systems relatively simple, and to make the inference fast. Actually, the main structural information which is traded is data-nesting, that defines the way in which types interact.

Data-nesting is a synonym of hierarchical information. Together with document order, it entails that the semantics of a subtree depends on its position in the document. Describing a collection of accessed and updated trees with a set of types-names, is equivalent to consider all possible trees with fixed tag alphabet, ignoring any constraint on the hierarchy of data. Unfortunately, in many cases this causes lack of precision because of the tree shape of XML data.

The goal of this work is to improve the precision of the method proposed in [BC09a] by exploiting schema information in a different way. Rather than inferring sets of type-names accessed by the expressions, our goal is to infer *sequences* of accessed types, called *chains*, that essentially correspond to navigational paths traversed by query and update expressions on valid schema instances. Chains will form the basis of an extremely precise independence analysis.

We present now formal definitions about chains. We begin by defining the notion of *type* and *chain* of types for a *node location* belonging to a store.

**Definition 4.1.1** (Node Type and Chain). *Given a store  $\sigma$ , the type of a location  $l \in \text{dom}(\sigma)$ , denoted by  $\text{type}(l)$  is defined as*

$$\text{type}(l) = \begin{cases} a & \text{if } \sigma(l) = a[L] \\ \text{String} & \text{otherwise} \end{cases}$$

*The chain associated to location  $l$ , denoted by  $c_l^\sigma$ , is defined by*

$$c_l^\sigma = \begin{cases} \text{type}(l) & \text{if } l \text{ has no parent} \\ c_{\text{parent}(l)}^\sigma.\text{type}(l) & \text{otherwise} \end{cases}$$

If  $t = (\sigma, l_t)$  and  $l \in \text{dom}(t)$  then we denote by  $c_l^t$  the chain  $c_l^t = c_l^\sigma$ . To illustrate, consider the tree  $t_2$  in Figure 4.1. Here the type of the root location  $l_{t_2}$  is  $\text{type}(l_{t_2}) = \text{bib}$ , and it also coincides with its own chain  $c_{l_{t_2}}^\sigma$ , because the root node has no parent. The type of  $l'_3$  and  $l'_4$  are the same, indeed  $\text{type}(l'_3) = \text{type}(l'_4) = \text{name}$ . However their chains differ since

$$c_{l'_3}^\sigma = \text{bib.book.author.name} \quad c_{l'_4}^\sigma = \text{bib.book.publisher.name}$$

Notice that since the store coincides with the tree, we have that  $c_l^\sigma = c_l^{t_2}$  for all  $l \in \text{dom}(\sigma)$ . As a convention in the notation, we use “.” to separate two consecutive labels in a chain. Notice also that this definition considers only *rooted* chains, where the starting symbol is always the root.

We stress that the above definition works for any XML tree and node location, independently from schema constraints. However, dealing with schemas is our ultimate goal. Hence in order to develop a static analysis which is correct we need to ensure that, if a document is compliant with a schema, then all node chains in the document are induced by the schema.

To formalize this property, we define the set of chains induced by a schema relying on a reachability relation between types.

**Definition 4.1.2** (Reachability and Chains). *Given a DTD  $d$ , we denote by  $a \Rightarrow_d b$  the reachability relation that holds for a pair  $a, b \in \Sigma_S$  if  $b$  occurs in the regular expression  $d(a)$ . A chain  $c$  over  $d$  is a non-empty sequence of labels*

$$a_1.a_2 \dots .a_n$$

*such that  $a_i \Rightarrow_d a_{(i+1)}$  for  $i = 1 \dots n-1$ .*

<pre> &lt;bib&gt;   &lt;book&gt;     &lt;title&gt;       Types and Programming Languages     &lt;/title&gt;     &lt;author&gt;       &lt;name&gt;B. Pierce&lt;/name&gt;     &lt;/author&gt;     &lt;publisher&gt;       &lt;name&gt;MIT Press&lt;/name&gt;     &lt;/publisher&gt;   &lt;/book&gt; &lt;/bib&gt; </pre>	$\sigma = \begin{cases} l_{t_2} \leftarrow \text{bib}[l_1] \\ l_1 \leftarrow \text{book}[l_2, l_3, l_4] \\ l_2 \leftarrow \text{title}[l'_2] \\ l'_2 \leftarrow \text{"Types and ..."} \\ l_3 \leftarrow \text{author}[l'_3] \\ l'_3 \leftarrow \text{name}[l''_3] \\ l''_3 \leftarrow \text{"B. Pierce"} \\ l_4 \leftarrow \text{publisher}[l'_4] \\ l'_4 \leftarrow \text{name}[l''_4] \\ l''_4 \leftarrow \text{"MIT Press"} \end{cases}$
---	--

Figure 4.1: Document  $t_2$  and store  $(\sigma, l_{t_2})$ 

The set of chains associated with the DTD  $d$  is denoted by  $C_d$ . Observe that chains in  $C_d$  are of finite length and may start with any DTD symbol. The set  $C_d$  is infinite only if  $d$  is a vertical-recursive schema (i.e., some type is defined in terms of the type itself). The empty chain is denoted by  $\epsilon$ , and we assume, as a convention, that it is never induced by any schema. Given two chains  $c_1$  and  $c_2$ , the concatenation of  $c_1$  and  $c_2$  is denoted  $c_1.c_2$ , of course  $c_1.\epsilon = c_1$ . We write  $c_1 \leq c_2$  to indicate that  $c_1$  is a prefix of  $c_2$ , that is  $c_2 = c_1.c$  for some chain  $c$ .

As an example of schema chains, consider the DTD  $d_3$  in Figure 3.1. The set  $C_{d_3}$  includes the chains

`bib.book.title.String`   `bib.book.author.name.String`   `bib.book.publisher.name.String`

and all their non-empty subsequences. For instance, `bib` and `author.name.String` are  $C_{d_3}$ -chains, while `bib.author` and `book.String` are not.

Relation  $\Rightarrow_d$  is both sound and complete. Soundness means that for any pair of parent-child locations belonging to a valid tree, there exists a correspondent pair of types in  $\Rightarrow_d$ . Completeness means that the relation is also minimal, thus removing one pair from  $\Rightarrow_d$  we break soundness. This is stated by the following lemmas. Proofs are reported in Chapter 9.

**Lemma 4.1.3** (Soundness of  $\Rightarrow_d$ ). *Let  $d$  be a DTD and  $t = (\sigma, l_t) \in d$  a valid tree. If  $(l, l') \in \text{Child}_t$  then  $\text{type}(l) \Rightarrow_d \text{type}(l')$ .*

**Lemma 4.1.4** (Completeness of  $\Rightarrow_d$ ). *For all DTD  $d$ , if  $a \Rightarrow_d b$  then there exists a valid tree  $t = (\sigma, l_t) \in d$ , and a pair of locations  $(l, l') \in \text{Child}_t$ , such that  $\text{type}(l) = a$  and  $\text{type}(l') = b$ .*

At this point, we can reconcile chains associated to node locations (Definition 4.1.1) and chains induced by a schema (Definition 4.1.2) by stating that, when we consider a document valid against a schema, all node-chains *are* schema-chains.

**Proposition 4.1.5.** *Given a valid tree  $t \in \mathbf{d}$  and a location  $l \in \text{dom}(t)$  we have  $c_l^t \in \mathbf{C}_d$ .*

*Proof.* Immediate by Lemma 4.1.3. □

The converse of this proposition is not true. Not all of the schema-chains are node-chains of a tree, for *all* valid schema instance. This is because some elements may be optional. Therefore, what can be proved is that for all *rooted* schema chain there *exists* a document instance including a node with that chain. This, as a corollary of Lemma 4.1.4.

At this point, we formally defined types and chains for documents and schemas. We have shown that schema chains can be safely used in order to abstract over tree chains (Proposition 4.1.5), as a corollary of validity. The sequel of the chapter is dedicated to define an independence analysis based on schema-chains.

## 4.2 Chain Inference

We begin by defining chain inference for XPath navigational steps, we then generalize to queries and updates.

### 4.2.1 Step Chain Inference

The definition of our chain inference system makes the assumption that the inference is made starting from an input set of chains  $\mathbf{C}$ . Given a DTD  $\mathbf{d}$ , this set can be either  $\mathbf{C}_d$  (possibly infinite analysis) or a finite subset of  $\mathbf{C}_d$  (finite analysis). We would like to stress that assuming a pre-computed chain set is only made to ease the formal presentation. Any reasonable implementation can avoid this, by inferring chains in  $\mathbf{C}$  on the fly (see Chapter 5).

The first ingredient for query/update chain inference is chain inference for a single XPath step. We first define chain inference for axes, and then for node tests. Axis chain inference aims at inferring all chains that can be generated by axis navigation, in a  $\mathbf{d}$  instance, starting from a node typed by a chain  $c \in \mathbf{C}$ . Chain inference for XPath axes is defined below, where we assume that  $c', c'' \neq \epsilon$ .

$$\begin{aligned}
 \mathbf{A}_C(c, \text{self}) &\stackrel{\text{def}}{=} \{ c \} \\
 \mathbf{A}_C(c, \text{child}) &\stackrel{\text{def}}{=} \{ c.a \mid c.a \in \mathbf{C} \} \\
 \mathbf{A}_C(c, \text{descendant}) &\stackrel{\text{def}}{=} \{ c.c' \mid c.c' \in \mathbf{C} \} \\
 \mathbf{A}_C(c, \text{parent}) &\stackrel{\text{def}}{=} \{ c' \mid c = c'.a \} \\
 \mathbf{A}_C(c, \text{ancestor}) &\stackrel{\text{def}}{=} \{ c' \mid c = c'.c'' \}
 \end{aligned}$$

Chain inference rules strictly mimic XPath semantics of axes. It is worth noticing that all navigations follows from the reachability relation  $\Rightarrow_d$ . Notice also that

$$\mathbf{A}_C(c, \text{axis-or-self}) = \mathbf{A}_C(c, \text{axis}) \cup \mathbf{A}_C(c, \text{self})$$

for  $\text{axis} \in \{\text{descendant}, \text{ancestor}\}$ .

We illustrate step chain inference over the DTD  $d_3$  in Figure 3.1, with the following examples.

$$\begin{aligned} \mathbf{A}_{C_{d_3}}(\text{bib.book}, \text{child}) &= \{ \text{bib.book.title} \\ &\quad \text{bib.book.author} \\ &\quad \text{bib.book.publisher} \} \\ \mathbf{A}_{C_{d_3}}(\text{bib.book.author}, \text{descendant}) &= \{ \text{bib.book.author.name} \\ &\quad \text{bib.book.author.name.String} \} \\ \mathbf{A}_{C_{d_3}}(\text{bib.book.author}, \text{parent}) &= \{ \text{bib.book} \} \end{aligned}$$

Chain inference for sibling axes has a slightly more involved definition. First we need to define a relation between sibling types. This is needed in order to deal with horizontal axis.

**Definition 4.2.1** (Sibling-type relation). *Two types  $a, b$  are said to be siblings wrt a regular expression  $r$ , denoted by  $a <_r b$  if there exists a word  $u$  belonging to the language generated by  $r$  in which a letter  $a$  occurs before a letter  $b$ . The relation  $<_r$  is inductively defined as follows*

$$\begin{aligned} <() &\stackrel{def}{=} \emptyset \\ <_a &\stackrel{def}{=} \emptyset \\ <_{r_1|r_2} &\stackrel{def}{=} <_{r_1} \cup <_{r_2} \\ <_{r_1, r_2} &\stackrel{def}{=} <_{r_1} \cup <_{r_2} \cup \text{Sym}(r_1) \times \text{Sym}(r_2) \\ <_{r_1^+} &\stackrel{def}{=} <_{r_1, r_1} \\ <_{r_1^*} &\stackrel{def}{=} <_{r_1, r_1} \end{aligned}$$

where  $\text{Sym}(r)$  is the set of symbols used in the regular expression  $r$ .

In the following rules, with a little abuse of notation, given a chain  $c$  on  $d$ , we use  $d(c)$  to indicate either the regular expression  $d(a)$ , when  $c = c'.a$ , or the empty regular expression  $\epsilon$ , when  $c = c'.\text{String}$ . Chain inference for preceding/following-sibling axes is defined as follows.

$$\begin{aligned} \mathbf{A}_C(c, \text{following-sibling}) &\stackrel{def}{=} \{ c'.b \in C \mid c = c'.a, a <_{d(c')} b \} \\ \mathbf{A}_C(c, \text{preceding-sibling}) &\stackrel{def}{=} \{ c'.a \in C \mid c = c'.b, a <_{d(c')} b \} \end{aligned}$$

It is worth noticing that this sibling typing is very precise, since the sibling relation is evaluated for a particular regular expression. To illustrate, let  $r_1, r_2$  be two regular expressions in  $\mathbf{d}$ . It may be that  $a <_{r_1} b$  but  $a \not<_{r_2} b$ , and this difference can be captured using chains. This is not the case for the type analysis in [BC09a] that would consider  $a$  and  $b$  as sibling types in all cases.

Step chain inference for sibling axes over DTD  $\mathbf{d}_3$  is exemplified below.

$$\begin{aligned} \mathbf{A}_{\mathbf{C}_{\mathbf{d}_3}}(\text{bib.book.title, following-sibling}) &= \{ \text{bib.book.author} \\ &\quad \text{bib.book.publisher} \} \\ \mathbf{A}_{\mathbf{C}_{\mathbf{d}_3}}(\text{bib.book.title, preceding-sibling}) &= \emptyset \end{aligned}$$

As a more involved example of the sibling relation  $<_r$  consider for instance

$$<_{a, (b|f)^+, g} = \{(b, b), (b, f), (f, b), (f, f), (a, b), (a, f), (a, g), (b, g), (f, g)\}$$

since

$$\begin{aligned} <_{a, (b|f)^+, g} &= <_{a, (b|f)^+} \cup <_g \cup \{a, b, f\} \times \{g\} \\ &= <_a \cup <_{(b|f)^+} \cup \{a\} \times \{b, f\} \cup \{a, b, f\} \times \{g\} \\ &= <_{(b|f), (b|f)} \cup \{a\} \times \{b, f\} \cup \{a, b, f\} \times \{g\} \\ &= <_{b|f} \cup \{b, f\} \times \{b, f\} \cup \{a\} \times \{b, f\} \cup \{a, b, f\} \times \{g\} \\ &= <_b \cup <_f \cup \{b, f\} \times \{b, f\} \cup \{a\} \times \{b, f\} \cup \{a, b, f\} \times \{g\} \\ &= \{b, f\} \times \{b, f\} \cup \{a\} \times \{b, f\} \cup \{a, b, f\} \times \{g\} \\ &= \{(b, b), (b, f), (f, b), (f, f), (a, b), (a, f), (a, g), (b, g), (f, g)\} \end{aligned}$$

It the development we assumed left associativity for union and concatenation, that implies  $r_1, r_2, r_3 = (r_1, r_2), r_3$ . It is worth noticing that the relation  $<_r$  considers all following siblings of a given type, like  $(a, g)$ , and not just the immediate ones. Also,  $<_a = <_b = <_f = <_g = \emptyset$ .

The following two lemmas state that the relation  $<_r$  is both sound and complete.

**Lemma 4.2.2** (Soundness of  $<_r$ ). *Let  $\mathbf{d}$  be a DTD and  $t = (\sigma, l_t) \in \mathbf{d}$  a valid tree. If  $(l, l') \in \text{FollowingSibling}_t$  and  $c_l^t = c_{l'}^t = c.a$  then  $\text{type}(l) <_{\mathbf{d}(\mathbf{c})} \text{type}(l')$ .*

**Lemma 4.2.3** (Completeness of  $<_r$ ). *For all DTD  $\mathbf{d}$ , If  $a <_{\mathbf{d}(\mathbf{c})} b$  then there exists a valid tree  $t = (\sigma, l_t) \in \mathbf{d}$ , and a pair of locations  $(l, l') \in \text{FollowingSibling}_t$ , such that  $\text{type}(l) = a$  and  $\text{type}(l') = b$ , with  $c_l^t = c.a$  and  $c_{l'}^t = c.b$ .*

Rules for node-test chain inference are easier and defined below.

$$\begin{aligned} T_C(c, \text{node}()) &\stackrel{\text{def}}{=} \{ c \} \\ T_C(c, a) &\stackrel{\text{def}}{=} \{ c \mid c = c'.a \} \\ T_C(c, \text{text}()) &\stackrel{\text{def}}{=} \{ c \mid c = c'.\text{String} \} \end{aligned}$$

Composition of axis and node-filters is noted as

$$T_C(\mathbf{A}_C(c, \text{axis}), \phi)$$

that stands for

$$\bigcup_{c' \in \mathbf{A}_C(c, \text{axis})} T_C(c', \phi)$$

To illustrate, consider the following examples over DTD  $d_3$

$$\begin{aligned} T_{C_{d_3}}(\mathbf{A}_{C_{d_3}}(\text{bib.book.author}, \text{descendant}), \text{text}()) &= \{ \text{bib.book.author.name.String} \} \\ T_{C_{d_3}}(\mathbf{A}_{C_{d_3}}(\text{bib.book.title}, \text{following-sibling}), \text{author}) &= \{ \text{bib.book.author} \} \end{aligned}$$

In this section we defined chain inference for XPath steps, the building blocks of our analysis. We will define next chain inference for query and update expressions. It is worth noting that in the above examples we considered only *rooted* chains, where the first symbol is the root type. This is because these chains are the ones that make sense for the access analysis, despite the fact that all definitions work also for non-rooted chains.

#### 4.2.2 Query Chain Inference

The analysis of query and update chains requires to further structure the inference process, in order to better describe the behavior of query and update expression, and prepare for a more precise independence analysis. For this aim, we introduce several classes of chains that we use to capture accessed and updated data.

As output of inference, our system produces chains of different kinds. The classification resembles that of Marian and Simeon in [MS03] for query path extraction, and of Ghelli et al. in [GRS08] for update path extraction, and is needed in order to reflect different way a query manipulates nodes during its evaluation. Reflecting this distinction is crucial for soundness and precision of the independence analysis.

In our framework, a query chain belongs to one of the following three disjoint classes.

- **Return chains** type input document nodes (*return nodes*) that are roots of subtrees outputted by the query. All descendants of a return node are in the query result, thus a return chain  $c$  implicitly embodies these descendants. Now, if a change made by an update  $u$  targets a *return node* or some of its ancestors or descendants, query-update independence is not guaranteed.



- **Used chains** type nodes (*used nodes*) belonging to the input document and participating to the query evaluation, without necessarily being part of the result itself. Clearly, if a change of an update  $u$  targets a *used node* or some of its ancestors, then query-update independence is not guaranteed. However, if it targets a used node descendant then independence may still hold.
- **Element chains** type newly constructed elements; an element chain is of the form  $\mathbf{a.c}$ , where  $\mathbf{a}$  is the tag of the constructed  $\mathbf{a}$  element, while  $\mathbf{c}$  traverses nodes belonging to the subtree rooted at the  $\mathbf{a}$  element. Extracting these chains is important to record the structure of forests built and inserted by update expressions.

For updates, we have one class of chains:

- The purpose of an **Update chain**, denoted by  $\mathbf{c} : \mathbf{c}'$ , is twofold:  $\mathbf{c}$  types nodes  $l$  whose content may be changed by the update and  $\mathbf{c}'$ , rooted at a child label of  $l$ , types descendants of  $l$  (either introduced or removed by the update) involved in the changes. For example, given  $\mathbf{c} : \mathbf{c}'$ , independence is not guaranteed if a query returns an element whose root is typed by  $\mathbf{c.a}$ , with  $\mathbf{a}$  a prefix of  $\mathbf{c}'$ .

Rules for query chain inference are presented in Table 4.1. These rules prove judgements of the form:

$$\Gamma \vdash_{\mathbf{C}} \mathbf{q} : (\mathbf{r}, \mathbf{v}, \mathbf{e})$$

meaning that starting from a variable environment  $\Gamma$  and a set of chains  $\mathbf{C}$ , the chain inference for query  $\mathbf{q}$  produces the sets  $\mathbf{r}$ ,  $\mathbf{v}$  and  $\mathbf{e}$ , containing the return, used and element chains for  $\mathbf{q}$ , respectively.  $\Gamma$  is a static variable environment that associates each query free-variable  $\mathbf{x}$  with a set  $\Gamma(\mathbf{x})$  of chains, typing nodes that can be assigned to the variable during query evaluation. When  $\mathbf{x}$  is bound to a single chain, we write  $\mathbf{x} \mapsto \mathbf{c}$  as a shorthand of  $\mathbf{x} \mapsto \{\mathbf{c}\}$ .

All the rules mimic query semantics [DFF<sup>+</sup>10, BC09b] introduced in Section 2.3. Rule (EMPTY) returns no chain, as the empty query performs no operation. Rule (CONCAT) returns the union of return, used and element chains for the two expressions.

Navigational steps are captured by rule (STEPF) and (STEPUH) dealing with forward and upward/horizontal navigations respectively. Rule (STEPF) produces only return chains, by using in turn step chain inference rules defined in the previous section. This is due to the fact that forward navigations only *extend* input chains, without losing track of visited ancestor types. Differently, rule (STEPUH) produces also used chains. These are input chains leading to a non-empty result for the axis inference. These chains are needed for developing a correct access analysis, since return chains produced by an horizontal/upward step may not contain as a prefix the input chain in  $\Gamma(\mathbf{x})$  from which they have been generated. Also, these input chains are kept as *used* chains because the descendants of the locations they are typing do not need to be accessed as well.

$\frac{}{\Gamma \vdash_C () : (\emptyset, \emptyset, \emptyset)} \text{ (EMPTY)}$	$\frac{\Gamma \vdash_C q_i : (r_i, v_i, e_i) \quad i=1, 2}{\Gamma \vdash_C q_1, q_2 : (r_1 \cup r_2, v_1 \cup v_2, e_1 \cup e_2)} \text{ (CONC)}$
$\frac{\begin{array}{l} \text{axis} \in \{\text{self}, \text{child}, \text{descendant-or-self}\} \\ r_c = \mathbf{T}_C(\mathbf{A}_C(c, \text{axis}), \phi) \text{ for any } c \in \Gamma(x) \end{array}}{\Gamma \vdash_C x/\text{axis}::\phi : (\bigcup_{c \in \Gamma(x)} r_c, \emptyset, \emptyset)} \text{ (STEPF)}$	
$\frac{\begin{array}{l} \text{axis} \in \{\text{parent}, \text{ancestor}, \text{ancestor-or-self}, \\ \text{following-sibling}, \text{preceding-sibling}\} \\ r_c = \mathbf{T}_C(\mathbf{A}_C(c, \text{axis}), \phi) \text{ for any } c \in \Gamma(x) \end{array}}{\Gamma \vdash_C x/\text{axis}::\phi : (\bigcup_{c \in \Gamma(x)} r_c, \bigcup_{\substack{c \in \Gamma(x) \\ r_c \neq \emptyset}} \{c\}, \emptyset)} \text{ (STEPUH)}$	
$\frac{\begin{array}{l} \Gamma \vdash_C q_1 : (r_1, v_1, e_1) \\ \Gamma[x \mapsto c] \vdash_C q_2 : (r_c, v_c, e_c) \text{ for any } c \in r_1 \end{array}}{\Gamma \vdash_C \text{for } x \text{ in } q_1 \text{ return } q_2 : (\bigcup_{c \in r_1} r_c, v_1 \cup \bigcup_{\substack{c \in r_1 \\ r_c \cup e_c \neq \emptyset}} (v_c \cup \{c\}), \bigcup_{c \in r_1} e_c)} \text{ (FOR)}$	
$\frac{\Gamma \vdash_C q_1 : (r_1, v_1, e_1) \quad \Gamma[x \mapsto r_1] \vdash_C q_2 : (r_2, v_2, e_2)}{\Gamma \vdash_C \text{let } x := q_1 \text{ return } q_2 : (r_2, r_1 \cup v_1 \cup v_2, e_2)} \text{ (LET)}$	
$\frac{\Gamma \vdash_C q_i : (r_i, v_i, e_i) \quad i=0..2}{\Gamma \vdash_C \text{if } (q_0) \text{ then } q_1 \text{ else } q_2 : (r_1 \cup r_2, \bigcup_{i=0..2} v_i \cup r_0, e_1 \cup e_2)} \text{ (IF)}$	
$\frac{}{\Gamma \vdash_C \text{"txt"} : (\emptyset, \emptyset, \{\text{String}\})} \text{ (TEXT)}$	
$\frac{\begin{array}{l} \Gamma \vdash_C q : (r, v, e) \\ e_0 = \{ a.a.c' \mid c.a \in r, c.a.c' \in \bar{r} \} \cup \{ a.c \mid c \in e \} \cup \{ a \mid r \cup e = \emptyset \} \end{array}}{\Gamma \vdash_C <a>q</a> : (\emptyset, \bar{r} \cup v, e_0)} \text{ (ELT)}$	

Table 4.1: Chain Inference Rules for Queries

The following examples are made on the DTD  $d_3$  in Figure 3.1, and we assume that variable `doc` is always bound to `bib` in  $\Gamma$ .

To illustrate rules (STEPF) and (STEPUH), consider , and the query

`doc//title/following-sibling::author`

Chain inference for this expression is done by analyzing the two steps in sequence. Rule (STEPF) infers chain `bib.book.title` for the step `doc//title`. By applying rule (STEPUH) on the step `following-sibling::author`, provided that  $\Gamma' = \Gamma \cup \{x \mapsto \text{bib.book.title}\}$ , we obtain the following derivation.

$$\frac{T_{C_{d_3}}(A_{C_{d_3}}(\text{bib.book.title}, \text{following-sibling}), \text{author}) = \{\text{bib.book.author}\}}{\Gamma' \vdash_{C_{d_3}} x/\text{following-sibling}::\text{author} : (\{\text{bib.book.author}\}, \{\text{bib.book.title}\}, \emptyset)}$$

The rule produces a return chain, namely `bib.book.author`, that is typing subtrees rooted at `author` returned by the query, but also a used chain, namely `bib.book.title`, that is typing locations accessed before getting to the result.

Rules (FOR) and (LET) are very similar, thus we mainly comment on the former. According to query semantics, rule (FOR) performs an iteration on the set of *return* chains inferred for  $q_1$  and, at each iteration step, it evaluates  $q_2$ . We illustrate this process through an example. Consider the DTD  $d_3$  and the following query

`for x in doc/descendant::name return x/parent::*`

The expression iterates on the sequence of document nodes labeled with *name*, and outputs the parent of the each node. Here, chain inference for the left subexpression `/descendant::name` produces two return chains  $c_a = \text{bib.book.author.name}$  and  $c_p = \text{bib.book.publisher.name}$ . Accordingly, chain inference for the right subexpression `x/parent::*` is done twice, with  $x$  bound to  $c_a$  and  $c_p$  respectively.

Moreover, three other operations are necessary in order to perform a correct and precise analysis.

First, *return* chains for  $q_1$  are converted into *used* chains for the whole expression. This is needed because chain inference is a bottom-up process: inside  $q_1$  a path expression is seen as a query producing a *result* (and as such it *locally* produces return chains), while it only selects nodes to be *used* in the outer iteration `for x in  $q_1$  return  $q_2$` . To illustrate, consider the query

`for x in doc//publisher return  
for y in doc//name return y`

In this case the return chain `bib.book.publisher` inferred for the step `//publisher`, should not be considered as a return chain for the whole expression, since the descendants of

nodes labeled with *publisher* are not accessed by the expression. Rather, it should be considered as a used chain for the whole expression.

Second, return chains for  $q_1$  that are irrelevant for  $q_2$  are filtered out. To illustrate, consider the query

```
for x in doc//node() return
if (x/title) then x/author else ()
```

Chain inference for the expression `doc//node()` returns all rooted chains of the schema. If these are all kept in the final result, the analysis is jeopardized since the query would be considered as dependent wrt almost every update. Our system avoids this, by discarding all return chains of `doc//node()` that do not further satisfy `x/title`. Chain inference for the whole expression only produces the used chain `bib.book.title` and the return chain `bib.book.author`.

Third, used chains for  $q_1$  become chains for the whole expression. This operation is needed in order to ensure a sound analysis. Consider the DTD  $d_3$  and following query

```
for x in doc//title/following-sibling::author return x
```

here the left query `doc//title/following-sibling::author` makes the system inferring a return chain `bib.book.author` and a used chain `bib.book.title`. Now, in order to have a sound analysis, the chain `bib.book.title` should be inferred as a used chain for the whole expression.

Rule (IF) is similar to rule (CONCAT). It invokes chain inference for each subexpression, and then outputs the proper union of sets inferred for the subexpressions. As an example consider the query

```
for x in doc//book return
if (x/price) then x/author else x/title
```

This case illustrates well that a correct analysis needs the chains inferred for all subexpressions, namely `bib.book.price`, `bib.book.author`, `bib.book.title`. Return chains for the then-branch (`bib.book.author`) and the else-branch (`bib.book.title`) become return chains for the whole expression. Differently, the return chain inferred for the boolean conditional-query (`bib.book.price`) is converted into a used chains for the whole expression, since this is just needed to perform boolean tests on the existence of some nodes, without requiring the whole subtrees.

Element queries `<a>q</a>` are dealt with by rule (ELT). In the element construction rule,  $\bar{\tau}$  denotes all descendant chains of chains in the set  $\tau$  wrt  $C$ .

$$\bar{\tau} \stackrel{def}{=} \{ c.c' \in C \mid c \in \tau \}$$

This rule infers element chains of the form `a.c`, where `c` is obtained from either an element or return chain of `q`. The rule also infers a set of used chains, composed by *i*) used

chains of  $q$  and *ii*) return chains of  $q$  converted into used chains for the whole expression. To this end  $\bar{r}$  is used to extend returned chains of the inner query. This return-to-used chain conversion is needed to correctly handle nested element construction. For instance, consider the following query  $q = \text{for } x \text{ in doc//book return } \langle \text{livre} \rangle q' \langle / \text{livre} \rangle$  where

$$\begin{aligned} q' = & \ x / \text{title} , \\ & \langle \text{auteur} \rangle \\ & \ x / \text{author} / \text{name} \\ & \langle / \text{auteur} \rangle \end{aligned}$$

The expression restructures the result of a query selecting some informations on the books in the bibliography, changing both the structure and the label of nodes of a tree labeled with *book*. Element chains for  $q$  are inferred in terms of chains for  $q'$ . So, element chains for  $q$  are *livre.title* and *livre.auteur.name*, given that for  $q'$  the return chain is *bib.book.title* (for  $x / \text{title}$ ) and the element chain is *auteur.name* (for  $\langle \text{auteur} \rangle x / \text{author} / \text{name} \langle / \text{auteur} \rangle$ ). In order to avoid ending up with a wrong element chain for  $q$ , i.e., *title.name* the return chain *bib.book.author.name* inferred for  $x / \text{author} / \text{name}$  does not have to belong to the set of return chains for  $q'$ . At the same time, we must record the fact that the whole subtree rooted at *name* is accessed. This is handled by the return-to-used conversion of the return chain *bib.book.author.name* when inferring chains for  $\langle \text{auteur} \rangle x / \text{author} / \text{name} \langle / \text{auteur} \rangle$ . The return-to-used conversion extends  $r$ , the set of return chains of the inner query  $x / \text{author} / \text{name}$ , with all possible continuations by means of  $\bar{r}$  and then it converts chains in  $\bar{r}$  to used chains for the whole expression.

It is worth stressing that if we just convert return chains to used ones without the extension  $\bar{r}$ , then we loose their semantic property of representing entire subtrees of data. Notice that this extension is needed for the purpose of the formal presentation although any efficient implementation can avoid performing these extensions by using intensional representations.

Rule (TEXT) deals with expressions building new text nodes. The rule simply infers *String* as an element chain for the built textual node. With a little abuse of notation, for simplicity, we preferred not to use another class of chains.

In this section we illustrated chain inference rules for query expressions. As update expressions are defined in terms of queries, we will see, in the next section, that chain inference rules for updates rely on the rules for query chain inference.

### 4.2.3 Update Chain Inference

Rules for update chain inference are presented in Table 4.2. These rules prove judgments of the form

$$\Gamma \vdash_C u : U$$

meaning that starting from a static variable environment  $\Gamma$  and a set of chains  $C$ , the chain inference for the update  $u$  produces the set of update chains  $U$ .

As seen before, update chains are of the form  $c : c'$ . Essentially, the *prefix*  $c$  types updated nodes, that are nodes whose children are modified by the update, while the *suffix*  $c'$  types modified children or new inserted/replaced descendants.

The most important inference rules are the ones corresponding to delete, rename, insert and replace operations, which are explained in the sequel.

```

      bib  ←  book*
      book ←  title?, author*, publisher*
author, publisher ←  name*
      title, name ←  String?

```

Figure 4.2: DTD  $d_4$

The following examples are made on the DTD  $d_4$ , and again we assume that variable **doc** is always bound to **bib** in  $\Gamma$ . We recall that when the update target query returns a sequence with more than one location, this has to be bound to a variable and then the update command has to be applied by means of iteration, otherwise a dynamic error is raised. However for the sake of conciseness, and w.l.o.g., in the following examples we do not make iteration explicit and assume that the update target query returns only one location.

Inference for delete expressions is defined by the (DELETE) rule, which simply puts the separator “:” just before the last symbol of each return chain of the target query. A delete chain  $c : a$  captures the fact that a node typed by  $c$  has a child labeled by  $a$  which may be deleted. As an example, consider the update

**delete doc//publisher**

By applying rule (DELETE) we obtain the following derivation

$$\frac{\Gamma \vdash_{C_{d_4}} \text{doc//publisher} : (\{\text{bib.book.publisher}\}, \emptyset, \emptyset)}{\Gamma \vdash_{C_{d_4}} \text{delete doc//publisher} : \{\text{bib.book} : \text{publisher}\}}$$

showing how return chains inferred for the target query become update chains.

For rename updates, the process is similar. Rule (RENAME) infers chains  $c : a$  where  $a$  is the tag of the target node before renaming, but it also produces chains  $c : b$ , where  $b$  is the tag of the target node after renaming. Recall that textual nodes are never target of rename commands. As an example, consider the DTD  $d_4$  and the rename update

**rename doc//author as publisher**

By applying rule (RENAME) we obtain the following derivation.

$$\frac{\Gamma \vdash_{\mathbf{c}_{d_4}} \text{doc//author} : (\{\text{bib.book.author}\}, \emptyset, \emptyset)}{\Gamma \vdash_{\mathbf{c}_{d_4}} \text{rename doc//author as publisher} : \{\text{bib.book} : \text{author}, \text{bib.book} : \text{publisher}\}}$$

In this case the update chains capture both the fact that subtrees rooted at **author** disappear and that subtrees rooted at **publisher** appear.

Chain inference for insert-into expressions (position ranges over into, first and last) is specified by the rule (INSERT-1). For any chain  $\mathbf{c} : \mathbf{c}'$  inferred, the prefix  $\mathbf{c}$  is a return chain of the target query  $\mathbf{q}_0$  (typing the insertion point), while the suffix  $\mathbf{c}'$  is either a return or element chain of the source expression  $\mathbf{q}$  typing a branch of a node element returned by  $\mathbf{q}$  itself; this element can either be a new one or a sub-element of the input document; in both cases the suffix chain corresponds to inserted data. As an example, consider the DTD  $\mathbf{d}_4$  and the update expression **insert** *<name/>* **into** *doc/self :: bib/book*. By applying rule (INSERT-1) we obtain the following.

$$\frac{\Gamma \vdash_{\mathbf{c}_{d_4}} \text{<name/>} : (\emptyset, \emptyset, \{\text{name}\}) \quad \Gamma \vdash_{\mathbf{c}_{d_4}} \text{doc/self :: bib/book} : (\{\text{bib.book}\}, \emptyset, \emptyset)}{\Gamma \vdash_{\mathbf{c}_{d_4}} \text{insert <name/> into doc/self :: bib/book} : \{\text{bib.book} : \text{name}\}}$$

Rule (INSERT-2) is similar, and deals with insert-before/after updates.

The two rules for insert updates rely on the fact that element chains can be extracted in source queries. The fact that we opportunely attach element chains of the source-expression to result chains of the target-expression, allow us to have a precise analysis. It is worth saying that related update analysis of [BC09a] and [GRS08] define as updated *everything* below the parent of the source location. This design choice has a strong impact for this case, since it is likely to make the analysis deem the update as dependent with almost all queries accessing information concerning the book, while it should only raise conflicts on *name* elements.

This last case suggests also that element chains are necessary for a precise independence analysis, when dealing with updates. To illustrate, consider the following update over the DTD in Figure 4.2.

```
for x in doc//book return
insert <author>q'</author> into x
```

Here the source expression is an element query, for which we infer element chains of the form *author.c'*, with  $\mathbf{c}'$  a chain inferred for  $\mathbf{q}'$ . The update chain *bib.book : author.c'* is obtained by concatenation of the chain *bib.book* associated with the target expression  $\mathbf{x}$ , and the chain *author.c'*. This allows one to conclude independence wrt the query

$$\begin{array}{c}
\frac{\Gamma \vdash_C q_0 : (r_0, v_0, e_0) \quad U = \{c : a \mid c.a \in r_0\}}{\Gamma \vdash_C \text{delete } q_0 : U} \text{(DELETE)} \quad \frac{\Gamma \vdash_C q_0 : (r_0, v_0, e_0) \quad U = \{c : a \mid c.a \in r_0\} \cup \{c : b \mid c.a \in r_0\}}{\Gamma \vdash_C \text{rename } q_0 \text{ as } b : U} \text{(RENAME)} \\
\\
\frac{\Gamma \vdash_C q : (r, v, e) \quad \Gamma \vdash_C q_0 : (r_0, v_0, e_0) \quad \text{pos} \in \{\text{into}, \text{first}, \text{last}\} \quad U = \{c : c' \mid c \in r_0, c'.c'' \in e, c' \neq e\} \cup \{c : a.c'' \mid c.a \in r_0, c'.a \in r, c'.a.c'' \in C\}}{\Gamma \vdash_C \text{insert } q \text{ pos } q_0 : U} \text{(INSERT-1)} \\
\\
\frac{\Gamma \vdash_C q : (r, v, e) \quad \Gamma \vdash_C q_0 : (r_0, v_0, e_0) \quad \text{pos} \in \{\text{after}, \text{before}\} \quad U = \{c : c' \mid c.a \in r_0, c'.c'' \in e, c' \neq e\} \cup \{c : b.c'' \mid c.a \in r_0, c'.b \in r, c'.b.c'' \in C\}}{\Gamma \vdash_C \text{insert } q \text{ pos } q_0 : U} \text{(INSERT-2)} \\
\\
\frac{\Gamma \vdash_C q : (r, v, e) \quad \Gamma \vdash_C q_0 : (r_0, v_0, e_0) \quad U = \{c : a \mid c.a \in r_0\} \cup \{c : b.c'' \mid c.a \in r_0, c'.b \in r, c'.b.c'' \in C\} \cup \{c : c' \mid c.a \in r_0, c'.c'' \in e, c' \neq e\}}{\Gamma \vdash_C \text{replace } q_0 \text{ with } q : U} \text{(REPLACE)} \\
\\
\frac{}{\Gamma \vdash_C u = () : \emptyset} \text{(TU-EMPTY)} \quad \frac{\Gamma \vdash_C u_1 : U_1 \quad \Gamma \vdash_C u_2 : U_2}{\Gamma \vdash_C u_1, u_2 : U_1 \cup U_2} \text{(TU-CONCAT)} \\
\\
\frac{\Gamma \vdash_C q : (r, v, e) \quad \Gamma[x \mapsto c] \vdash_C u : U_c \text{ for any } c \in r}{\Gamma \vdash_C \text{for } x \text{ in } q \text{ return } u : \bigcup_c U_c} \text{(TU-FOR)} \quad \frac{\Gamma \vdash_C q : (r, v, e) \quad \Gamma[x \mapsto r] \vdash_C u : U}{\Gamma \vdash_C \text{let } x := q \text{ return } u : U} \text{(TU-LET)} \\
\\
\frac{\Gamma \vdash_C q : (r, v, e) \quad \Gamma \vdash_C u_i : U_i \quad i = 1, 2}{\Gamma \vdash_C \text{if } (q) \text{ then } u_1 \text{ else } u_2 : U_1 \cup U_2} \text{(TU-IF)}
\end{array}$$

Table 4.2: Chain Inference Rules for Updates



`doc//title`, whose unique return chain is `bib.book.title` (forasmuch as *title* element is never a descendant of an *author* element): the update chain is not a prefix of the query chain and vice-versa.

Now, let us do the analysis without considering element chains: for the source expression `<author>q'</author>`, the best that can be done is to infer the chain

`bib.book : _`

telling that something happens beneath book elements. As a consequence, we would not deduce the independence. As said in the introduction, for similar reasons [BC09a] wrongly excludes independence for similar cases.

In the presence of nested element construction, the same remark holds. In the previous example, if `q'` is

`<first>Umberto</first>, <second>Eco</second>`

then by composing element chains during the inference, we end up with the following update chains `bib.book : author.first.String`, `bib.book : author.second.String`, `bib.book : author`, `bib.book : author.first` and `bib.book : author.second`. Indeed, this is necessary to exclude independence wrt the query `doc//author/email` (assuming the DTD allows for email elements into author elements).

The rule for replace expressions (REPLACE) is built on the same principles as (INSERT-2) and (DELETE) rules. The chains produced by the expression capture the fact that *replaced* subtrees are now missing, and that *replacing* subtrees can now be retrieved. As an example, consider the DTD `d4` and the rename update

`replace doc//author/name with doc//publisher/name`

By applying rule (REPLACE) we obtain the following.

$$\frac{\begin{array}{l} \Gamma \vdash_{C_{d_4}} \text{doc//author/name} : (\{\text{bib.book.author.name}\}, \emptyset, \emptyset) \\ \Gamma \vdash_{C_{d_4}} \text{doc//publisher/name} : (\{\text{bib.book.publisher.name}\}, \emptyset, \emptyset) \end{array}}{\Gamma \vdash_{C_{d_4}} \text{replace doc//author/name with doc//publisher/name} : \left\{ \begin{array}{l} \text{bib.book.author : name,} \\ \text{bib.book.author : name.String} \end{array} \right\}}$$

Once again the update chain is obtained by the concatenation of the target and the source expression. Notice that the update chain contains also the expansion of the return chain for the source query. This is done in order to completely trace inserted structures.

Remaining rules for the empty update, concatenation, iteration, let-binding and conditional expressions are formulated in the same spirit as for query chain inference.

In this section we presented rules for update chain inference, showing also that they are built on top of rules for query chain inference. We illustrated the precision of the

chain inference with examples compared with related work. In the next section, we will formally prove that query and update chain inference rules are correct, in the sense that they allow to infer types capturing all data which are accessed and updated by query and update expression.

### 4.3 Correctness and Precision of Chain Inference

In this section we show the correctness of the chain inference system for queries and updates. We begin by showing correctness of XPath step chain inference, for which completeness also hold. Then, we state the correctness for queries and updates. The proofs of all statements of this section can be found in the appendix 9.

Soundness of axis and node-test chain inference wrt query semantics is stated below.

**Lemma 4.3.1** (Soundness of Step Chain Inference). *Let  $\mathbf{d}$  be a DTD,  $t = (\sigma, l_t) \in \mathbf{d}$  be a valid tree, and  $l \in \text{dom}(\sigma)$  be a location of the tree. If*

$$\sigma, (\mathbf{x} := l) \models \mathbf{x}/\mathbf{axis} :: \phi \Rightarrow \sigma, L$$

*then for each location  $l' \in L$  we have that*

$$c_{l'}^t \in \mathcal{T}_C(\mathcal{A}_C(c_l^t, \mathbf{axis}), \phi)$$

A *complete* chain inference demands also the set of inferred chains to be minimal, so that each inferred chain is typing a result location belonging to some valid tree instance. The completeness of  $\Rightarrow_{\mathbf{d}}$  and  $<_{\mathbf{d}(\mathbf{a})}$  entails that step chain inference is minimal.

**Lemma 4.3.2** (Completeness of Step Chain Inference). *Let  $\mathbf{d}$  be a DTD and  $\mathbf{c} \in \mathcal{C}_{\mathbf{d}}$  a chain. If*

$$\mathbf{c}' \in \mathcal{T}_{\mathcal{C}_{\mathbf{d}}}(\mathcal{A}_{\mathcal{C}_{\mathbf{d}}}(\mathbf{c}, \mathbf{axis}), \phi)$$

*then there exists a valid tree  $t = (\sigma, l_t) \in \mathbf{d}$  and a location  $l \in \text{dom}(t)$  such that  $c_l^t = \mathbf{c}$  and, assuming that  $\sigma, (\mathbf{x} := l) \models \mathbf{x}/\mathbf{axis} :: \phi \Rightarrow \sigma, L$  we have that  $c_{l'}^t = \mathbf{c}'$ , for some  $l' \in L$ .*

Showing that chain inference is complete for navigational axes is easy, since chains carry a navigational context that leads inference rules to strictly mimic axis semantics even with backwards and horizontal navigations.

### 4.3.1 Soundness of Query Chain Inference

Proving soundness of query chain rules consists of proving that, for any schema instance, any node used or built by the query  $q$  is captured (typed) by the chains inferred for  $q$ . The proof relies on the notion of projection for XML [MS03, BCCN06].

A tree  $t'$  is said to be a *projection* of  $t$ , denoted by  $t' \leq t$ , if  $t'$  is obtained from  $t$  by removing some subtrees of  $t$ . According to our formalism, we define the projection of a tree  $t$  induced by a set of locations  $\mathcal{L} \subseteq \text{dom}(t)$ , provided that  $\mathcal{L}$  is non-empty and upward closed with respect to the node parent-child relationship in  $t$ . This simply means that if  $l \in \mathcal{L}$  and  $l'$  is the parent of  $l$  in  $\sigma$  then  $l' \in \mathcal{L}$ .

Given a sequence of locations  $L$ , we define the filtering of  $L$  wrt a set of node locations  $\mathcal{L}$ , denoted by  $L|_{\mathcal{L}}$ , as the subsequence of  $L$  that contains only identifiers in  $\mathcal{L}$  and that preserves the order of  $L$ . The projection of a tree  $t = (\sigma, l_t)$  wrt a set of locations  $\mathcal{L} \subseteq \text{dom}(\sigma)$ , is defined as  $t|_{\mathcal{L}} = (\sigma|_{\mathcal{L}}, l_t)$ , where

$$\sigma|_{\mathcal{L}} \stackrel{\text{def}}{=} \{ l \leftarrow a[L|_{\mathcal{L}}] \mid l \in \mathcal{L}, (l \leftarrow a[L]) \in \sigma \} \cup \{ l \leftarrow \text{"txt"} \mid l \in \mathcal{L}, (l \leftarrow \text{"txt"}) \in \sigma \}$$

That is the store obtained by removing the use of all element and text node locations not in  $\mathcal{L}$ .

Given a query  $q$ , we can define a notion of projection wrt  $q$ . We say that a store  $\sigma|_{\mathcal{L}}$  is a  $q$ -projection of  $\sigma$  if, provided that  $\sigma, \gamma \models q \Rightarrow \sigma, L$  and  $\sigma|_{\mathcal{L}}, \gamma \models q \Rightarrow \sigma|_{\mathcal{L}}, L'$  we can conclude

$$(\sigma, L) \cong (\sigma|_{\mathcal{L}}, L')$$

This definition enforces  $\mathcal{L}$  to contain all locations needed in order to evaluate  $q$ , without altering the result of  $q$ . Given two trees  $t = (\sigma, l)$  and  $t' = (\sigma', l)$ , if  $\sigma'$  is a  $q$ -projection of  $\sigma$  then  $t'$  is a  $q$ -projection of  $t$ .

A tree  $t'$  is said to be a *minimal*  $q$ -projection of  $t$  if for all  $t''$  *strict* projections of  $t'$  we have that  $t''$  is not a  $q$ -projection of  $t'$  (or, equivalently, of  $t$ ). Here strict means that  $\text{dom}(t'') \subsetneq \text{dom}(t')$ . Notice that, in general, a projection is not minimal. In the sequel we also show that a minimal projection is not unique, due to the query language considered.

To illustrate the above concepts, consider again the tree  $t_2$  in Figure 4.1 that we report below. A projection of  $t_2$  can be obtained by choosing the set of locations  $\mathcal{L} = \{l_{t_2}, l_1, l_2, l_3, l_4\}$ , that induces the projection of the book element with its own children nodes. However the sets  $\mathcal{L}_1 = \{l_1, l_2\}$  or  $\mathcal{L}_2 = \{l_{t_2}, l_2\}$  do not induce projections because they are not closed wrt the node parent-child relation, since  $\mathcal{L}_1$  misses  $l_{t_2}$  and  $\mathcal{L}_2$  misses  $l_1$ . Neither  $\mathcal{L} = \emptyset$  induces a projection because it contains no location.

Now, let us consider the query  $q = //title$ , selecting all titles of book in the store. By evaluating  $q$  over  $t_2$ , we get the subtree rooted at  $l_2$ , since there is only one book in  $t_2$ .

<pre> &lt;bib&gt;   &lt;book&gt;     &lt;title&gt;       Types and Programming Languages     &lt;/title&gt;     &lt;author&gt;       &lt;name&gt;B. Pierce&lt;/name&gt;     &lt;/author&gt;     &lt;publisher&gt;       &lt;name&gt;MIT Press&lt;/name&gt;     &lt;/publisher&gt;   &lt;/book&gt; &lt;/bib&gt; </pre>	$\sigma = \begin{cases} l_{t_2} \leftarrow \text{bib}[l_1] \\ l_1 \leftarrow \text{book}[l_2, l_3, l_4] \\ l_2 \leftarrow \text{title}[l'_2] \\ l'_2 \leftarrow \text{"Types and ..."} \\ l_3 \leftarrow \text{author}[l'_3] \\ l'_3 \leftarrow \text{name}[l''_3] \\ l''_3 \leftarrow \text{"B. Pierce"} \\ l_4 \leftarrow \text{publisher}[l'_4] \\ l'_4 \leftarrow \text{name}[l''_4] \\ l''_4 \leftarrow \text{"MIT Press"} \end{cases}$
---	--

Figure 4.3: Document  $t_2$  and store  $(\sigma, l_{t_2})$ 

Also, in order to reach the result, we need all ancestor locations of  $l_2$ , namely  $l_t$  and  $l_1$ . The set of locations  $\mathcal{L}_q = \{l_{t_2}, l_1, l_2, l'_2\}$  can be used in order to build a  $\mathbf{q}$ -projection of  $t_2$ . The result of  $\mathbf{q}$  over  $t_2|_{\mathcal{L}_q}$  is equal to  $\mathbf{q}$  over  $t_2$ . In fact, all other locations in the store are not needed by the query because they concern complementary information such as authors and publishers of a book.

It is worth noticing that in the definition of projection we employed a notion of forest-equivalence  $\cong$ , that checks only if the two resulting forests are isomorphic, regardless of the node identifiers. Above, we have an example of a stronger equivalence, since not only the structure of the results is the same, but also the node identifiers coincide.

To illustrate minimality of projection, consider again the tree  $t_2$ , a dynamic environment  $\gamma = \{\text{doc} \mapsto l_{t_2}\}$  and the query

```

q = for x in doc//book return
    if (x/author, x/publisher) then x/title else ()

```

that returns the title of a book either if this has an author or a publisher child node (that is, if the sequence  $(\mathbf{x}/\text{author}, \mathbf{x}/\text{publisher})$  is non-empty). For this expression, we have a number of  $\mathbf{q}$ -projections, namely

$$\mathcal{L}_1 = \{l_{t_2}, l_1, l_2, l'_2, l_3, l_4\} \quad \mathcal{L}_2 = \{l_{t_2}, l_1, l_2, l'_2, l_3\} \quad \mathcal{L}_3 = \{l_{t_2}, l_1, l_2, l'_2, l_4\}$$

and all backward-closed supersets of  $\mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_3$ . It could seem surprisingly that  $\mathcal{L}_2$  and  $\mathcal{L}_3$  are projections, because they both discard a node location that is used for building the sequence used by the conditional test. Anyway this is correct, since we need only one of  $l_3$  and  $l_4$  in order to ensure the conditional expression to be non-empty. In this case we have that  $\mathcal{L}_1$  is not minimal simply because  $\mathcal{L}_2, \mathcal{L}_3 \subsetneq \mathcal{L}_1$ . Differently,  $\mathcal{L}_2$  and  $\mathcal{L}_3$

are minimal, since by removing any of their locations we may either break backward-closeness or alter the query result (by turning the value of the conditional expressions to empty). This shows that a minimal projection is not unique.

The goal of the chain analysis is to capture all nodes needed by all minimal projections of the query wrt all documents valid wrt a schema. To this end, we define how to connect chains with projections. Given a set of chains  $\tau$ , the set  $\mathcal{L}_\tau^t$  of locations in  $t$  typed by chains in  $\tau$  is defined as follows.

$$\mathcal{L}_\tau^t \stackrel{\text{def}}{=} \{ l \in \text{dom}(t) \mid c_l^t.c \in \tau \}$$

Dynamic and static variable environments have to link data and types in a consistent way, in order to ensure correctness of chain inference. This is stated next.

**Definition 4.3.3** (Consistent environments). *Let  $\mathbf{d}$  be a DTD,  $t = (\sigma, l_t) \in \mathbf{d}$  a valid tree,  $\gamma$  a dynamic environment binding variables to locations of  $\sigma$  and  $\Gamma$  a static environment binding variables to chains of  $\mathbf{C}_\mathbf{d}$ . The environments  $\gamma$  and  $\Gamma$  are said to be consistent if  $l \in \gamma(\mathbf{x})$  implies  $c_l^\sigma \in \Gamma(\mathbf{x})$ , for all  $l \in \text{dom}(\sigma)$ . We denote this by  $\sigma \models_\mathbf{d} \gamma : \Gamma$ .*

The following theorem formally states that chains inferred for a query  $\mathbf{q}$  cover the structure of data relevant for the query, and newly constructed elements.

**Theorem 4.3.4** (Soundness of Query Chain Inference). *Let  $\mathbf{d}$  be a DTD,  $t = (\sigma, l_t) \in \mathbf{d}$  a valid tree,  $\mathbf{q}$  a query,  $\gamma$  a dynamic environment and  $\Gamma$  a static environment such that  $\sigma \models_\mathbf{d} \gamma : \Gamma$ . Provided  $\sigma, \gamma \models \mathbf{q} \Rightarrow \sigma_{\text{new}}, L$  and  $\Gamma \vdash_{\mathbf{C}_\mathbf{d}} \mathbf{q} : (\mathbf{r}, \mathbf{v}, \mathbf{e})$  the following hold.*

1. *if  $t^{\text{min}}$  is a minimal  $\mathbf{q}$ -projection of  $t$  then  $t^{\text{min}} \leq t|_{\mathcal{L}_{\mathbf{r}, \mathbf{v}}^t} \leq t$*
2. *if  $t^{\text{new}} = (\sigma_{\text{new}}, l_{\text{new}})$  with  $l_{\text{new}} \in L \setminus \text{dom}(\sigma)$  is a fresh subtree then  $\text{dom}(t^{\text{new}}) = \mathcal{L}_{\mathbf{e}}^{t^{\text{new}}}$*

The first item of Theorem 4.3.4 states that chain inference is sound for used and return chains: a projection of any valid input tree made in terms of used and return chains includes every minimal  $\mathbf{q}$ -projection, hence preserves query semantics (the projection contains all the query needs for its evaluation). The second item is dedicated to element chain inference which is one of the key feature of our query-update analysis as already illustrated. Intuitively, this statement says that if element chains are used to project newly constructed elements (notice that  $l' \in L \setminus \text{dom}(\sigma)$ ) no node is pruned out, so element chains cover all possible chains in new elements of the query result. We conclude that our type inference system gives, as a side effect, a type-projector for the query [BCCN06].

### 4.3.2 Soundness of Update Chain Inference

Proving update chain soundness consists of establishing a link between  $i$ ) nodes in the stores  $t$  and  $\mathbf{u}(t)$  that are *involved* in the changes (deletion, insertion, renaming, and replacements) made by  $\mathbf{u}$  and  $ii$ ) nodes in these trees which are typed by the chains statically inferred for  $\mathbf{u}$ .

We stress that the evaluation of an update  $\mathbf{u}$  on a tree  $t = (\sigma, l_t) \in \mathbf{d}$  is a two phase process, that first generates the PUL  $\omega$  together with an extended  $\sigma_\omega$  containing the arguments of atomic insert and replace updates, for which the following relations hold.

$$\sigma, \gamma \models \mathbf{u} \Rightarrow \sigma_\omega, \omega \qquad \sigma_\omega \models \omega \rightsquigarrow \sigma_{\mathbf{u}}$$

Recall that  $\sigma_\omega$  includes  $\sigma$ , and that  $\mathbf{u}(t)$  denotes the tree  $(\sigma_{\mathbf{u}}@l_t, l_t)$ .

**Definition 4.3.5** (Involved Location). *We say that the update  $\mathbf{u}$  involves the location  $l \in \text{dom}(\sigma_\omega)$  if one of the following hold.*

- $l$  is the target location of an elementary delete command in  $\omega$
- $l$  is the target location of an elementary rename command in  $\omega$
- $l$  is a location, or the descendant of a location, in the source list  $L$  of a command  $\text{ins}(L, \_, \_)$  in  $\omega$
- $l$  is the target location of an elementary replace command in  $\omega$
- $l$  is a location, or the descendant of a location, in the source list  $L$  of a command  $\text{repl}(\_, L)$  in  $\omega$

The static analysis for update expressions is correct if all involved locations are typed by some update chains. In the following, we present examples of involved locations and how they are typed by the analysis.

The following examples are made over the tree  $t_2$  in Figure 4.1, which is valid against the DTD  $\mathbf{d}_4$  in Figure 4.2. As before, we assume that variable `doc` is bound to location  $l_{t_2}$  during dynamic evaluation and bound to type `bib` during static analysis.

Consider the following sequence of delete and rename updates

`delete doc//title, rename doc//publisher as “author”`

The update produces a pending list composed of two atomic update commands

`ren( $l_4$ , “author”), del( $l_2$ )`

meaning that the label of  $l_4$  (the publisher node) has to be renamed as “author” and that the location  $l_2$  has to be deleted. In this case both  $l_2$  and  $l_4$  are involved locations, since

they are target locations of elementary update commands. Note first that an involved location may belong to the initial tree  $t$  but not to the updated tree  $u(t)$ , for instance  $l_2$ , and a location may also belong to both trees like for instance  $l_4$ .

The static analysis for the delete updates over the DTD  $d_4$  produces the chain  $c = \text{bib.book} : \text{title}$ , that is the chain of the deleted location  $l_2$ , in fact  $c_{l_2}^{t_2} = c$ . Concerning the rename expression, two chains are inferred  $c_1 = \text{bib.book} : \text{publisher}$  and  $c_2 = \text{bib.book} : \text{author}$ . These two chains type the renamed node in the original and in the updated tree respectively, since

$$c_{l_4}^{t_2} = c_1 \quad c_{l_4}^{u(t_2)} = c_2$$

Involved locations for insert expressions are illustrated by the following example. Consider the tree  $t_2$  and the update

**insert** *<author/>* **into** *doc/self::bib/book*

In this case, the first phase of the evaluation produces a PUL composed of the command **ins**( $l'$ , **into**,  $l_1$ ) together with the store  $\sigma$  extended with a fresh element node labeled as *author*, which is referenced by a fresh location  $l'$ . Then, applying the PUL attaches location  $l'$  as children of  $l_1$  (which refers to the book node). The critical location for the PUL is  $l'$ . The static analysis for the update over the DTD  $d_4$  returns the chain  $c = \text{bib.book} : \text{author}$  and this is correct since  $c_{l'}^{u(t_2)} = c$ .

Involved locations for replace expressions are illustrated by the following example. Consider tree  $t_2$  and the update

**replace** *doc//author/name* **with** *doc/publisher/name*

here the update first build a fresh tree referenced by  $l_c$ , which is a copy of the one rooted at  $l'_3$ . Then  $l_c$  replaces  $l'_4$  (referencing to the publisher name). This results in the atomic update command

**repl**( $l'_4, l_c$ )

The critical locations are  $l'_4, l_c$  and all descendants of  $l_c$ .  $l'_4$  is a critical location because it is the target of the replace command.  $l_c$  and all of its descendants are critical locations, because  $l_c$  is the source locations of the update. We stress on the fact that descendants of  $l_c$  are fresh locations as well.

The static analysis wrt the DTD  $d_4$  returns chain  $c = \text{bib.book.author} : \text{name}$  typing the replaced location  $l'_4$ . Chain inference outputs also  $c = \text{bib.book.publisher} : \text{name}$  and  $c = \text{bib.book.publisher} : \text{name.String}$  that are the chains for  $l_c$  and for the child of  $l_c$ .

At this point we can state the correctness of the chain analysis for update expressions. The theorem below states that all locations involved by the update  $u$  are typed by chains inferred from  $u$ .

**Theorem 4.3.6** (Soundness of Update Chain Inference). *Let  $\mathbf{d}$  be a DTD,  $t = (\sigma, l_t) \in \mathbf{d}$  a valid tree,  $\mathbf{u}$  an update,  $\gamma$  a dynamic environment and  $\Gamma$  a static environment such that  $\sigma \models_{\mathbf{d}} \gamma : \Gamma$ . Provided that*

$$\sigma, \gamma \models \mathbf{u} \Rightarrow \sigma_{\omega}, \omega \quad \sigma_{\omega} \models \omega \rightsquigarrow \sigma_{\mathbf{u}} \quad \Gamma \vdash_{\mathbf{C}_{\mathbf{d}}} \mathbf{u} : \mathbf{U}$$

*it holds that*

- *if  $l \in \text{dom}(\sigma)$  is a location in  $t$ , and the update  $\mathbf{u}$  involves  $l$  then there exists  $\mathbf{c} : \mathbf{c}' \in \mathbf{U}$  such that  $\mathbf{c}_l^{\sigma} = \mathbf{c}.\mathbf{c}'$ , where  $\mathbf{c}' \neq \epsilon$ .*
- *if  $l$  is a location in  $\mathbf{u}(t)$ , i.e.  $l \in \text{dom}(\sigma_{\mathbf{u}} @ l_t)$ , and the update  $\mathbf{u}$  involves  $l$  then there exists  $\mathbf{c} : \mathbf{c}' \in \mathbf{U}$  such that  $\mathbf{c}_l^{\sigma_{\mathbf{u}}} = \mathbf{c}.\mathbf{c}'$ , where  $\mathbf{c}' \neq \epsilon$ .*

In the above statement, in case a location  $l$  belongs both to  $t$  and  $\mathbf{u}(t)$ , it may be that the chain typing  $l$  in  $t$  is different from the chain typing  $l$  in  $\mathbf{u}(t)$  (e.g., due to renaming). Although we made the assumption that the update expression is preserving the schema, it is worth noticing that Theorem 4.3.6 holds also for updates violating schema constraints ( $\mathbf{u}(t) \notin \mathbf{d}$ ), since chains corresponding to deleted or inserted nodes are always traced by the system regardless of correctness wrt the schema.

In this section we showed the correctness of chain inference for query and updates. This will be used in the remaining of the chapter to define a notion of independence, and show its correctness.

## 4.4 Independence Analysis

In this section we define a notion of query update independence based on chain overlapping. We show that this static notion of independence is correct and implies independence on all schema instances. We show that, for recursive schemas, the notion of independence based on chains does not have a trivial terminating algorithm. To overcome this issue, we propose a finite analysis which is complete wrt the infinite analysis and terminates in all cases.

### 4.4.1 Infinite Analysis

The notion of query-update independence  $\mathbf{q} \perp_{\mathbf{d}} \mathbf{u}$  (Definition 2.5.2) is based on the semantics of  $\mathbf{q}$  and  $\mathbf{u}$ , and involves all possible  $\mathbf{d}$  instances. The static counterpart of this notion is now proposed and is of course based on query and update chain inference. As chain inference depends on a set  $\mathbf{C}$  of chains, we first introduce a general static notion of  $\mathbf{C}$ -independence. Given two sets of chains  $\tau_1$  and  $\tau_2$ , the set of conflicting pairs of chains for  $\tau_1, \tau_2$  is defined as follows.

$$\text{confl}(\tau_1, \tau_2) \stackrel{\text{def}}{=} \{ (\mathbf{c}_1, \mathbf{c}_2) \mid \mathbf{c}_1 \in \tau_1, \mathbf{c}_2 \in \tau_2, \mathbf{c}_1 \leq \mathbf{c}_2 \}$$



**Definition 4.4.1** (C-independence). *Let  $d$  be a DTD,  $q$  a query and  $u$  an update, with at most only one free variable  $x$ , and  $\Gamma = \{ x \mapsto s_d \}$  a static environment. The query  $q$  and the update  $u$  are C-independent, denoted by  $q \perp_C u$ , if provided that  $\Gamma \vdash_C q : (r, v, e)$  and  $\Gamma \vdash_C u : U$  we have*

$$\text{confl}(r, U) = \emptyset \quad \text{confl}(U, r) = \emptyset \quad \text{confl}(U, v) = \emptyset$$

$C_d$ -independence can be decided in a finite amount of time when chain inference produces a finite set of chains i.e. when  $d$  has no vertical recursion. The next section provides a technique to capture  $C_d$ -independence by means of a finite analysis.

```

      bib ← book*
      book ← title?, author*, publisher*, price?
author, publisher ← name*
title, name, price ← String?

```

Figure 4.4: DTD  $d_5$

We illustrate the definition of C-independence throughout some examples. We consider DTD  $d_5$  of Figure 4.4. In the queries and updates below, we assume as before that variable `doc` is bound to `bib`. We check independence of query

```

q = for x in doc//book return
    if (x/price) then x/author else x/publisher

```

against the following updates

```

u1 = delete doc//title
u2 = rename doc//publisher as "author"
u3 = insert <price> 50 <price/> into doc//book
u4 = replace doc//price/text() with "10"

```

The static analysis infers, from the query,  $r_q = \{ \text{bib.book.author}, \text{bib.book.publisher} \}$  and  $v_q = \{ \text{bib.book}, \text{bib.book.price} \}$ . The following static independences hold.

$q$  is  $C_{d_5}$ -independent of  $u_1$ : the expressions are independent because the update does not delete any node that the query accesses or outputs. The static analysis infers

from the update  $U = \{ \text{bib.book} : \text{title} \}$ . Indeed, there is no pair of conflicting query and update chains.

query		update
$\text{bib.book.price}$	$\not\leq$	$\text{bib.book} : \text{title}$
$\text{bib.book.author}$	$\not\leq$	$\text{bib.book} : \text{title}$
$\text{bib.book.publisher}$	$\not\leq$	$\text{bib.book} : \text{title}$

Hence it results that  $\text{confl}(r, U) = \text{confl}(U, r) = \text{confl}(U, v) = \emptyset$ .

$q$  is not  $C_{d_5}$ -independent of  $u_2$ : the expressions are conflicting because the update may change the result of the query. The static analysis infers from the update  $U = \{ \text{bib.book} : \text{publisher}, \text{bib.book} : \text{author} \}$ . Indeed, query return-chains conflict with update chains.

query		update
$\text{bib.book.publisher}$	$\leq$	$\text{bib.book} : \text{publisher}$
$\text{bib.book.publisher}$	$\geq$	$\text{bib.book} : \text{publisher}$
$\text{bib.book.author}$	$\leq$	$\text{bib.book} : \text{author}$
$\text{bib.book.author}$	$\geq$	$\text{bib.book} : \text{author}$

Hence it results that  $\text{confl}(r, U) = \text{confl}(U, r) \neq \emptyset$

$q$  is not  $C_{d_5}$ -independent of  $u_3$ : the expressions may conflict because the data inserted by the update are accessed by the conditional query. The static analysis infers from the update  $U = \{ \text{bib.book} : \text{price}, \text{bib.book} : \text{price.String} \}$ . Indeed, a used chain conflicts with an update chain.

update		query
$\text{bib.book} : \text{price.String}$	$\not\leq$	$\text{bib.book.price}$
$\text{bib.book} : \text{price}$	$\leq$	$\text{bib.book.price}$

Hence it results that  $\text{confl}(U, v) \neq \emptyset$ .

$q$  is  $C_{d_5}$ -independent of  $u_4$ : the expressions are independent because the update does not affect the result of the query. In fact, the update impacts the textual node below the price node which is accessed by the query, without changing the boolean value of the conditional expression.

update		query
$\text{bib.book.price} : \text{String}$	$\not\leq$	$\text{bib.book.price}$

Hence it results that  $\text{confl}(r, U) = \text{confl}(U, r) = \text{confl}(U, v) = \emptyset$ .

We conclude this section by formally stating the soundness of the independence analysis based on chains. When  $C$  is taken as the set  $C_d$  of chains generated for the DTD  $d$ ,  $C$ -independence implies  $\perp_d$  independence, as stated by the following theorem.

**Theorem 4.4.2** (Soundness of  $C_d$  independence). *Let  $d$  be a DTD,  $q$  a query and  $u$  an update, with only one free variable. Then,*

$$q \perp_{C_d} u \text{ implies } q \perp_d u$$

As already stated, updates are assumed to preserve the schema. The above theorem needs this assumption in order to correctly use *query* chains in the independence analysis. Actually, if deletions violate the schema (a mandatory node is deleted), the notion of static independence  $\perp_{C_d}$  is still sound. The problem comes from insertions creating *new* chains (not belonging to  $C_d$ ) because they are not considered during chain inference for queries. As a consequence, the analysis made to check  $\perp_{C_d}$  could miss conflicting chains.

In this section we defined a notion of static independence based on chains which is sound. In the next section we will show how independence can be verified when recursive query and schemas may entail the inference of infinite sets of chains.

#### 4.4.2 Finite Analysis

The notion of  $C_d$ -independence (Definition 4.4.1) cannot be directly used to define a terminating decision algorithm, because for DTDs with vertical recursion the sets of inferred chains can be infinite. In this section we show how to finitely approximate sets of inferred chains so that  $C_d$ -independence can be detected in finite time.

One feature of chains generated by a recursive DTD  $d$  is that some of them contain multiple occurrences of (recursively defined) tags. So one way to characterize a finite set of  $d$ -chains is to restrict to chains having at most  $k$  occurrences of each tag. Hereafter, these chains are called  $k$ -chains, and for any set of chains  $\tau$ , its subset of  $k$ -chains is denoted by  $\tau^k$ . Thus,  $C_d^k$  denotes the set of  $k$ -chains generated by  $d$ .

As illustrated next, a *multiplicity value*  $k$  can be inferred from the query  $q$  and update  $u$ , so that independence according to chains in  $C_d$  is equivalent to independence according to inferred chains in  $C_d^k$ . The value  $k$  is derived by a two-steps static analysis.

Given an expression  $e$ , being either a query  $q$  or an update  $u$ , the first step associates a value  $k_e$  to  $e$  such that the set of  $k_e$ -chains inferred for  $e$  is *representative* of all possible inferred chains for  $e$ . Intuitively, the representative set of inferred chains for an expression synthesizes all possible inferred chains: any possible inferred chain can be mapped to a chain in the representative set by some folding transformations, according to recursive definitions in the DTD.

The second step infers a value  $k$  from the values  $k_q$  and  $k_u$ , such that the search of conflicting chains decisive for statically detecting  $q$ - $u$  independence can be safely done in the finite set of inferred  $k$ -chains.

Inferring the values  $k_q$  and  $k_u$  mainly depends on navigational properties of the XPath expressions occurring in the query and update. Thus, we start the discussion by focusing on XPath expressions  $p$ , and then consider FLWR expressions.

In the discussion, most of the examples are made on the following recursive DTD  $d_6$ . As before, variable `doc` in the expressions is meant to be bound to the root of  $d_6$ , that is `store`.

```

store   ← component*
component ← description, complist?, uselist?
uselist ← component+
complist ← component+
description ← String

```

The inference of the  $k$  value depends on the navigational axes employed in the XPath steps.

### Dealing with child, self and parent

The case where paths employ only the child, self or parent axis is a special one. In fact, such paths always generate a finite number of chains, hence a finite number of representative chains, and no folding is needed. What remains to do is to find a  $k$  value that fits all representative chains. The particularity of these paths is that they explicitly indicate how many times a single tag has to be matched: a good choice for  $k_p$  is the maximal tag frequency in the path  $p$ .

To illustrate, consider the following downward path  $p$

`doc/self :: store/component/complist/component/description`

there are four tags in  $p$ , all but `component` have frequency 1, while `component` has frequency 2. Therefore, the maximal tag frequency is 2, and indeed 2-chains include the representative chain

`store.component.complist.component.description`

which is the only chain inferred for this path (for clarity, we over line the tag with  $k$ -frequency).

The same holds for the navigational path with backward axis

`doc/self :: store/component/complist/component/parent :: complist`

Here, the chains inferred from the path are 2-chains

```

store.component.complist.component
store.component.complist.component.complist

```

Similarly, for the path

$\text{doc/self} :: \text{store/component/complist}/*$

we infer  $k_p = 2$ , since the wildcard  $*$  stands for any label.

It is worth noticing that, for the analysis to be correct, the query/update free variables have to be always bound to the schema root type. If we relax this assumption, the chain analysis may turn out to be unsound for recursive schemas. To see this, consider for instance the query

$\text{doc/component}$

with  $\text{doc}$  initially bound to the chain  $\text{store.component.description}$ . The chain inferred for the query is  $\text{store.component.description.component}$ , which is a 2-chain. However, since for this expression  $k = 1$ , the analysis infers the empty set of used, return and element chains for the query, thus deeming it as independent of any update, while it should not.

### Dealing with descendant and ancestor

When a path  $p$  makes use of the **descendant** axis, the number of chains inferred for an expression may not be finite, and the length of inferred chains totally unrelated to the length of  $p$ . For instance, the path

$\text{doc/descendant} :: \text{description}$

makes the system inferring an infinite number of chains, of unbounded length.

This is what led us to reason in terms of tag frequency rather than path length. To generate a finite set of chains, the value  $k_p$  is determined by taking into account the number of occurrences of descendant axes in  $p$ .

To illustrate this, we still consider the schema  $d_6$ , and observe that the type **component** is defined in terms of **description**, **uselist** and **complist**, and the two latter are defined in terms of **component**. These mutual recursive definitions entail that, in a valid document instance, a node with type **uselist** can be a descendant of a node with type **complist**, and vice versa, along the same chain of the tree. In addition, a chain connecting **uselist** and **complist** nodes always contains an intermediate **component** label, which also occurs before the first occurrence of a **description**, **uselist**, and **complist** label.

As a consequence, for the following path

$p = \text{doc/descendant} :: \text{complist/descendant} :: \text{uselist/descendant} :: \text{description}$

the shortest chain that is inferred is the 3-chain

$\text{store.component.complist.component.uselist.component.description}$

Simple tag frequency, like for the previous cases, would lead to  $k_p = 1$ . This is not satisfactory because no 1-chain is inferred for  $p$ . To reflect the fact that each recursive axis may permit any tag to repeat once in inferred chains, the correct maximal tag frequency we have to consider for the path  $p$  is 3; in fact, 3-chains do allow to infer a non-empty set of representative chains - other inferred chains can be obtained by unfolding recursion.

Of course, an XPath expression may combine both recursive and non-recursive axes. In this case, for a path  $p$  we obtain  $k_p$  as the sum of two components computed independently: the maximal tag frequency for non-recursive steps, and the number of recursive steps in  $p$ . As an example, for

$$p = \text{doc}/\text{descendant} :: \text{component}/\text{uselist}/\text{component}$$

we have  $k_p = 2$  since the maximal tag frequency for the path suffix  $/\text{uselist}/\text{component}$  is 1, and there is 1 descendant step  $/\text{descendant} :: \text{component}$ .

Recursive backward axes are handled similarly. Here we have to pay attention to the fact that chains navigated by an **ancestor** step are prefixes of some chains generated by previous steps. Consider

$$p = \text{doc}/\text{descendant} :: \text{uselist}/\text{ancestor} :: \text{complist}$$

Here  $k_p$  has to be such that the *used* chain

$$\text{store}.\overline{\text{component}}.\text{complist}.\overline{\text{component}}.\text{uselist}$$

can be generated. Thus, we enforce **ancestor** steps to increment the tag frequency by 1. This is reminiscent of what we have seen before in the case of **descendant**; the way  $p$  is processed can be compared to the way would be processed the navigational path

$$\text{doc}/\text{descendant} :: \text{complist}/\text{descendant} :: \text{uselist}$$

and therefore chains containing **complist** as ancestor of **uselist** need to be generated for  $p$ .

Concerning paths  $p$  employing either **descendant-or-self** or **ancestor-or-self**,  $k_p$  is computed as for the self-less axes. The same holds for shortcut expressions like  $p = //\phi$ , that stands for  $p' = /\text{descendant-or-self} :: */\phi$ , when  $\phi \neq *$ . It follows that  $k_p = k_{p'}$  simply because the sets of chains inferred for the two expressions are exactly the same.

It is worth noticing that by assuming that the root-type of a DTD is non-recursive we can provide a uniform the definition of the  $k$  value, and simplify proofs. To illustrate, consider the DTD  $d = \{ a \leftarrow a? \}$ , the query  $q = x/\text{child} :: a$ , and the static environment  $\Gamma = \{ x \mapsto a \}$ . Our system infers from  $q$  the 2-chain  $a.a$ . However, since  $k_q = \mathcal{F}(a, x/\text{child} :: a) = 1$  the chain  $a.a$  is ruled-out from inference, while it should not. The same happens if  $q = x/\text{descendant} :: a$ . This of course can be captured by setting

$\mathcal{F}(a, \mathbf{x}/\mathbf{axis} :: \phi) = 2$  if *i*)  $\mathbf{x}$  is bound to the singleton root-type, *ii*) the root-type is recursive, *iii*) the root-type satisfies  $\phi$ , and of course *iv*)  $\mathbf{axis}$  is downward. This would make the formal definition for the base case quite inelegant, while assuming a non-recursive root-type we have an uniform formalization of the tag frequency. This discrepancy is also caused by the fact that, as discussed in Chapter 2, we do not have the document-node in our data-model.

### Dealing with sibling axes

Sibling axes are managed as child and parent axes. Consider the following recursive DTD  $\mathbf{d}_7$

```

store    ← component*
component ← description, component?
description ← String

```

and the navigational path

$\mathbf{doc}/\mathbf{descendant} :: \mathbf{description}/\mathbf{following-sibling} :: \mathbf{component}$

For this path, the used 1-chain  $\mathbf{store.component.description}$  and the return 2-chain

$\overline{\mathbf{store.component.component}}$

are the needed chains. The presence of the recursive step  $\mathbf{/descendant} :: \mathbf{description}$  entails  $k = 2$ .

### Dealing with FLWR expressions and Updates

Based on concepts previously illustrated, we provide now formal definitions to deal with the general case of a FLWR expression  $\mathbf{e}$ .

As seen before, the computation of  $k_{\mathbf{e}}$  is decomposed into two tasks. The first one determines via the function  $\mathcal{F}(a, \mathbf{e})$  the frequency of each tag  $a \in \Sigma$  on the whole expression, in order to derive the maximal frequency. The second task computes via the function  $\mathcal{R}(\mathbf{e})$  the maximal number of recursive steps used by a navigational path embodied in the whole expression. The value  $k_{\mathbf{e}}$  is the sum of these two values. Formally:

$$k_{\mathbf{e}} \stackrel{def}{=} \max\{ \mathcal{F}(a, \mathbf{e}) \mid a \in \Sigma \} + \mathcal{R}(\mathbf{e})$$

The functions  $\mathcal{F}(a, \mathbf{e})$  and  $\mathcal{R}(\mathbf{e})$  are defined by structural induction in Table 4.3 and 4.4, respectively.

When  $\mathbf{e}$  is a for or a let expression, the value  $k_{\mathbf{e}}$  is specified by summing the sub-expression values. This is captured by rules ( $\mathcal{F}$ -FOR/LET) and ( $\mathcal{R}$ -FOR/LET). This

$\frac{}{\mathcal{F}(a, ()) = 0} (\mathcal{F}\text{-EMPTY})$	$\frac{}{\mathcal{F}(a, \text{"txt"}) = 0} (\mathcal{F}\text{-TEXT})$	$\frac{f = \mathcal{F}(a, \mathbf{q})}{\mathcal{F}(a, \langle a \rangle \mathbf{q} \langle /a \rangle) = f + 1} (\mathcal{F}\text{-ELEMENT})$
$\frac{\text{axis is recursive}}{\mathcal{F}(a, \text{axis} :: \phi) = 0} (\mathcal{F}\text{-STEP1})$	$\frac{\phi \notin \{a, \text{node}()\}}{\mathcal{F}(a, \text{axis} :: \phi) = 0} (\mathcal{F}\text{-STEP2})$	$\frac{\begin{array}{c} \text{axis is not recursive} \\ \phi \in \{a, \text{node}()\} \end{array}}{\mathcal{F}(a, \text{axis} :: \phi) = 1} (\mathcal{F}\text{-STEP3})$
$\frac{f_1 = \mathcal{F}(a, \mathbf{q}_1) \quad f_2 = \mathcal{F}(a, \mathbf{q}_2)}{\mathcal{F}(a, \mathbf{q}_1, \mathbf{q}_2) = \max(f_1, f_2)} (\mathcal{F}\text{-CONCAT})$	$\frac{f_0 = \mathcal{F}(a, \mathbf{q}_0) \quad f_1 = \mathcal{F}(a, \mathbf{q}_1) \quad f_2 = \mathcal{F}(a, \mathbf{q}_2)}{\mathcal{F}(a, \text{if } (\mathbf{q}_0) \text{ then } \mathbf{q}_1 \text{ else } \mathbf{q}_2) = \max(f_0, f_1, f_2)} (\mathcal{F}\text{-IF})$	
$\frac{f_1 = \mathcal{F}(a, \mathbf{q}_1) \quad f_2 = \mathcal{F}(a, \mathbf{q}_2)}{\mathcal{F}(a, \text{for } \mathbf{x} \text{ in } \mathbf{q}_1 \text{ return } \mathbf{q}_2) = f_1 + f_2} (\mathcal{F}\text{-FOR})$	$\frac{f_1 = \mathcal{F}(a, \mathbf{q}_1) \quad f_2 = \mathcal{F}(a, \mathbf{q}_2)}{\mathcal{F}(a, \text{let } \mathbf{x} := \mathbf{q}_1 \text{ return } \mathbf{q}_2) = f_1 + f_2} (\mathcal{F}\text{-LET})$	
$\frac{f = \mathcal{F}(a, \mathbf{q})}{\mathcal{F}(a, \text{delete } \mathbf{q}) = f} (\mathcal{F}\text{-DELETE})$		
$\frac{f = \mathcal{F}(a, \mathbf{q})}{\mathcal{F}(a, \text{rename } \mathbf{q} \text{ as } a) = 1 + f} (\mathcal{F}\text{-RENAME1})$	$\frac{f = \mathcal{F}(a, \mathbf{q}) \quad a \neq b}{\mathcal{F}(a, \text{rename } \mathbf{q} \text{ as } b) = f} (\mathcal{F}\text{-RENAME2})$	
$\frac{f = \mathcal{F}(a, \mathbf{q}) \quad f_0 = \mathcal{F}(a, \mathbf{q}_0)}{\mathcal{F}(a, \text{insert } \mathbf{q} \text{ pos } \mathbf{q}_0) = f + f_0} (\mathcal{F}\text{-INSERT})$	$\frac{f_0 = \mathcal{F}(a, \mathbf{q}_0) \quad f = \mathcal{F}(a, \mathbf{q})}{\mathcal{F}(a, \text{replace } \mathbf{q}_0 \text{ with pos } \mathbf{q}) = f_0 + f} (\mathcal{F}\text{-REPLACE})$	
$\frac{f_1 = \mathcal{F}(a, \mathbf{u}_1) \quad f_2 = \mathcal{F}(a, \mathbf{u}_2)}{\mathcal{F}(a, \mathbf{u}_1, \mathbf{u}_2) = \max(f_1, f_2)} (\mathcal{F}\text{-UCONCAT})$	$\frac{f_0 = \mathcal{F}(a, \mathbf{q}_0) \quad f_1 = \mathcal{F}(a, \mathbf{u}_1) \quad f_2 = \mathcal{F}(a, \mathbf{u}_2)}{\mathcal{F}(a, \text{if } (\mathbf{q}_0) \text{ then } \mathbf{u}_1 \text{ else } \mathbf{u}_2) = \max(f_0, f_1, f_2)} (\mathcal{F}\text{-UIF})$	
$\frac{f = \mathcal{F}(a, \mathbf{q}) \quad f' = \mathcal{F}(a, \mathbf{u})}{\mathcal{F}(a, \text{for } \mathbf{x} \text{ in } \mathbf{q} \text{ return } \mathbf{u}) = f + f'} (\mathcal{F}\text{-UFOR})$	$\frac{f = \mathcal{F}(a, \mathbf{q}) \quad f' = \mathcal{F}(a, \mathbf{u})}{\mathcal{F}(a, \text{let } \mathbf{x} := \mathbf{q} \text{ return } \mathbf{u}) = f + f'} (\mathcal{F}\text{-ULET})$	

Table 4.3:  $\mathcal{F}$  definition.



$\frac{}{\mathcal{R}(() = 0)} (\mathcal{R}\text{-EMPTY})$	$\frac{}{\mathcal{R}(\text{"txt"}) = 0} (\mathcal{R}\text{-TEXT})$	$\frac{r = \mathcal{R}(q)}{\mathcal{R}(<a>q</a>) = r} (\mathcal{R}\text{-ELEMENT})$
$\frac{\text{axis is recursive}}{\mathcal{R}(\text{axis} :: \phi) = 1} (\mathcal{R}\text{-STEP1})$	$\frac{\text{axis is not recursive}}{\mathcal{R}(\text{axis} :: \phi) = 0} (\mathcal{R}\text{-STEP2})$	
$\frac{r_1 = \mathcal{R}(q_1) \quad r_2 = \mathcal{R}(q_2)}{\mathcal{R}(q_1, q_2) = \max(r_1, r_2)} (\mathcal{R}\text{-CONCAT})$	$\frac{r_0 = \mathcal{R}(q_0) \quad r_1 = \mathcal{R}(q_1) \quad r_2 = \mathcal{R}(q_2)}{\mathcal{R}(\text{if } (q_0) \text{ then } q_1 \text{ else } q_2) = \max(r_0, r_1, r_2)} (\mathcal{R}\text{-IF})$	
$\frac{r_1 = \mathcal{R}(q_1) \quad r_2 = \mathcal{R}(q_2)}{\mathcal{R}(\text{for } x \text{ in } q_1 \text{ return } q_2) = r_1 + r_2} (\mathcal{R}\text{-FOR})$	$\frac{r_1 = \mathcal{R}(q_1) \quad r_2 = \mathcal{R}(q_2)}{\mathcal{R}(\text{let } x := q_1 \text{ return } q_2) = r_1 + r_2} (\mathcal{R}\text{-LET})$	
$\frac{r = \mathcal{R}(q)}{\mathcal{R}(\text{delete } q) = r} (\mathcal{R}\text{-DELETE})$	$\frac{r = \mathcal{R}(q)}{\mathcal{R}(\text{rename } q \text{ as } b) = r} (\mathcal{R}\text{-RENAME})$	
$\frac{r = \mathcal{R}(q) \quad r_0 = \mathcal{R}(q_0)}{\mathcal{R}(\text{insert } q \text{ pos } q_0) = r + r_0} (\mathcal{R}\text{-INSERT})$	$\frac{r_0 = \mathcal{R}(q_0) \quad r = \mathcal{R}(q)}{\mathcal{R}(\text{replace } q_0 \text{ with pos } q) = r_0 + r} (\mathcal{R}\text{-REPLACE})$	
$\frac{r_1 = \mathcal{R}(u_1) \quad r_2 = \mathcal{R}(u_2)}{\mathcal{R}(u_1, u_2) = \max(r_1, r_2)} (\mathcal{R}\text{-UConcat})$	$\frac{r_0 = \mathcal{R}(q_0) \quad r_1 = \mathcal{R}(u_1) \quad r_2 = \mathcal{R}(u_2)}{\mathcal{R}(\text{if } (q_0) \text{ then } u_1 \text{ else } u_2) = \max(r_0, r_1, r_2)} (\mathcal{R}\text{-UIF})$	
$\frac{r = \mathcal{R}(q) \quad r' = \mathcal{R}(u)}{\mathcal{R}(\text{for } x \text{ in } q \text{ return } u) = r + r'} (\mathcal{R}\text{-UFor})$	$\frac{r = \mathcal{R}(q) \quad r' = \mathcal{R}(u)}{\mathcal{R}(\text{let } x := q \text{ return } u) = r + r'} (\mathcal{R}\text{-ULet})$	

Table 4.4:  $\mathcal{R}$  definition.

is motivated by the fact that, for instance, for-expressions are usually used to encode nested iterations performed by XPath paths, like in the query

```
for x in doc/a for y in x/b return y
```

This definition leads in some cases to an overestimation of the value  $k_e$  that would be actually sufficient for a finite analysis. For instance, for the query

```
q = for x in doc/a/a for y in doc/a/b return x,y
```

we have  $\mathcal{F}(a, q) = 3$ , while the value 2 would be sufficient. It is worth saying that this overestimation raises the  $k_q$  value. More precision can be obtained by tracing variable bindings in the definition of  $\mathcal{F}$ . The same argument holds for  $\mathcal{R}$ . However, this would make the formalization cumbersome without being a decisive factor for the analysis. Thus, our choice has been guided by simplicity and conciseness of  $\mathcal{F}$  and  $\mathcal{R}$  definitions.

The rules for element construction, ( $\mathcal{F}$ -ELEMENT) and ( $\mathcal{R}$ -ELEMENT), deserve some comment. Note that tags of constructed elements are taken into account. Indeed, these elements can be inserted by an update as children of existing elements, thus generating new chains pointing nodes that can be accessed by a query. Consider the recursive schema  $d_7$  and the following update  $u$ .

```
for x in doc/self::store/component return
  insert
    <component>
      <component>
        <description/>
      </component>
    </component>
into x
```

As already outlined, precision of the independence analysis relies (among other things) on the chains generated for element construction. The rules in Table 4.3 lead to infer  $k_u = \mathcal{F}(u, component) = 3$ , and thus the chain

```
store.component : component.component.description
```

is inferred for the finite analysis. The inferred value is the same, if  $u$  is a replace update instead of an insert update.

Tag frequency for delete update expression is determined as for queries. For rename expressions tag frequency is determined by observing that for recursive schemas, after renaming, the frequency of the renaming label may increase of 1. To illustrate, consider the DTD  $d_7$  update

```
u = rename doc//description as "component"
```

In this case, if we settle  $\mathcal{F}(\mathbf{u}, \text{component}) = 1$  to capture the fact that a new node labeled with *component* appears. Together with  $\mathcal{R}(/description) = 1$ , this makes  $k_{\mathbf{u}} = 2$ , and indeed this is the multiplicity value needed to infer the 2-chains

$$\text{store.component} : \text{description} \quad \text{and} \quad \text{store.component} : \text{component}$$

Notice that if the frequency of *component* is not kept into account then  $k_{\mathbf{u}} = 1$  and chain  $\text{store.component} : \text{component}$  is wrongly excluded from inference.

### Finite independence analysis

We see now how to use the values  $k_{\mathbf{q}}$  and  $k_{\mathbf{u}}$  in order to determine a  $k$  value such that  $\mathbf{C}_d$ -independence can be detected by restricting the analysis to  $k$ -chains.

Consider the following query and update pair over the DTD  $\mathbf{d}_6$ .

$$\mathbf{q} = /descendant :: \text{complist} \quad \mathbf{u} = \text{delete } /descendant :: \text{uselist}$$

The two expressions are dependent because *uselist* is an ancestor of *complist* and they are both navigated by a recursive axis. It results that  $k_{\mathbf{q}} = 1$  and  $k_{\mathbf{u}} = 1$ . At this point, we could argue that a sound choice is

$$k = \max(k_{\mathbf{q}}, k_{\mathbf{u}})$$

which allows the finite analysis to infer the query chain

$$\text{store.component.complist}$$

and the update chain

$$\text{store.component} : \text{uselist}$$

Unfortunately, these chains do not conflict, and rule out dependence. The problem here comes from the fact that the update may change a descendant of a query returned node, and that  $k = \max(k_{\mathbf{q}}, k_{\mathbf{u}})$  does not permit to capture this in the finite analysis. It may also occur that a query navigates descendants of nodes inserted/created by the update, and again such a definition rules out independence. To avoid this problem, it is necessary that representative chains that are inferred for the update cover query returned nodes, and vice versa. To this end, while inferring chains for the update  $\mathbf{u}$ , structural properties of the query  $\mathbf{q}$  have also to be taken into account. This is obtained by setting

$$k = k_{\mathbf{q}} + k_{\mathbf{u}}$$

Now with  $k = 2$  it is possible to infer 2-chains witnessing the conflict such as

$$\text{store.component.complist} \leq \text{store.component.complist.component} : \text{uselist}$$

### 4.4.3 Soundness of the Finite Chain Analysis

In this section we state one of our main results, soundness of the finite analysis.

**Theorem 4.4.3** (Soundness of  $C_d^k$  Independence). *Let  $d$  be a DTD,  $q$  a query and  $u$  an update, with only one free variable. Let  $k = k_q + k_u$  as defined above, then*

$$q \perp_{C_d^k} u \quad \text{implies} \quad q \perp_d u$$

We develop the main steps of the proof of the Theorem. The thesis follows by showing that  $q \perp_{C_d} u$  is equivalent to  $q \perp_{C_d^k} u$ . Proving that  $q \perp_{C_d} u$  implies  $q \perp_{C_d^k} u$  is straightforward since  $C_d^k \subseteq C_d$ . Also, we reason in terms of *dependence*, rather than independence. We prove that  $C_d$ -dependence implies  $C_d^k$ -dependence (these notions directly follow from Definition 4.4.1) by showing that from any pair of chains in  $C_d$ , witness of dependence, it is possible to identify a pair of  $k$ -chains in  $C_d^k$ , witness of dependence.

The proof is composed of three steps. First, we show that there exists a folding from query chains to  $k$ -query-chains, for any query  $q$  (Lemma 4.4.4). Then, we show that there exists a folding from query chains to  $k$ -query-chains also preserving the prefix relation  $\preceq$ , for any pair of queries  $(q, q')$  (Lemma 4.4.5). Finally, we show that such a folding exists for chains inferred for any query-update pair  $(q, u)$  (Theorem 4.4.3). Proofs are reported in Chapter 9.

Given a DTD  $d$ , we define a folding relation  $\hookrightarrow_d \subseteq C_d \times C_d$  as

$$\hookrightarrow_d \stackrel{\text{def}}{=} \{ (c_1, c_2) \mid c_1 = c.a.c'.c'' \wedge c_2 = c.a.c'' \}$$

Notice that, above, the symbol  $a$  is a recursive type of the schema. We denote by  $\hookrightarrow_d^*$  the reflexive and transitive closure of  $\hookrightarrow_d$ .

We state now the correctness of the  $k$ -folding for a single query. In the statements below we consider that, if  $c$  is a used (respectively return or element) chain, then  $c'$  is a used (respectively return or element) chain.

**Lemma 4.4.4** (Folding). *Let  $d$  be a DTD,  $q$  a query with only one free variable  $x$ , and  $\Gamma = \{ x \mapsto s_d \}$  a static environment. Assume that  $\Gamma \vdash_{C_d} q : (r, v, e)$  and  $\tau = r \cup v \cup e$ . For each chain  $c \in \tau$  there exists a chain  $c' \in \tau$  such that  $c \hookrightarrow_d^* c'$  and  $c'$  is a  $k_q$ -chain.*

As an example, consider DTD  $d_6$  and the query  $q = //component$ . In this case  $k_q = 1$  and we have, for instance, the following 2-chain folding onto a 1-chain.

$$\text{store.component.uselist.component} \quad \hookrightarrow_d \quad \text{store.component}$$

When  $q$  and  $u$  are  $C_d$ -dependent, at least one of  $\text{confl}(U, v)$ ,  $\text{confl}(r, U)$ ,  $\text{confl}(U, r)$  is nonempty (see Definition 4.4.1). This implies that there exists a conflicting pair of inferred chains, witness of the  $C_d$ -dependence of  $q$  and  $u$ . As updates are defined *in terms of* queries, the next lemma which focuses on “conflicting” query chains is needed to conclude the proof of Theorem 4.4.3.

**Lemma 4.4.5** (Folding and Conflict Preservation). *Let  $d$  be a DTD,  $q_1, q_2$  two queries with only one free variable  $x$ , and  $\Gamma = \{ x \mapsto s_d \}$  a static environment. Assume that  $\Gamma \vdash_{C_d} q_i : (r_i, v_i, e_i)$  and  $\tau_i = r_i \cup v_i \cup e_i$ . For each pair of chains  $(c_1, c_2) \in \tau_1 \times \tau_2$  such that  $c_1 \leq c_2$ , there exists  $(c'_1, c'_2) \in \tau_1 \times \tau_2$  such that  $c'_1 \leq c'_2$  and  $c_i \mapsto_d^* c'_i$  with  $c'_i$  a  $(k_{q_1} + k_{q_2})$ -chain ( $i = 1, 2$ ).*

As an example, consider DTD  $d_6$  and the query  $q = //uselist$  and the update  $u = \text{delete } //complist$ . In this case  $k = k_q + k_u = 2$  and we have, for instance, the following conflicting chains folding onto  $k$ -chains, and preserving the conflict. Below, for the sake of concision,  $st = \text{store}$ ,  $co = \text{component}$ ,  $ul = \text{uselist}$ , and  $cl = \text{complist}$ .

$$\begin{array}{ccc} st.co.ul.co.ul.co.ul & \mapsto_d^* & st.co.ul \\ \leq & & \leq \\ st.co.ul.co.ul.co.ul.co.cl & \mapsto_d^* & st.co.ul.co.cl \end{array}$$

While, if we have that  $k < 2$  then the folding does not preserve the conflicts.

$$\begin{array}{ccc} st.co.ul.co.ul.co.ul & \mapsto_d^* & st.co.ul \\ \leq & & \not\leq \\ st.co.ul.co.ul.co.ul.co.cl & \mapsto_d^* & st.co.cl \end{array}$$

## Conclusions

In this chapter we presented two of the main contributions of this thesis: a type system for chain inference and a static notion of independence based on chains. We proved also that both the type system and the independence analysis are correct. Moreover, we illustrated how to make the static notion of independence verifiable in a finite time, even in the case where inference entails an infinite number of query and update chain. In the next chapter we will discuss algorithms and data structures for efficiently implementing the chain-based independence.

## Chapter 5

# Implementing the Chain-based Independence Analysis

In this chapter we present algorithms for efficiently implementing the chain-based independence analysis. Chain inference poses significant challenges from the computational point of view. Because of the expressivity of the schema and of the query/update languages considered, it is easy to have exponential blowups of chains during inference. In light of this, we present a data-structure for storing chains based on DAGs, and algorithms implementing the chain-based independence analysis, that allow to run the independence analysis in polynomial space and time, whilst still retaining an high precision.

### 5.1 Introduction

Chain inference is a challenging problem not only from a formal point of view, but also from a computational point of view. In fact, the number of distinct chains inferred for a single expression may be exponential in the size of the input DTD. This happens for DTDs that make heavy use of recursive definitions, but also for non-recursive ones, like for DTD  $d_8$ , that is defined by the following set of productions.

$$a \leftarrow b_1, f_1 \quad \dots \quad b_{(i-1)} \leftarrow (b_i, f_i)? \quad f_{(i-1)} \leftarrow (b_i, f_i)?$$

with  $1 \leq i \leq n$ .

The DTD  $d_8$  employs  $2n + 1$  types, and a simple navigational query like  $//b_n$  makes the system inferring  $\Omega(2^n)$  chains. Namely, all chains of the form

$$a \ x_1 \ x_2 \ \dots \ x_{(n-1)} \ b_n$$

with  $x_i \in \{b_i, f_i\}$ . This happens because types are heavily reused in the schema. If no type is reused, it turns out that, the DTD is non-recursive and the number of chains becomes linear in the size of the schema.

This already demonstrates the expressive power of the constraints that can be defined by a non-recursive DTD. The situation is even exacerbated by the use of recursion, where it is possible to reproduce the previous situation with a schema featuring just three types. Consider the DTD  $\mathbf{d}_9$  defined as follows

$$\begin{aligned} \mathbf{a} &\leftarrow \mathbf{b}, \mathbf{f} \\ \mathbf{b} &\leftarrow (\mathbf{b}, \mathbf{f})? \\ \mathbf{f} &\leftarrow (\mathbf{b}, \mathbf{f})? \end{aligned}$$

In this case, the chains of length  $h + 1$  inferred for a navigational step like  $//b$ , are again all chains of the form  $\mathbf{a} x_1 x_2 \dots x_{(h-1)} \mathbf{b}$  with  $x_i \in \{\mathbf{b}, \mathbf{f}\}$ . Again, we have an exponential blowup since the cardinality of the result is in  $\Omega(2^h)$ . As we show later (see Proposition 5.1.5), when the multiplicity value  $k$  is considered instead of the chain length  $h$ , the quantity of inferred chains is still exponential.

These examples outline the hardness of chain inference already in terms of cardinality of the results. If chains are not stored efficiently, this may have a direct consequence on space occupancy. This is problematic because if the space complexity for intermediate computations has an exponential upper-bound, then time complexity is at least of the same order, thus compromising tractability of the analysis.

In the remainder of the chapter we will see that the crux in implementing the chain analysis is the choice of the data-structure to use for efficiently storing and retrieving chains. Unsurprisingly, the blowups outlined in the previous example imply that storing chains requires exponential space for instance if we use a trivial data-structure such as lists, or a less trivial one, such as trees. In light of this, we propose a compressed data structure based on DAGs that permits to store chains in polynomial space, paving the way to the tractability of the analysis.

### Cardinality of Schema Chains

In this section, we investigate the cardinality of the set of chains induced by a schema, and relative issues concerning their space occupancy. We study both non-recursive and recursive DTDs, and storing chains on lists and trees. Our results show clearly that the static notion of independence we defined in Chapter 4 may require an exponential time to be checked.

Recall that given two functions  $f$  and  $g$ , we say that  $f$  is in  $O(g)$  or, equivalently, that  $g$  is in  $\Omega(f)$  if  $f \leq Cg + C'$ , for two positive constants  $C, C'$  and all values of the variables. We say that  $f$  is in  $\Theta(g)$  if  $f$  is in  $O(g)$  and  $f$  is in  $\Omega(g)$ .

The cases when our system produces the larger number of chains are those relative to saturated schema (defined in Chapter 2).

**Proposition 5.1.1.** *For a saturated non-recursive DTD  $\mathbf{d}$  with  $n$  types, the cardinality of  $\mathbf{C}_{\mathbf{d}}$  is  $\Theta(2^n)$ .*

*Proof.* Because  $\mathbf{d}$  is non-recursive and saturated,  $\Rightarrow_{\mathbf{d}}$  is a strict total order. It is total, because all pairs of types in  $\mathbf{d}$  are adjacent in  $\mathcal{G}_{\mathbf{d}}$ . It is strict, because  $\mathbf{d}$  is non-recursive, and thus  $\Rightarrow_{\mathbf{d}}$  is irreflexive. This implies that any *subset* of types in  $\mathbf{d}$  together with the order of  $\Rightarrow_{\mathbf{d}}$  corresponds to a chain in  $\mathbf{C}_{\mathbf{d}}$ . The size of  $\mathbf{C}_{\mathbf{d}}$  comes from the size of the powerset of types in  $\mathbf{d}$ . This is  $2^n - 1$ , without counting the empty chain. Moreover, the number of *rooted chains* induced by  $\mathbf{d}$  coincides with the number of subsets of types of  $\mathbf{d}$  that contain at least the root-type. These are exactly  $2^{n-1}$ .  $\square$

We compare the case of using *lists* and *prefix-trees* for storing chains.

**List** Lists are defined as sequences of objects. Storing chains on lists implies to store each time the full chain. For instance, chains `bib.book.author` and `store.component` are represented as `[bib; book; author]` and `[store; component]` respectively.

**Prefix-Trees** Prefix-trees are trees where all chain prefixes are stored only once. Each path in the tree is a chain. This optimizes chain inference wrt using lists. For instance, given two chains `bib.book.title` and `bib.book.author` the prefix-tree allow to store only the prefix `bib.book` on a single edge, even if it is shared by both chains.

**Proposition 5.1.2.** *The set of rooted chains induced by a saturated non-recursive DTD  $\mathbf{d}$  of  $n + 1$  types requires, to be stored on lists,  $\Theta(n2^n)$  space.*

*Proof.* The proof follows from some combinatorics. As discussed in the proof of Theorem 5.1.1,  $\mathbf{d}$  being saturated and non-recursive, the power set of types in  $\mathbf{d}$  is the equivalent of  $\mathbf{C}_{\mathbf{d}}$ . Since we count rooted chains, we consider all subsets of types containing the root, and hence we count  $n$  types instead of  $n + 1$ . We model space occupancy by making a weighted sum of all subsets of types of size  $1, 2, \dots, n$ . The number of sets of types of size  $i$  in  $\mathbf{d}$  can be counted by using the number of combinations of  $i$  elements chosen from a set of  $n$  elements  $\binom{n}{i}$ . Then we consider that, by using lists, a subset of size  $i$  would need space proportional to  $i + 1$  (by re-considering also the root). Therefore, list space occupation is

$$\sum_{i=0}^n \binom{n}{i} (i + 1) = \sum_{i=0}^n \binom{n}{i} i + \sum_{i=0}^n \binom{n}{i} = n2^{n-1} + 2^n = (n + 2)2^{n-1}$$

and thus we obtain  $\Theta(n2^n)$ .  $\square$

Next statement shows that, even if we use an optimized data-structure such as prefix-trees for storing chains, we still need exponential space for storing all schema chains.



**Proposition 5.1.3.** *The set of rooted chains induced by a saturated non-recursive DTD  $\mathbf{d}$  of  $n$  types, can be stored on a prefix tree in  $\Theta(2^n)$  space.*

*Proof.* The result follows from defining the size of a tree as a recurrence relation. We consider the worst case where there is a maximum type dependency for the non-recursive DTD  $\mathbf{d}$ . Let us denote by  $Size(i)$  the size of the prefix tree containing all rooted chains of a DTD with  $i$  types. For simplicity, we define the space occupancy of a node and of an edge as 1. The growth of the size of the tree is depicted in Figure 5.1. In general, for  $n > 0$ , we have the recurrence relation

$$Size(n) = 2 \cdot Size(n-1) + 1$$

from which we conclude that the size of the tree is in  $\Theta(2^n)$ .  $\square$

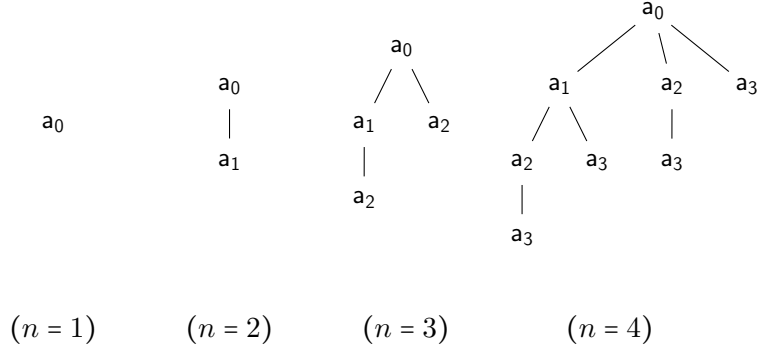


Figure 5.1: Prefix-tree for saturated non-recursive schemas with  $n$  types

From the previous results we have that trees are just linearly more succinct than lists for storing chains.

By considering recursive DTDs, the number of chains inferred for an expression augments substantially. We investigate the cardinality of sets of inferred chains for recursive schemas. To this end we adopt multisets, defined as follows.

**Definition 5.1.4** (Multiset). *A multiset  $(M, \mu)$  is a set of (distinct) elements  $M$  together with a function  $\mu$  that assigns to each element in  $M$  a non-negative integer value of multiplicity.*

For instance,  $(\{a, b\}, \{(a, 1), (b, 3)\})$  denotes the multiset in which the element  $a$  has one occurrence and  $b$  three occurrences. We say that  $(M', \mu')$  is a sub-multiset of  $(M, \mu)$  if  $M' \subseteq M$  and  $\mu'(m) \leq \mu(m)$ , for all  $m \in M$ . For instance, let  $M = \{a, b\}$ ,  $\mu = \{(a, 1), (b, 3)\}$  and  $\mu' = \{(a, 0), (b, 1)\}$ , the multiset  $(M, \mu')$ , corresponding to the set  $\{b\}$ , is a sub-multiset of  $(M, \mu)$ . In short, we denote by  $M_k$  the multiset  $(M, \mu)$  where each element of  $M$  has multiplicity  $k$ , so that  $\mu(a) = k$  for all  $a \in M$ .

**Proposition 5.1.5.** *A saturated recursive DTD  $\mathbf{d}$  with  $n + 1$  types induces  $\Omega(\frac{(nk)!}{(k!)^n})$  and  $O(\frac{(nk+1)!}{(k!)^n})$  rooted  $k$ -chains.*

*Proof.* Because  $\mathbf{d}$  is saturated and recursive, a rooted  $k$ -chain is a sequence of types, whose length is upper bounded by  $1 + nk$ . We model this by using multisets. Let  $M$  be the set of  $n$  mutual recursive types in  $\mathbf{d}$ , and let  $M_k$  be the multiset of  $M$  elements with multiplicity  $k$  each. A permutation of a multiset  $M$  is an ordered arrangement of objects of  $M$ . The set of rooted  $k$ -chains of  $\mathbf{d}$  can be seen as the *permutations with repetition* of all sub-multisets of  $M_k$ . For instance, let  $(M, \mu')$  be the multiset with  $M = \{m_1, \dots, m_n\}$  and

$$\mu' = \{ (m_1, k-1), (m_2, k), \dots, (m_n, k) \}$$

We have that  $(M, \mu')$  is a sub-multiset of  $M_k$ . The permutations associated with  $(M, \mu')$  can be seen as the set of all chains of length  $nk$ , where all types repeat exactly  $k$  times, but type  $m_1$  repeats exactly  $k - 1$  times. Of course, if  $b^1, b^2$  are the first and second occurrence of  $b$ , respectively, then  $ab^1b^2$  and  $ab^2b^1$  denote the same permutation.

The number of  $k$ -chains for  $\mathbf{d}$  can be computed exactly with the basic combinatorial formula for permutations with repetition, as the sum of all possible configurations of multiplicity values  $k_1, \dots, k_n$  (assigned to type-names  $m_1, \dots, m_n$ ) that are ranging from 0 to  $k$ . This is

$$\sum_{0 \leq k_1, \dots, k_n \leq k} \frac{(k_1 + \dots + k_n)!}{k_1! \dots k_n!}$$

also counting the singleton root-type.

A lower bound for this quantity can be given by computing just part of the inferred chains, namely the permutations of  $M_k$  of maximum length  $nk$  (we ignore subsets for the moment). This can be done with the basic formula and results in

$$\text{perm}(M_k; nk) = \frac{(nk)!}{k_1! \dots k_n!}$$

with  $k_1 = \dots = k_n = k$ . This quantity is exactly the number of rooted  $k$ -chains of length  $nk$ , and thus we have the lower bound  $\Omega(\frac{(nk)!}{(k!)^n})$ .

An upper bound for the set of  $k$ -chains induced by  $\mathbf{d}$  is obtained in the following way. First, notice that the whole set of rooted  $k$ -chains of  $\mathbf{d}$  is the set of prefixes of the chains of maximal length we computed above. Therefore, consider that each chain of length  $1 + nk$  has  $nk$  prefixes. So as to obtain

$$nk \frac{(nk)!}{(k!)^n}$$

which gives  $O(\frac{(nk+1)!}{(k!)^n})$  as upper bound of inferred chains. This upper bound is not tight, because we count twice a prefix shared by two chains, but still helpful in order to understand the result size. It is worth observing that  $\text{perm}(M_k; nk)$  is strictly smaller than the number of permutations with repetitions of length  $nk$  of  $n$  types, that is  $n^{nk}$ .  $\square$

These results show that chain inference is much harder in the presence of recursive schemas. Once again, it is possible to show that by using both lists and prefix-trees for storing chains, the space occupancy is exponential in the size of the DTD and the value  $k$ ; prefix trees are still polynomially more succinct than lists. Because space complexity has a direct impact on time complexity, this suggests that implementing the chain inference defined in Chapter 4 with one of these two data-structures may compromise scalability, as shown by the following result.

**Theorem 5.1.6.** *Let  $d$  be a DTD with  $n$  types and  $e$  be a query or update expression of size  $m$ , exact  $k$ -chain inference for  $e$  is in EXPTIME.*

*Proof.* By Lemma 5.1.5 the maximal number of chains one can infer from DTD  $d$  is in  $O(\frac{(nk+1)!}{(k!)^n})$ . Consider now the case where *each* subexpression of  $e$  needs to process all of such chains. All inference rules perform a linear number of operations on a single chain, but rules for insert and replace updates perform a quadratic number of operations. Hence we say that each inference step does at most a polynomial number of operations on the set of input chains. In turn, each operation on a chain can be done in  $O(nk)$ , that corresponds to the traversal of the whole chain. We conclude that exact chain inference for the whole expression of size  $m$  can be done in  $O(m \cdot nk \cdot \text{poly}(\frac{(nk+1)!}{(k!)^n}))$ .  $\square$

This theorem gives an upper bound to the complexity of the chain inference problem and it shows that exact chain inference may have an enormous worst-case cost by using data-structures such as lists or trees. Of course this does not mean that for any instance of the problem, chain inference is expensive. However, in practice it turns out that for schemas featuring cliques of mutual recursive types, the  $k$  analysis works efficiently only if  $k$  is small, and rapidly becomes intractable. From this we can deduce also that, once chains are inferred from queries and updates (in exponential time), checking for a conflict may require an exponential time, because of the number of chains one has to consider.

## 5.2 Storing Chains on DAGs

To overcome the limits outlined in the previous section, we present now a data-structure for storing chains based on directed acyclic graphs (DAGs) which is exponentially more succinct than trees, and allows to store chains in polynomial space and makes the analysis tractable in practice.

The intuition behind the use of DAGs instead of trees is that rather than storing only once common-prefixes of chains, in many cases we can also store only once common suffixes of chains. This can be done in a particularly efficient way for XML schemas that, being regular tree grammars and generate regular trees-of-chains featuring a number of repeating subtrees.

To illustrate, consider the chains induced by the non-recursive DTD  $d_8$ , which is defined by the following set of productions (below, we assume  $1 \leq i \leq 3$ ).

$$a \leftarrow b_1, f_1 \quad \dots \quad b_{(i-1)} \leftarrow (b_i, f_i)? \quad f_{(i-1)} \leftarrow (b_i, f_i)?$$

In Figure 5.2.a we show how the chains of  $d$  are stored as a prefix-tree. From now on, we call CTREE a prefix-tree storing chains. The CTREE shows clearly that by unfolding schema chains we are unfolding regular subtrees. We stress this fact by marking repeated subtrees with different colors.

In Figures 5.2.(b-d) we show that repeated subtrees can be recursively compressed until obtaining a DAG that is equivalent to the CTREE. From now on, we call CDAG a DAG storing chains. In fact, when there is a path (equivalently, a chain) in the CTREE, there is a path in the CDAG, and vice-versa.

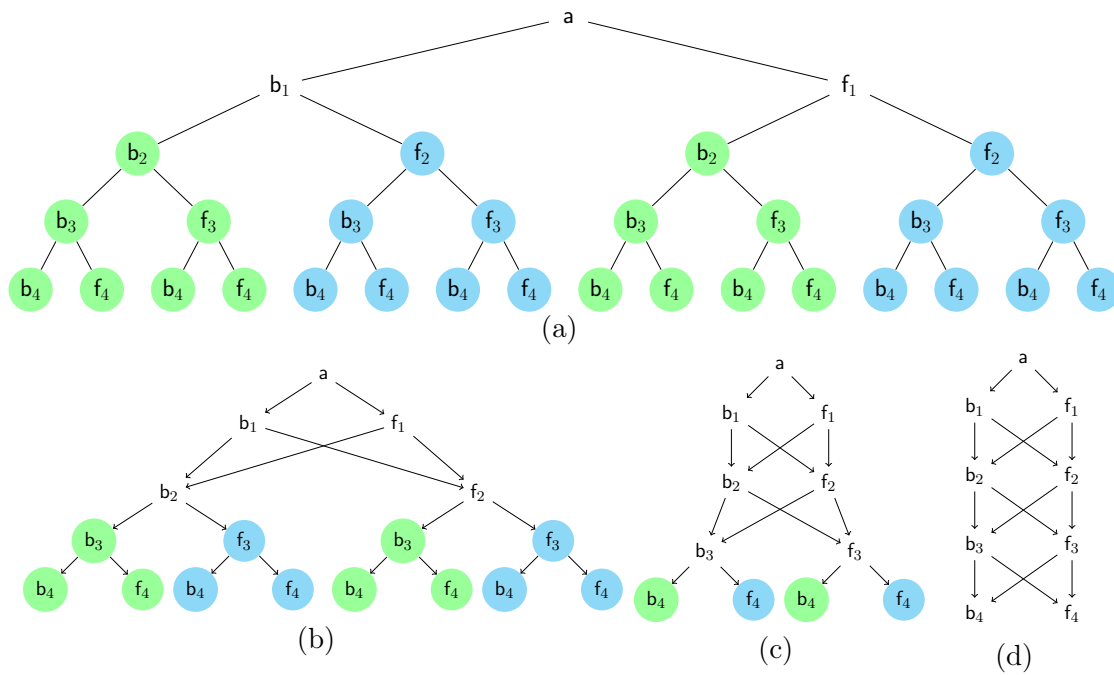


Figure 5.2: Compression of a CTREE to an equivalent CDAG

As suggested by the example, it is more convenient to perform the inference of the CDAG rather than the tree. Of course, we avoid to first infer prefix-trees and then compress them, because this again lead us to infer structures of exponential size. In the following, when speaking about CDAGs and CTREES, we assume that they are inferred on a given DTD  $d$ .

A CDAG  $G$  is a directed acyclic graph rooted at the schema root-type, that contains no self-loop and meets the following property.

$$\text{There is at most one node labeled } a \text{ at depth } h \text{ in } G \quad (5.1)$$

This means that, if two chains happen to have the same type in position  $h$ , they will share a common node in the CDAG at depth  $h$ . This property is for instance enjoyed by the tree in Figure 5.2.d. A CDAG is thus organized per levels.

Property (5.1) makes the *width* of the CDAG upper-bounded by the number of types in the schema size. Also, the *k*-analysis makes the *height* of a CDAG linear in the size of the schema. Therefore the number of vertex of a CDAG is always finite, and polynomial in  $d$  and  $q$ .

Let us illustrate property (5.1) with some examples over recursive and non-recursive schemas.

Consider the saturated non-recursive schema  $d_2$  previously defined.

$$\begin{aligned} a_1 &\leftarrow a_2, a_3, a_4, a_5 \\ a_2 &\leftarrow a_3, a_4, a_5 \\ a_3 &\leftarrow a_4, a_5 \\ a_4 &\leftarrow a_5 \end{aligned}$$

In Figure 5.3 we compare a CTREE and a CDAG storing all chains of  $d_2$ . Once again, we see the unfolding of the regular structure of the schema with the CTREE. By using DAGs instead of trees we can compress common chain-prefixes and store them only once. This is the case, for instance, for the subtrees rooted at  $a_4$  and  $a_5$  at level 2. In this case one could even compress more, by ignoring levels. However, this would substantially complicate not only the formal development, but also the implementation phase, thus resulting error prone for implementors.

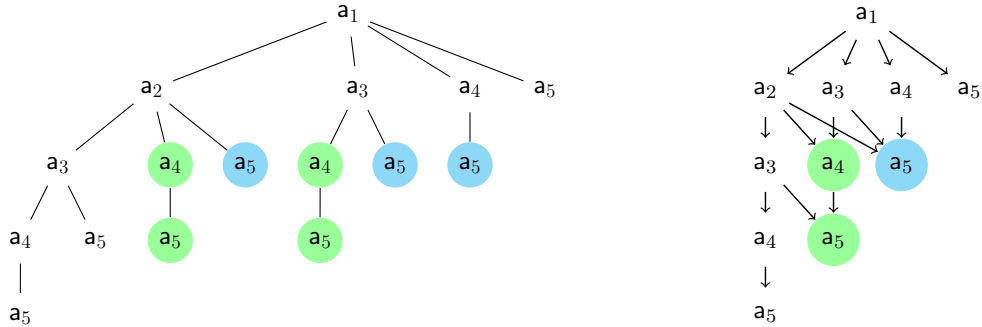
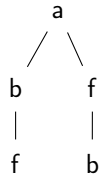


Figure 5.3: CTREE and CDAG storing chains of  $d_2$

Consider now the saturated recursive schema  $d_9$  defined at the beginning of the section.

$$\begin{aligned} a &\leftarrow b, f \\ b &\leftarrow b^*, f? \\ f &\leftarrow b?, f^* \end{aligned}$$

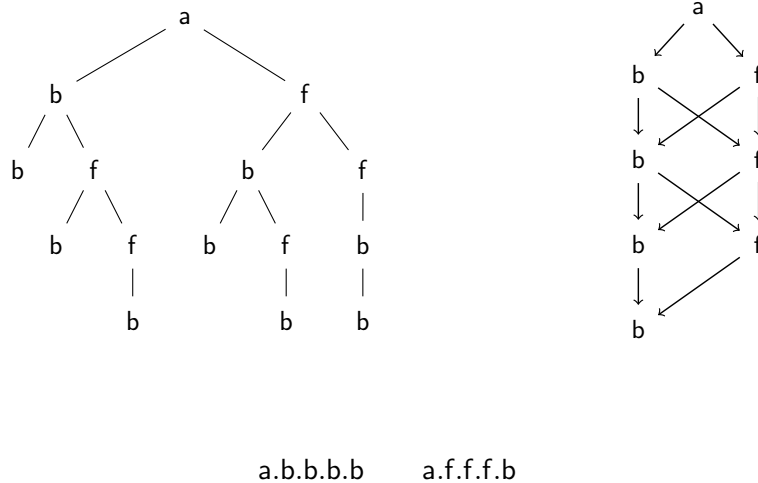
Let us consider all 1-chains induced by  $d_9$ . The CTREE and CDAG representing such chains are the same, as depicted in the following. Recall that each rooted path is a chain,



and 1-chains for  $d_9$  are  $\{a, a.b, a.f, a.b.f, a.f.b\}$ .

Now, let us consider all 2-chains for  $d$  terminating with type  $b$ . This is, for instance, the specification of chains matching a recursive navigation made by the query  $//b$ . Below we compare a CTREE and a CDAG storing such chains.

A chain in the graph is a rooted path ending with type  $b$ . Notice that a chain belongs to the CDAG if it belongs to the CTREE but, in this case, there are also some chains belonging to the CDAG that do not belong to the CTREE. The following chains are obtained as artifacts from the compression.



While this could suggest that the analysis imprecise (in fact, by inferring more chain than needed we may have false-negatives) it does not. The reason is that, because of the regularity of tree grammars, chains  $a.b.b.b.b$  and  $a.f.f.f.b$ , are indeed 4-chains matching  $//b$  over  $d_9$ . Therefore in this case chain inference for  $//b$  with CDAGs is both sound and complete, in the sense that it infers all 2-chains matching  $//b$  (soundness) and only chains matching  $//b$  (completeness).

While at first sight one could argue that, for the purpose of independence detection, a refined notion of completeness (wrt  $k$ -chains) may be needed, this is not the case. Indeed a  $k'$ -chain matching an expressions, with  $k' > k$ , cannot generate false negatives.

To conclude the discussion on the size of CDAGs, we prove that they are always polynomial in  $\mathbf{d}$ .

**Lemma 5.2.1.** *Let  $\mathbf{d}$  be a non-recursive DTD with  $n + 1$  types. The size of a CDAG inferred for all chains of  $\mathbf{d}$  is in  $O(n^3)$ .*

*Proof.* Assume that  $\mathbf{d}$  is a saturated non-recursive schema featuring  $n + 1$  types. The CDAG storing all  $\mathbf{d}$  chains is organized in  $n + 1$  levels where a) the root-level has one node and b) for the remaining  $n$  levels, level  $i$  has  $i$  nodes ( $1 \leq i \leq n$ ). Therefore we conclude that the CDAG has  $1 + \frac{n(n+1)}{2}$  nodes. Concerning the number of edges, we notice first, that, at the root level, the root type is connected with all other  $n$  types. Then, as we said, level  $i$  has  $i$  types. Assume that  $\mathbf{a}_2, \dots, \mathbf{a}_i$  is the sequence of types at level  $i$ , ordered by the number of outgoing edges (ascending order). Because the schema is saturated, we have that type  $\mathbf{a}_j$  at level  $i$  is connected with  $j - 1$  types at level  $i + 1$ . We conclude that the size of the edge-set is  $n + \sum_{i=1}^n \sum_{j=0}^{i-1} j$ , which is in  $O(n^3)$ .  $\square$

**Lemma 5.2.2.** *Let  $\mathbf{d}$  be a recursive DTD with  $n$  types and  $k$  a multiplicity value. The size of a CDAG inferred for  $k$ -chains of  $\mathbf{d}$  is in  $O(kn^3)$ .*

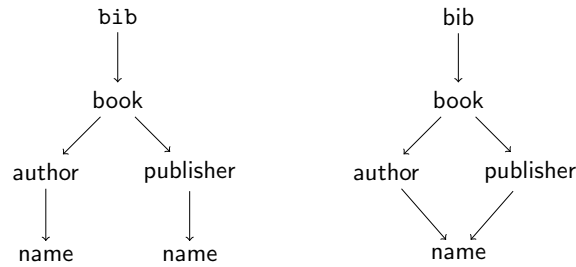
*Proof.* Assume that  $\mathbf{d}$  is a saturated recursive schema featuring  $n+1$  types. First we notice that the CDAG storing the  $k$ -chains of  $\mathbf{d}$  has at most  $kn$  levels. Therefore, because each of the  $n$  types can repeat in all  $kn$  levels, the CDAG has  $kn^2$  nodes. Because there exists at most  $n^2$  edges between two consecutive graph levels, the CDAG has  $O(kn^3)$  edges. We conclude since no edge crosses multiple levels.  $\square$

### Chain Inference and Edge-labeling

CDAGs do not allow only to reduce space occupancy, but also to reduce the number of operations needed by the inference. To illustrate, consider the query  $q = //book/*name$  and the DTD  $d_3$  previously defined as

bib	←	book*
book	←	title, author+, publisher+
author, publisher	←	name
title, name	←	String

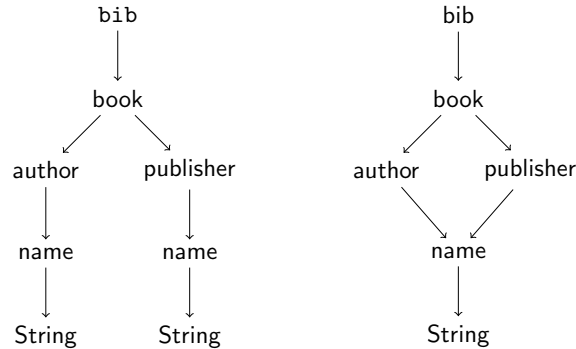
The CTREE and the CDAG storing chains of  $q$  over  $d_3$  are reported below.



Now, assume a further inference step has to be done, such as for query

$$q' = q / \text{text}()$$

Now, in order to extend the CTREE, two operations are needed (one for each node **name**) while, in order to extend the CDAG, only one operation is needed.



This shows that, with CDAGs, we can also save computations by doing chain inference by groups of chains. The number of groups of chains is polynomial in the size of a schema and of a query (in fact, a group is represented by all chains terminating in one node). Thus this makes polynomial the number of input nodes for an intermediate inference operation.



Another ingredient needed in the CDAG representation is a means to represent chains inferred for unrelated sub-expressions. Intuitively two expressions are related if one draws data from the other (e.g. by means of variable binding). If this does not happen, the expressions are unrelated. To illustrate, consider the bibliographic DTD  $d_3$  and the following sequence

$$q_1, q_2 = //book//name/text(), /bib/book/publisher/name/foo$$

Here the whole expression is constituted of two unrelated queries  $q_1$  and  $q_2$ . It turns out that these may produce unneeded chains as a consequence of overlapping.

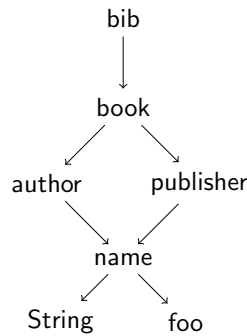
The chains inferred for  $q_1$  over  $d_3$  are

bib.book.author.name.String      bib.book.publisher.name.String

The chain inferred for  $q_2$  over  $d_3$  is

bib.book.publisher.name.foo

The naive CDAG representation of the chains for the two subexpressions is

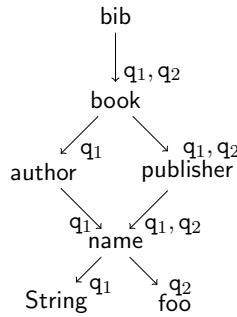


If the chains inferred for the two expressions are stored together, without any distinction, then the CDAG also encodes the chain

bib.book.author.name.foo

which is not a chain that should be inferred neither for  $q_1$  nor  $q_2$ .

In order to overcome ambiguity in the CDAG, we introduce some marking for distinguishing among two chains inferred for distinct sub-expressions in the CDAG representation. We do this by putting on any edge connecting two CDAG-nodes a *label* identifying the query/update expressions that produced the chain during inference. These identifiers are necessary to precisely perform chain inference, and independence checking as well. The CDAG representation of chains for  $q_1, q_2$  is the following.



We see that the labelled CDAG allows for a more precise analysis, in this case for instance the chain `bib.book.author.name.foo` cannot be traced anymore. (through  $q_1$  or  $q_2$  identifiers).

Labels help also in dealing with backward navigations. We can observe that if in the above query  $q_1, q_2$  we replace  $q_2$  by

$$q_2 = /bib/book/publisher/foo/ancestor::*$$

when inferring chains for the last step of  $q_2$ , thanks to edge-labeling, avoids to navigate upward parts of the CDAG that have not been generated by  $q_2$  (i.e., a `author` node), thus preserving the precision of the analysis in most of the cases.

Notice that backtracking on unvisited nodes does not affect the correctness of the analysis, but only its precision. The same observation holds for tracing unvisited chains. Removing labels, is a possibility if one wants to further lower the complexity of the analysis, that would result in dealing with a dimension less.

### 5.3 Chain Inference with CDAGs

In this section we present algorithms for implementing the chain inference with CDAGs. We first give some auxiliary definitions, and then we present chain inference. We begin with operators for XPath single steps, then treat navigational queries, element construction and, finally, we consider updates.

#### Preliminary definitions

**Definition 5.3.1** (CDAGs). *A CDAG is a directed acyclic rooted graph  $G = (V, E, s)$ , where  $V \subseteq \Sigma_S \times \mathbb{N}$  is a set of nodes that store a type in  $\Sigma_S$  and an integer value denoting the distance of the node from the root in  $G$ , where node  $(s, 0) \in V$  denotes the root type, and  $E \subseteq V^2 \times \mathbb{N}$  is a set of edges among nodes recording the code of the query/update subexpression that produced that edge.*

In order to implement our analysis we need some auxiliary tools for managing chains inferred from queries and updates.

We first need an *encoding* of query and update expressions that is used for separating the chains inferred from unrelated subexpressions. The concept of subexpression is standard:  $e'$  is a subexpression of  $e$ , denoted by  $e' \in e$ , if the parse tree of  $e'$  is a subtree of the parse tree of  $e$ .

**Definition 5.3.2** (Query encoding). *The encoding of a query or update expression  $e$  is an injective function  $id()$  mapping each subexpression of  $e$  to an integer in  $\mathbb{N}$ . The encoding of a set of expressions  $\{e_1, \dots, e_n\}$  is a function obtained by the union of functions encoding each  $e_i$ , with  $1 \leq i \leq n$ , provided that this union is still injective.*

The definition says that basically any coding would work for the inference task, provided that we are able to uniquely identify each subexpression. This is why the function is required to be injective. As an example, consider the sequence

$$q_0 = (q_1, q_2)$$

for which we can have the following encoding.

$$id(e) = \begin{cases} 1 & \text{if } e = q_1 \\ 2 & \text{if } e = q_2 \\ 12 & \text{if } e = q_0 \end{cases}$$

provided that there are no other subexpressions in  $q_1, q_2$ .

The second tool we need is a relation for capturing related expressions. As already said, two expressions are related if one draws data from the other. As an example, for the query

`for x in //book return x/author`

expressions `//book` and `x/author` are related because the first provides nodes that are bound to  $x$ , and the second accesses the content of  $x$ . As the example suggests, in the language we considered this relationship is established by the use of variables, and allows to reconstruct all navigational-paths embodied by a query or an update. We assume wlog that queries are written by using distinct variable-names for distinct variable-definitions.

To formalize this, we need a way to retrieve, given a variable  $x$ , the set of navigational steps that provide data to  $x$ . Assume that

$$e = \text{for/let } x \text{ in/} := e_1 \text{ return } e_2$$

then the set of navigational expressions providing data to  $x$  is defined as  $ret(e_1)$ , where

$$ret(e) \stackrel{def}{=} \begin{cases} \{e\} & \text{if } e = y/\text{step} \\ ret(e_2) & \text{if } e = \text{for/let } y \text{ in/} := e_1 \text{ return } e_2 \\ ret(e_1) \cup ret(e_2) & \text{if } e = \text{if } (e_0) \text{ then } e_1 \text{ else } e_2 \text{ or } e = e_1, e_2 \\ \emptyset & \text{otherwise} \end{cases}$$

To illustrate function  $ret()$ , consider the following query

$q = \text{for } x \text{ in } q_1 \text{ return } q_2$

with  $q_1 = (//title, //author)$  and  $q_2 = x//book$ . In this case we have that

$$ret(q) = \{ x//book \} = ret(q_2) \quad ret(q_1) = \{ //title, //author \}$$

We now define expressions related by the use of variables.

**Definition 5.3.3** (Variable-Reference). *Let  $e$  be a query or update expression, and  $x$  a variable defined in  $e$ . The set of navigational steps referenced by  $x$ , denoted by  $vr(x)$ , is defined as  $vr(x) = ret(e_x)$ , provided that  $e_x$  is a for/let expression introducing  $x$ .*

To illustrate this definition, consider the query

```
q =  for x in
      (for y in //a, //d return y//b)
    return
      let z :=
        if (//b) then x else x/parent :: *
      return x
```

The variable-reference relation, and its transitive closure, are graphed below.

$$\begin{array}{ll}
 y \xrightarrow{vr(y)} //a, //d & x \xrightarrow{vr(x)} y//b \xrightarrow{vr(y)} //d \\
 x \xrightarrow{vr(x)} y//b \xrightarrow{vr(y)} //a & z \xrightarrow{vr(z)} x, x/parent :: * \\
 \\ 
 z \xrightarrow{vr(z)} x \xrightarrow{vr(x)} y//b \xrightarrow{vr(y)} //a & z \xrightarrow{vr(z)} x/parent :: * \xrightarrow{vr(x)} y//b \xrightarrow{vr(y)} //a \\
 z \xrightarrow{vr(z)} x \xrightarrow{vr(x)} y//b \xrightarrow{vr(y)} //d & z \xrightarrow{vr(z)} x/parent :: * \xrightarrow{vr(x)} y//b \xrightarrow{vr(y)} //d
 \end{array}$$

With a little abuse of notation, we denote by  $vr(e)$ ,  $vr(i)$  and  $vr(x)$  the same set, provided that  $e = x/step$  and  $id(e) = i$ . The transitive closure of  $vr(x)$ ,  $vr(e)$  and  $vr(i)$  are denoted by  $vr^+(x)$ ,  $vr^+(e)$  and  $vr^+(i)$ , respectively.

The following property shows that an exponential number of sequences induced by the reference variable relation may arise.

**Proposition 5.3.4.** *Let  $\mathbf{e}$  be a query or update expression and  $\mathbf{e}_1, \dots, \mathbf{e}_n$  be a sequence of subexpressions of  $\mathbf{e}$  such that  $\mathbf{e}_{(i+1)} \in vr(\mathbf{e}_i)$ . The following holds.*

1. *the sequence size is bounded by the query size,  $n \sim O(|\mathbf{e}|)$*
2. *the number of such sequences is in the worst case exponential  $O(2^n)$*
3. *if no concatenation or conditional expression is nested in the left-branch of a for or let expression, then the number of such sequences is linear  $O(n)$ .*

The first item states that the *length* of a navigational path embodied in a query is upper bounded by the query size. The second item shows the succinctness of the query and update languages. It states that a query of size  $n$  can embody  $2^n$  navigational paths. The third item states that this succinctness is given by two operators: query concatenation and conditional expressions.

Proposition 5.3.4 suggests also that chain inference may lead to perform an exponential number of operations if the whole relation  $vr^+(\cdot)$  has to be visited. For instance, this may happen in the presence of backward navigations. Nevertheless, differently from trees and lists, CDAGs allow to design algorithms that are sound, that run in polynomial time and that are still precise, as we illustrate next.

**Auxiliary Data-Structures** The CDAG inference algorithms make use of some global data structures that we define below. The role of such structures will be clarified during the presentation of inference algorithms.

- $RET()$  is an *inverted index* mapping an expression id  $i$  to a set of nodes  $M$  in the CDAG. The set of nodes  $N$  corresponds to terminal of *return* chains inferred for  $\mathbf{e}$ .
- $UPD()$  is an *inverted index* mapping an expression  $\mathbf{e}$  to a set of nodes  $M$  in the CDAG. The set of nodes  $N$  corresponds to nodes in the suffix of an update chain inferred for  $\mathbf{e}$ .
- $CLEAN()$  is an *inverted index* mapping a step  $\mathbf{e}$  to a set nodes representing dangling chains for which inference for  $\mathbf{e}$  failed.

The input of the chain inference algorithm is a query or update expression  $\mathbf{e}$ , a schema  $\mathbf{d}$  and a multiplicity value  $k$ . The system outputs a CDAG  $G_{\mathbf{r}}$  representing used and return chains, a CDAG  $G_{\mathbf{e}}$  representing element chains and a  $G_{\mathbf{U}}$  representing update chains inferred from  $\mathbf{e}$ , respectively. Note that used and returned chains are put together in the same CDAG.

Given a CDAG  $G$ , and node  $v$  belonging to  $G$  is a pair  $v = (\mathbf{a}, h)$ , where  $\mathbf{a}$  is the type of the node and  $h$  is the level of the node in the  $G$ . The type of a node  $v$  is denoted by  $v.type$  while its height is denoted by  $v.level$ . In the examples, a node is also denoted by  $v_{[type, height]}$ .

## Chain Inference Algorithms

**Algorithm 1:** Procedure CHILD

---

**Input** : Node-set  $M \subseteq V_{rv}$ ,  $\phi$ , id  $i$   
**Output**: Infer chains for  $\text{child} :: \phi$ , side-effect on  $\text{RET}()$ ,  $\text{CLEAN}()$ ,  $V_{rv}$  and  $E_{rv}$

```

1 foreach  $v \in M$  do
2    $\text{new\_edge} \leftarrow \text{false}$ 
3   if  $v.\text{level} + 1$  does not exceeds the chain space then
4     foreach  $a$  child-type of  $v.\text{type}$  do
5       if  $a$  satisfies  $\phi$  then
6         if  $(a, v.\text{level} + 1) \notin V_{rv}$  then
7            $V_{rv} \leftarrow V_{rv} \cup \{(a, v.\text{level} + 1)\}$ 
8            $v' \leftarrow (a, v.\text{level} + 1)$ 
9            $E_{rv} \leftarrow E_{rv} \cup (v, v', i)$ 
10           $\text{new\_edge} \leftarrow \text{true}$ 
11           $\text{RET}(i) \leftarrow \text{RET}(i) \cup \{v'\}$ 
12   if  $\text{!new\_edge}$  then
13      $\text{CLEAN}(i) \leftarrow \text{CLEAN}(i) \cup \{v\}$ 

```

---

The first procedure we present is the one to infer chains for a navigational expression of the form  $\mathbf{x}/\text{child} :: \phi$ . Procedure CHILD is the most important procedure of the system, since it is the one responsible of extending the CDAGs. Other operators presented later on are defined in terms of the procedure CHILD.

The procedure takes as input a set of nodes  $M \subseteq V_{rv}$ , and tries to connect each  $v \in M$  with a child node  $v'$ , provided that  $v'$  satisfies the filtering condition  $\phi$ . If such  $v'$  is not a node of the CDAG, then  $v'$  is built and added to the vertex-set of the CDAG (line 7). Then  $v'$  becomes a terminal node of the expression invoking the **child** axis. The step expression  $\mathbf{x}/\text{child} :: \phi$ , for which the procedure infers chains, is referred by identifier  $i$ , which is also part of the input.

In line 3, we have a condition that ensures the termination of the inference. It says that a child navigation could be performed, provided that chains do not exceed a certain length. The maximal length of a chain we infer is  $nk$ , provided that that  $n$  is the number of types in  $\mathbf{d}$ , and  $k$  is the multiplicity value in input. This is a sufficient condition to ensure that all  $k$  chains are correctly inferred. Note also that this avoids to keep an hash table recording the frequency of each tag while doing the chain inference, thus ensuring better performances.

Finally, in the cases that the child extension for the input node  $v$  produces no new edges, then  $v$  is added to the list of possibly dangling nodes to backtrack (Lines 12-13).

This happens, for instance, if  $v$  exceeds the chain space, or none of the children on  $v$  satisfy  $\phi$ . A dangling node is later removed, provided that itself or one of its descendants are terminals of used or return chains in the CDAG. This essentially implements the chain filtering defined in chain inference rules of Chapter 4. For the sake of conciseness, we do not model this operation here.

Let us illustrate the workflow of the procedure with an example. Consider chain inference for the query  $\mathbf{x}/book$  over DTD  $\mathbf{d}_3$ , with  $\mathbf{x} \mapsto \mathbf{bib}$ . Inference for the step takes in input the following CDAG

$$G_{rv} = (\{ v_{[\mathbf{bib},0]} \}, \emptyset)$$

and a set  $M = \{ v_{[\mathbf{bib},0]} \}$ . The procedure first checks that the chain space would not be exceeded by attaching a child to  $v_{[\mathbf{bib},0]}$ , which is trivially true since

$$v_{[\mathbf{bib},0]}.level < nk$$

where we recall that  $n$  is the number of types in  $\mathbf{d}_3$  (in this case  $n = 8$ ) and  $k$  is the multiplicity value (in this case  $k = 1$ ).

At this point, because **book** is the only child-type of **bib**, and **book** satisfies  $\phi=book$ , the node  $v_{[\mathbf{book},1]}$  is added to  $V$  and connected to  $v_{[\mathbf{bib},0]}$ . Below, we assume  $id(\mathbf{x}/\mathbf{child} :: \phi) = i$ . Node  $v_{[\mathbf{book},1]}$  is now the terminal of chains inferred for the query. The resulting  $G_{rv}$  and  $RET()$  are the following.

$$\begin{array}{lll} \begin{array}{c} \mathbf{bib} \\ \downarrow i \\ \mathbf{book} \end{array} & \begin{array}{l} G_{rv} = (V_{rv}, E_{rv}) \\ V_{rv} = \{ v_{[\mathbf{bib},0]}, v_{[\mathbf{book},1]} \} \\ E_{rv} = \{ (v_{[\mathbf{bib},0]}, v_{[\mathbf{book},1]}, i) \} \end{array} & RET(i) = \{ v_{[\mathbf{book},1]} \} \end{array}$$

The second procedure defined is the one implementing the **self** axis. Procedure **SELF** receives as input a set of nodes  $M \subseteq V_{rv}$  and only checks if the type-name of each input node satisfies condition  $\phi$ . If this is true, the input node becomes a terminal for the expression. This is captured by line 3.

---

**Algorithm 2:** Procedure **SELF**


---

**Input** : Node-set  $M \subseteq V_{rv}$ ,  $\phi$ , id  $i$

**Output**: Infer chains for **self** axis, side-effect on  $RET()$

```

1 foreach  $v \in M$  do
2   if  $v.type$  satisfies  $\phi$  then
3      $RET(i) \leftarrow RET(i) \cup \{v\}$ 

```

---

Notice that, if a node does not satisfy  $\phi$ , then no operation is performed here. This is made on purpose, nodes that do not satisfy the axis test should be dealt with by the external procedures calling SELF.

Consider chain inference for query  $y/\text{self} :: *$  over DTD  $d_3$ , provided that  $y$  is bound to the result of  $x/\text{book}$  and  $x \mapsto \text{bib}$ , computed as in the previous example. Chain inference for  $y/\text{self} :: *$  results in the following CDAG. As before, we assume that  $id(x/\text{book}) = i$  and  $id(y/\text{self} :: *) = i'$ .

$$\begin{array}{lll}
 \begin{array}{c} \text{bib} \\ \downarrow i \\ \text{book} \end{array} & G_{rv} = (V, E) & \begin{array}{l} \text{RET}(i) = \{ v_{[\text{book},1]} \} \\ \text{RET}(i') = \{ v_{[\text{book},1]} \} \end{array} \\
 & V = \{ v_{[\text{bib},0]}, v_{[\text{book},1]} \} & \\
 & E = \{ (v_{[\text{bib},0]}, v_{[\text{book},1]}, i) \} & 
 \end{array}$$

Note that procedure SELF adds no edge in the CDAG, but it just affects RET(). In this case because node  $v_{[\text{book},1]}$  trivially satisfies  $\phi = *$  we have  $\text{RET}(i') = \{ v_{[\text{book},1]} \}$ .

The algorithm implementing the descendant axis is described in the following. We just describe the algorithm implementing axis **descendant** since **descendant-or-self** is similar.

DESCENDANT-REC is a recursive procedure that computes chains for navigational steps of the form  $x/\text{descendant} :: \phi$ . The descendant axis is computed by invoking the CHILD procedure. Because of this, in the pseudo-code, we use a fresh query-id  $i_v$  for labeling edges obtained from the iteration of the procedure, that started from node  $v$ . This is simply needed to correctly perform set operations and clean the graph from artificial edges.

At each iteration step, the procedure tries to extend of one level each input node  $v$ , with all of its children-types (line 5), by calling procedure CHILD. Before calling procedure CHILD, the abstract predicate satisfies uses schema information to check if some descendant-type of  $v.type$  do satisfy the filter  $\phi$  (line 4). When this is not true, any operation can be skipped and node  $v$  can be added to the backtrack list for the expression.

If some descendant-type of  $v.type$  do satisfy the filter  $\phi$  but the result of the child-navigation with  $*$ -filter is empty, it simply means that the node is at the frontier of the chain space. In this case the node is added to the list of possibly dangling nodes (line 14). If the result of the child navigation is non-empty (line 6), then we record that some new labeled-edges has been created (line 7). All nodes in  $\text{RET}(i_v)$  are saved in the set  $M'$  (line 8), and the descendant procedure is iterated on  $M'$  (line 15). Also, all nodes in  $\text{RET}(i_v)$  that satisfy  $\phi$ , are set as terminal nodes for the expression, by means of the



**Algorithm 3:** Procedure DESCENDANT-REC

---

**Input** : Node-set  $M \subseteq V_{rv}$ ,  $\phi$ , id  $i$   
**Output**: Infer chains for **descendant** axis, side-effect on  $\text{RET}()$ ,  $\text{CLEAN}()$ ,  $V_{rv}$  and  $E_{rv}$

---

```

1  $M' \leftarrow \emptyset$ 
2 foreach  $v \in M$  do
3    $\text{new\_edge} \leftarrow \text{false}$ 
4   if some descendant-type of  $v$ .type satisfies  $\phi$  then
5      $\text{CHILD}(\{v\}, *, i_v)$ 
6     if  $\text{RET}(i_v) \neq \emptyset$  then
7        $\text{new\_edge} \leftarrow \text{true}$ 
8        $M' \leftarrow M' \cup \text{RET}(i_v)$ 
9        $\text{SELF}(\text{RET}(i_v), \phi, i)$ 
10      foreach  $(v, v', i_v) \in E$  do
11         $E \leftarrow E \setminus (v, v', i_v)$ 
12         $E \leftarrow E \cup (v, v', i)$ 
13    if  $\text{new\_edge}$  then
14       $\text{CLEAN}(i) \leftarrow \text{CLEAN}(i) \cup \{v\}$ 
15  $\text{DESCENDANT-REC}(M', \phi, i)$ 

```

---

SELF procedure (line 9). At this point, all edges produced by the CHILD procedure with the input code  $i_v$  are re-labeled as  $i$  edges. Notice that  $i_v$  is needed to perform this operation.

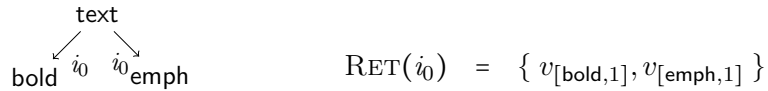
We illustrate chain inference for step **x/descendant** :: *bold* with **x** bound to **text** over the recursive DTD  $d_{12}$  defined as

```

text  ← (bold | emph)*
bold  ← emph?, String
emph  ← bold?, String

```

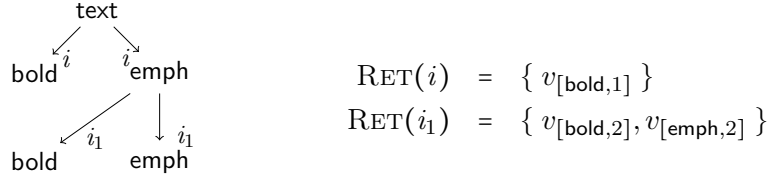
In this case  $k_q = 1$  and thus we want to infer 1-chains for the query, namely **text.bold** and **text.emph.bold**. To this aim, procedure DESCENDANT-REC is iterated three times. In the first iteration, the procedure takes as input the triplet  $(\{v_{[\text{text},0]}\}, \phi, i)$ . A child navigation from the root is then performed. This gives the following intermediate result. Below,  $i_0$  is used to denote  $i_{v_{[\text{text},0]}}$ .



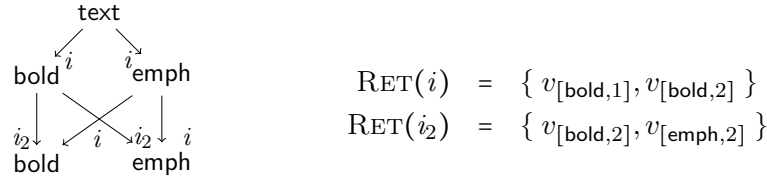
Now, because node  $v_{[\text{bold},1]}$  satisfies  $\phi$ , it is set as terminal node for the query, and we have  $\text{RET}(i) = \{ v_{[\text{bold},1]} \}$ . Then edges labeled as  $i_0$  are relabeled with  $i$ . We assume also that the terminal-index  $\text{RET}(i_0)$  is set to empty. However, since this is not mandatory for the correctness of the analysis we do not add it in the pseudocode. The result of the first iteration is the following.



In the second iteration, the procedure takes as input the triplet  $(\{v_{[\text{bold},1]}, v_{[\text{emph},1]}\}, \phi, i)$ . The procedure iterates over the two input nodes. Let us consider first node  $v_{[\text{emph},1]}$ . Below,  $i_1$  is used as a shorthand for  $i_{v_{[\text{emph},1}]}$ . Then, node  $v_{[\text{emph},1]}$  is extended with its children types, in the following way.

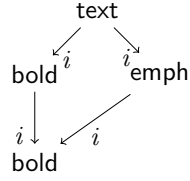


By following line 9, node  $v_{[\text{bold},2]}$  is set as terminal node since it satisfies test  $\phi$ . Let us consider now node  $v_{[\text{bold},1]}$ . Below we assume  $i_2 = i_{v_{[\text{bold},1}]}$ . The child navigation creates an edge between  $v_{[\text{bold},1]}$  and  $v_{[\text{bold},2]}$ . As we already discussed this is not problematic because the chain **text.bold.bold** indeed matches **q**. Notice also that there is a group of chains ending at  $v_{[\text{emph},2]}$  which does not match the query. They will be removed in the next iteration.



The third recursive call of the procedure is the one that terminates the inference. In this iteration step, the input of the procedure is the triplet  $(\{v_{[\text{bold},2]}, v_{[\text{emph},2]}\}, \phi, i)$ . It turns out that both input nodes have some descendant types that satisfy  $\phi$ , but they cannot be further extended, without exceeding the chain space. The maximal length of chains to consider is  $1 + nk = 3$ , because the number of types in the schema is 2 (plus the root) and  $k = 1$ .

Both nodes  $v_{[\text{bold},2]}$  and  $v_{[\text{emph},2]}$  are thus added to the backtracking list. Because node  $v_{[\text{bold},2]}$  is a terminal node for a query return chain it will not be removed. Because



$$\text{RET}(i) = \{ v_{[\text{bold},1]}, v_{[\text{bold},2]} \}$$

node  $v_{[\text{emph},2]}$  is not a terminal node, and it is not shared by any other expression, it is removed, and also all of its incoming edges. The resulting CDAG is the following.

Procedure PARENT computes chains for backward steps of the form  $\mathbf{x}/\text{parent} :: \phi$ . In analogy with the child and descendant axes, this procedure is also used to compute the ancestor axis. The algorithm backtracks on each parent  $v'$  of an input node  $v$ , when the two nodes are connected by an edge  $i'$  belonging to the set  $I$  of ids that can be backtracked (which is part of the input). If the procedure is called to infer chains for a step  $\mathbf{e} = \mathbf{x}/\text{parent} :: \phi$  then  $I$  is  $I = vr^+(\mathbf{e})$ . If the procedure is iterated to infer chains for an ancestor then  $I$  is defined wrt the ancestor step, as discussed next.

To illustrate, consider chain inference over DTD  $\mathbf{d}_3$  and the following query

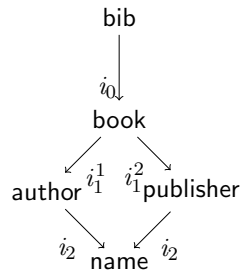
```

for  $\mathbf{x}_0$  in  $\mathbf{x}/\text{book}$ 
for  $\mathbf{x}_1$  in  $\mathbf{x}_0/\text{author}$ ,  $\mathbf{x}_0/\text{publisher}$ 
for  $\mathbf{x}_2$  in  $\mathbf{x}_1/\text{name}$ 
for  $\mathbf{x}_3$  in  $\mathbf{x}_2/\text{parent} :: *$ 
for  $\mathbf{x}_4$  in  $\mathbf{x}_3/\text{parent} :: \text{book}$ 
return  $\mathbf{x}_4$ 

```

where  $\mathbf{x}$  is bound to **bib**.

Below we assume the following coding  $id(\text{//book}) = i_0$  and  $id(\mathbf{x}_0/\text{author}) = i_1^1$  and  $id(\mathbf{x}_0/\text{publisher}) = i_1^2$  and for the remaining steps  $id(\mathbf{x}_i/\text{step}_i) = i_{(i+1)}$ . Chain inference for the first three lines of the query gives the following CDAG.



$$\begin{aligned} \text{RET}(i_0) &= \{ v_{[\text{book},1]} \} \\ \text{RET}(i_1^1) &= \{ v_{[\text{author},2]} \} \\ \text{RET}(i_1^2) &= \{ v_{[\text{publisher},2]} \} \\ \text{RET}(i_2) &= \{ v_{[\text{name},3]} \} \end{aligned}$$

Processing of step  $\mathbf{x}_2/\text{parent} :: *$  is now illustrated. Recall that  $id(\mathbf{x}_2/\text{parent} :: *) = i_3$ .

**Algorithm 4:** Procedure PARENT

---

**Input** : Node-set  $M \subseteq V_{rv}$ ,  $\phi$ , id  $i$ , id set  $I$   
**Output**: Infer chains for **parent** axis, side-effect on **RET()**, **CLEAN()**

```

1 foreach  $v \in M$  do
2    $some\_parent \leftarrow \text{false}$ 
3   foreach  $(v', v, i') \in E$  do
4     if  $v'.type$  satisfies  $\phi$  then
5       if  $i' \in I$  then
6          $RET(i) \leftarrow RET(i) \cup \{v'\}$ 
7          $some\_parent \leftarrow \text{true}$ 
8   if  $!some\_parent$  then
9      $CLEAN(i) \leftarrow CLEAN(i) \cup \{v\}$ 

```

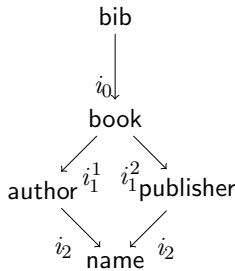
---

The procedure iterates over chains bound to  $\mathbf{x}_2$ , thus over node  $v_{[name,3]}$ . Because  $i_2 \in vr^+(i_3)$ , the procedure backtracks to the edges labeled as  $i_2$ , namely  $(v_{[author,2]}, v_{[name,3]}, i_2)$  and  $(v_{[publisher,2]}, v_{[name,3]}, i_2)$ .

Let us consider first the edge  $(v_{[author,2]}, v_{[name,3]}, i_2)$ . Because  $v_{[author,2]}$  satisfies  $\phi$ , it is set as terminal node for the parent navigation, so  $RET(i_3) = \{v_{[author,2]}\}$ . The second iteration step on edge  $(v_{[publisher,2]}, v_{[name,3]}, i_2)$  is similar.

Processing of step  $\mathbf{x}_3/\text{parent}::book$  is now illustrated. Recall that  $id(\mathbf{x}_3/\text{parent}::book)=i_4$ .

The procedure iterates over chains bound to  $\mathbf{x}_2$ , thus  $v_{[author,2]}$  and  $v_{[publisher,2]}$ . It backtracks on both edges labeled as  $i_1^1$  and  $i_1^2$ , because  $i_1^1, i_1^2 \in vr^+(i_3)$ . Let us consider first edge  $(v_{[book,1]}, v_{[author,2]}, i_1^1)$ . Because  $v_{[book,1]}$  satisfies  $\phi = book$ , it is set as terminal node for the step, so  $RET(i_4) = \{v_{[book,1]}\}$ . The second iteration step on edge  $(v_{[book,1]}, v_{[publisher,2]}, i_1^2)$  is similar, and the resulting CDAG is the following.



$RET(i_0)$	$= \{v_{[book,1]}\}$
$RET(i_1^1)$	$= \{v_{[author,2]}\}$
$RET(i_1^2)$	$= \{v_{[publisher,2]}\}$
$RET(i_2)$	$= \{v_{[name,3]}\}$
$RET(i_3)$	$= \{v_{[author,2]}, v_{[publisher,2]}\}$
$RET(i_4)$	$= \{v_{[book,1]}\}$

Note that the procedure does not add new edges to the CDAG but it just affects **RET()**.

Algorithm 4 is simple, sound and ensure a quite precise analysis at once. It is worth outlining that, in some rare cases, the algorithms can infer unneeded chains, that is a

chain that would not be inferred by the type system presented in Chapter 4. This may happen for expressions that navigate twice the same part of the schema, as illustrated by the following example over the DTD  $d_3$ .

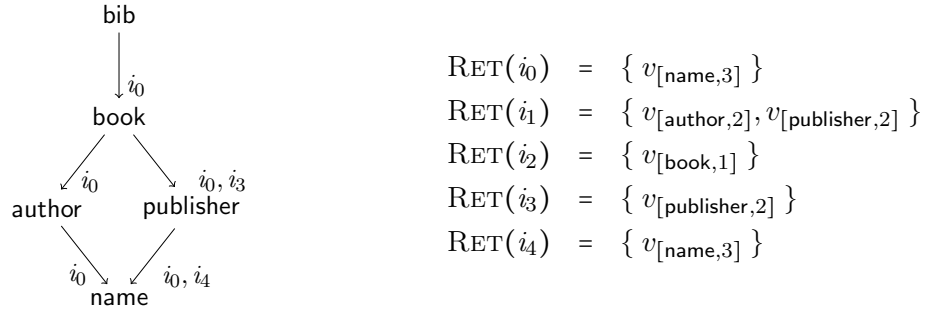
```

let  $x_0$  :=  $x//name$ 
let  $x_1$  :=  $x_0/parent :: *$ 
let  $x_2$  :=  $x_1/parent :: *$ 
let  $x_3$  :=  $x_2/publisher$ 
let  $x_4$  :=  $x_3/name$ 
return  $x_4/parent :: *$ 

```

where we assume  $x$  is bound to **bib**.

Below we assume  $id(x//name) = i_0$  and  $id(x_i/step) = i_{i+1}$ . In this case, the last backward step of the query  $x_4/parent :: *$  produces an unneeded chain. To see this, consider the CDAG obtained by inference over all precedent steps.



When processing step  $x_4/parent :: *$ , the procedure iterates over chains bound to  $x_4$ , thus  $v_{[name,3]}$ . It backtracks on both edges labeled as  $i_0$  and  $i_4$ , because  $i_0, i_4 \in vr^+(i_5)$  and  $i_5$  is the id of  $x_4/parent :: *$ . At this point the procedure sets as terminal nodes both  $v_{[author,2]}$  and  $v_{[publisher,2]}$ , while it should backtrack only on  $v_{[publisher,2]}$ .

In general, for a pair of codes  $i_1, i_2$  such that  $i_1 \in vr^+(i_2)$ , by running Procedure 4 it is possible to visit edges labeled as  $i_1$  before those labeled as  $i_2$ , and this can lead to infer chains as false negatives. However, it is indeed by avoiding to follow the transitive closure of relation  $vr()$  in the proper order, that we can compute backward navigations in polynomial time.

The way backward navigations are used in the above example resembles that of a conditional expression. They allow to verify that node elements labeled with *book* do have descendants labeled with *name* and, if this is the case, descendant nodes of *book* ones labeled with *publisher* are further navigated. Notice that this is not the usual way of expressing such a property in an XQuery program. The above expression could be for instance reformulated by using a conditional expression such as

```

if ( $x//name$ ) then  $x/publisher/name/parent :: *$  else ()

```

**Algorithm 5:** Procedure ANCESTOR-REC

---

**Input** : Node-set  $M \subseteq V_{rv}$ , node-set  $S \subseteq V_{rv}$ ,  $\phi$ , id  $i$   
**Output**: Infer chains for **ancestor** axis, side-effect on **RET()**, **CLEAN()**

---

```

1  $M' \leftarrow \emptyset$ 
2 foreach  $v \in M$  do
3   if some ancestor-type of  $v.type$  satisfies  $\phi$  then
4      $PARENT(\{v\}, *, i_v, vr^+(i))$ 
5      $S \leftarrow S \cup RET(i_v)$ 
6      $SELF(RET(i_v), \phi, i)$ 
7      $M' \leftarrow M' \cup RET(i_v)$ 
8   else
9      $CLEAN(i) \leftarrow CLEAN(i) \cup \{v\}$ 
10  $ANCESTOR-REC((M' \setminus S), S, \phi, i)$ 

```

---

For this expression, chain inference over the DTD  $d_3$  is sound and complete wrt the inference system presented in Chapter 4, because edges inferred for the conditional query  $x//name$  would not be mixed with edges inferred for the rest of the expression.

As illustrated, cases of incompleteness exist but are quite rare in practice and may have an equivalent formulation for which the analysis is complete. The precision of our algorithm have been tested over the wide XMark/XPathMark benchmarks as shown in Chapter 6, without incurring in a situation analogous to the one presented here, and by showing high improvements wrt the state-of-the-art method.

Algorithms implementing ancestor axes are presented in the following. We just describe the algorithm implementing axis **ancestor** since the case for **ancestor-or-self** is similar. Procedure ANCESTOR-REC computes chains for steps of the form  $x//\mathbf{ancestor} :: \phi$ .

The procedure performs a bottom-up visit of the CDAG, by simply iterating the procedure for the parent axis (Line 4). The set of codes that can be visited by procedure PARENT is set as  $vr^+(i)$ , as discussed before. The parents of a node  $v \in M$ , that are recorded in  $RET(i_v)$ , are then filtered by means of procedure SELF to check if they satisfy the ancestor navigation.

Furthermore, the procedure records a set of visited nodes  $S$ , so as to perform only once the parent navigation to each node of the  $G$ . This optimization makes the procedure running in linear time in the size of the CDAG. A node that has no parent satisfying  $\phi$  represent a dangling chain that may be removed.

We illustrate the procedure with an example. Consider the DTD  $d_{11}$  defined as follows

```

        bib ← book*
        book ← title, author, publisher
        author, publisher ← name
        name ← first, last, String
        title, first, last ← String

```

and the query sequence

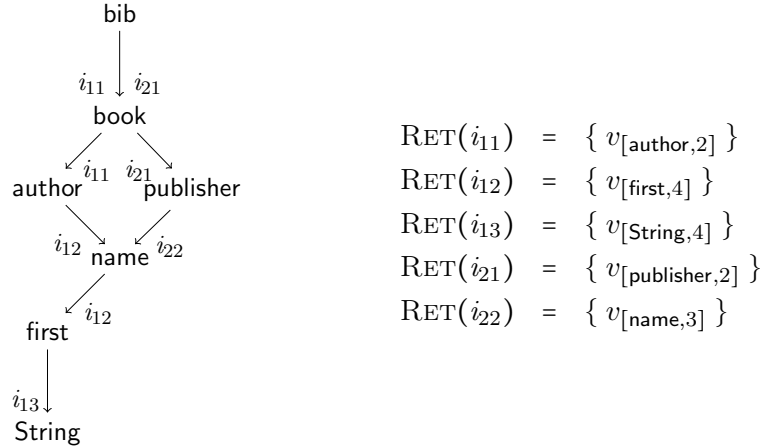
```

q =  for x1 in x//author      for y1 in y//author
      for x2 in x1//first ,   for y2 in y1/name
      return x2/text()        return y2/ancestor :: *

```

with both  $x$  and  $y$  bound to **bib**.

Below, we assume that  $id(x//author) = i_{11}$ ,  $id(x_1//first) = i_{12}$ ,  $id(x_2/text()) = i_{13}$ ,  $id(y//author) = i_{21}$ ,  $id(y_1/name) = i_{22}$  and  $id(y_2/ancestor :: book) = i_{23}$ . Provided that both  $x$  and  $y$  are bound to **bib**, chain inference for  $q$ , yields the following CDAG.



Now, inferring chains for  $y_2/ancestor :: *$  consist of visiting all ancestors of nodes of  $v_{[name,3]}$  on  $G_{rv}$  with a label in  $vr^+(i_{22})$ , that means  $i_{21}$  and  $i_{22}$ . The set of terminal nodes for  $y_2/ancestor :: *$  is thus  $RET(i_{22}) = \{ v_{[publisher,2]}, v_{[book,1]}, v_{[bib,0]} \}$ . Notice that thanks to the edge labeling we avoid to backtrack also to node  $v_{[author,2]}$ .

Procedure FOLLOWING-SIBLING computes chains for the preceding-sibling navigational axis, by means of the PARENT and CHILD procedures. Notice that the other procedure for sibling axis PRECEDING-SIBLING just differ from this one by the test made in Line 5.

**Algorithm 6:** Procedure FOLLOWING-SIBLING

---

**Input** : Node set  $M \subseteq V_{rv}$ ,  $\phi$ , id  $i$   
**Output**: Infer chains for **preceding-sibling**, **following-sibling** axes,  
side-effect on  $RET()$ ,  $CLEAN()$ ,  $V$  and  $E$

---

```

1 foreach  $v \in M$  do
2    $new\_sibling \leftarrow \text{false}$ 
3    $PARENT(\{v\}, *, i_v, vr^+(i))$ 
4   foreach  $v' \in RET(i_v)$  do
5     foreach a following-sibling type of  $v.type$  that satisfies  $\phi$  do
6        $CHILD(\{v\}, a, i)$ 
7        $new\_sibling \leftarrow \text{true}$ 
8   if  $!new\_sibling$  then
9      $CLEAN(i) \leftarrow CLEAN(i) \cup \{v\}$ 

```

---

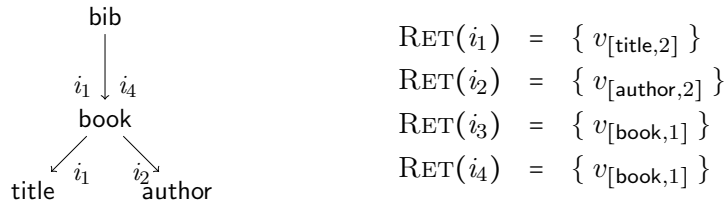
To illustrate the procedure, consider chain inference for the following query  $q$  over DTD  $d_3$  with  $x \mapsto \text{bib}$ .

```

q =  for  $x_1$  in  $x_0//title$ 
      for  $x_2$  in  $x_1//following-sibling::author$ 
      for  $x_3$  in  $x_2//parent::*$ 
      return  $x_3//following-sibling::book$ 

```

Below, we assume that  $id(x_i/step) = i_{i+1}$ . Chain inference for  $q$  gives the following CDAG.



Chain inference for step  $x_2//following-sibling::author$  is performed by backtracking on edges labeled as  $i_1 \in vr^+(i_3)$  created by the previous descendant navigation, and then by adding node  $v_{[author,2]}$  in the CDAG. After this, the parent navigation further backtracks to node  $v_{[book,1]}$ . Chain inference for step  $x_3//following-sibling::book$  is performed again by backtracking on edges labeled with  $i_1 \in vr^+(i_4)$ .

Notice that the backward navigation in Line 3 makes this procedure inherently sound but not complete. Also, in the case where the schema enforces that a type  $a$  is both a following-sibling and a preceding-sibling of a type  $b$  (e.g.,  $r \leftarrow a, b^+, a^?$ ) we have that the result of chain inference for both the sibling axis is the same, if we start from some CDAG node labeled with  $b$ . Again this approximation is sound, but it may be not complete, because expressions like  $doc//b/following-sibling::a$  and



**Algorithm 7:** Procedure INFER-QUERY

---

**Input** : Query  $q$ ,  $G_{rv}$  and  $G_e$   
**Output**: Infer chains for  $q$ , side effect on  $G_{rv}$ ,  $G_e$  and  $RET()$

```

1  $i \leftarrow id(q)$ 
2 switch ( $q$ ) do
3   case  $q = x/axis :: \phi$ 
4      $\lfloor$  call procedure computing  $x/axis :: \phi$ 
5   case  $q = q_1, q_2$  or  $q = \text{if } (q_0) \text{ then } q_1 \text{ else } q_2$ 
6     foreach  $j \leq 2$  do
7        $\lfloor$  INFER-QUERY( $q_j$ )
8        $\lfloor$   $i_j \leftarrow id(q_j)$ 
9      $\lfloor$   $RET(i) \leftarrow RET(i_1) \cup RET(i_2)$ 
10  case  $q = \text{for/let } x \text{ in/} := q_1 \text{ return } q_2$ 
11     $\lfloor$  INFER-QUERY( $q_1$ )
12     $\lfloor$  INFER-QUERY( $q_2$ )
13     $\lfloor$   $RET(i) \leftarrow RET(i_2)$ 
14  case  $q = \text{"txt"}$ 
15     $\lfloor$  BUILD(String,  $\_$ ,  $i$ )
16  case  $q = \langle a \rangle q' \langle /a \rangle$ 
17     $\lfloor$  INFER-QUERY( $q'$ )
18     $\lfloor$   $i' \leftarrow id(q')$ 
19     $\lfloor$  BUILD( $a$ ,  $RET(i')$ ,  $i$ )
20   $\lfloor$  CLEAN-GRAPH( $i$ )

```

---

`delete doc//a/preceding-sibling :: b` are deemed dependent, while they are not. However, this does not depend on how chain inference is implemented, but on the fact that chains abstract away from structural constraints concerning the horizontal order of data.

Procedure INFER-QUERY infers chains for whole query expressions. The procedure does essentially two things. It recursively infer chains for each subexpression of the query, and it computes used, return, and element chains for the whole query.

The base case of procedure INFER-QUERY consist of calling the proper axis navigation for a step of the form  $x/axis :: \phi$ . Recall that for parent axis, the procedure also needs to pass  $vr^+(x)$  as input. Return chains for  $q$  are set as return chains for the proper subexpression of  $q$ . Note that we do not need any particular extra structure for used chains of  $q$ , since they result to be marked as return chains of each  $q'$  subexpression of  $q$ .

We assume that in the first call of the procedure all CDAGs  $G_{rv}$  and  $G_e$  are initialized, and the root node is added in  $G_{rv}$ . Finally, after chain inference for each subexpression

**Algorithm 8:** Procedure BUILD

---

**Input** : Type  $a \in \Sigma_S$ , node-set  $M \subseteq V_e$ , expression id  $i$   
**Output**: Infer element chains for  $q$ , side effect on  $V_e$  and  $E_e$

---

```

1  $v_a \leftarrow (a, 0)$ 
2  $V_e \leftarrow V_e \cup \{ v_{[a,0]} \}$ 
3 foreach  $v \in M$  do
4    $v' \leftarrow \text{COPY}(v, 1, i)$ 
5    $E_e \leftarrow E_e \cup (v_a, v', i)$ 
6   if  $v \in V_{rv}$  then
7      $\text{ELTRET}(i) \leftarrow \text{ELTRET}(i) \cup \{ v' \}$ 
8  $\text{RET}(i) \leftarrow \text{RET}(i) \cup \{ v_a \}$ 

```

---

the CDAGs is done, dangling chains are removed by procedure CLEAN-GRAPH. This last operation correspond to the chain filtering made by the rule for iteration presented in the type system of Chapter 4.

Element chains are inferred by procedure BUILD. This procedure extends the CDAG of element chains  $G_e$ . The procedure is called either for string query “txt”, or for an expression of the form  $\langle a \rangle q \langle /a \rangle$ . Because we assume that constructed elements are never navigated (see Chapter 2), element chains in the  $G_e$  are simply stored as trees. This comports no relevant overhead.

The procedure first creates a fresh node representing the root of the constructed structure. Then, all nodes in  $M$  are copied (Line 4) and connected with  $v_a$  (Line 5).

Function  $\text{COPY}(v)$ , despite not being explicitly modeled here, does three things. First, it distinguishes if the input node belongs to  $G_{rv}$  or to  $G_e$ . In the first case, it simply copies the node. In the second case, it copies the node and all the structure below. Second, it sets as 1 the height of the copied nodes, and accordingly of all copied descendants. Third, it marks edges of the copied structure with  $i$ .

Finally, according to the rule for element construction presented in Chapter 4, to record the fact that that a node  $v'$  copy of  $v$  should be meant as a whole subtree, we record  $v'$  in  $\text{ELTRET}(i)$ . This has to be done only if the original node  $v$  belongs to the CDAG  $G_{rv}$ .

**Algorithm 9:** Procedure INFER-UPDATE

---

**Input** : Update  $u$ , CDAG  $G_U$   
**Output**: Infer chains for  $u$ , side effect on  $G_U$

```

1  $i \leftarrow id(u)$ 
2 switch ( $u$ ) do
3   case  $u = \text{delete } q_0$ 
4     INFER-QUERY( $q_0$ )
5      $i_0 \leftarrow id(q_0)$ 
6      $G_U \leftarrow G_{rv}$ 
7      $UPD(i) \leftarrow UPD(i_0)$ 
8      $DESCUPD(i) \leftarrow UPD(i_0)$ 
9   case  $u = \text{rename } q_0 \text{ as } a$ 
10    INFER-QUERY( $q_0$ )
11     $i_0 \leftarrow id(q_0)$ 
12     $M' \leftarrow RET(i_0)$ 
13     $SIBLING(RET(i_0), a, i_0)$ 
14     $RET(i_0) \leftarrow RET(i_0) \cup M'$ 
15     $G_U \leftarrow G_{rv}$ 
16     $UPD(i) \leftarrow UPD(i_0)$ 
17     $DESCUPD(i) \leftarrow UPD(i_0)$ 
18  case  $u = \text{insert } q \text{ pos } q_0 \text{ or } u = \text{replace } q_0 \text{ with } q$ 
19    INFER-QUERY( $q$ )
20    INFER-QUERY( $q_0$ )
21     $i_0 \leftarrow id(q_0)$ 
22     $i_1 \leftarrow id(q)$ 
23     $G_U \leftarrow G_{rv}$ 
24    if  $\text{pos} \in \{\text{into}, \text{as first}, \text{as last}\}$  then
25      BUILD-UPDATE( $RET(i_1), RET(i_0), i_0, G_e$ )
26    else
27       $M' \leftarrow RET(i_0)$ 
28       $RET(i_0) \leftarrow \emptyset$ 
29       $I \leftarrow vr^+(ret(q_0))$ 
30       $PARENT(M', *, i_0, I)$ 
31      BUILD-UPDATE( $RET(i_0), RET(e_1), i_0, G_e$ )
32     $UPD(i) \leftarrow UPD(i_0)$ 
33    if  $u$  is a replace then
34       $DESCUPD(i) \leftarrow UPD(i_0)$ 

```

---

Procedure **INFER-UPDATE** infers chains from an update expression  $u$ . As updates are defined in terms of queries, this mainly consists of calling procedure **INFER-QUERY** for the update subexpressions.

For delete expressions, update chain inference is the same as query chain inference for the target expression  $q_0$ . The update CDAG  $G_U$  is set as the query CDAG  $G_{rv}$ . Here, CDAG nodes in the prefix  $c'$  of an update chain of the form  $c : c'$  are recorded by the inverted index  $UPD()$ . Furthermore it makes use of another index  $DESCUPD()$  that records the fact that the whole subtrees rooted at  $RET(i_0)$  has to be considered as impacted by the update. This permits to avoid to infer all chains corresponding to descendants of nodes in  $RET(i_0)$ .

Chain inference is similar for rename expressions, with the difference that chains capturing renamed node are inferred by means of a further navigation on the siblings of the return nodes for  $q_0$ . Procedure **SIBLING** performs a following-sibling and a preceding-sibling navigation on the terminal nodes of  $q_0$ .

Chain inference for insert and replace expressions are similar. They both combine chains inferred from the source and for the target expression. The case where an insert update uses either position **after** or **before** is the same as for the replace command. In both cases it is necessary to perform a parent step, on the result locations of  $q_0$ . In order to allow a backward navigation, the set of edges that can be visited has to be computed. This is done by first retrieving return subexpressions of  $q_0$  by means of function  $ret()$  and  $vr^+$ , as already outlined.

Chains recording the structure of inserted subtrees are inferred by procedure **BUILD-UPDATE**. The procedure puts an edge between terminal nodes of the target and the source query expressions. Nodes belonging to the source query are copied, setting also the level and the label of edges.

---

**Algorithm 10:** Procedure **BUILD-UPDATE**


---

**Input** : Target node-set  $T \subseteq V_U$ , source node-set  $S \subseteq V_e$ , expression id  $i$ , CDAG  $G_e$

**Output**: Infer element chains for  $q$ , side effect on  $V_U$  and  $E_U$

```

1 foreach  $v \in T$  do
2   foreach  $v' \in S$  do
3      $v'' \leftarrow \text{COPY}(v', v.\text{level} + 1, i)$ 
4      $V_U \leftarrow V_U \cup \{v''\}$ 
5      $E_U \leftarrow E_U \cup (v', v'', i)$ 

```

---

**Algorithm 11:** Function CHECKINDEPENDENCE

---

**Input** : CDAG for query  $G_{rv}$  and CDAG for update  $G_U$ , node sets  $M_r, M_v \subseteq V_{rv}$   
and  $M_U, M'_U \subseteq V_U$

**Output:** Check for a common chain between two graphs

```

1 foreach  $h = 0.. \min(h_{G_{rv}}, h_{G_U})$  do
2   foreach  $(v, v', i) \in E_{rv}$  and  $(v, v', j) \in E_U$  do
3     if  $((i', j') \in I^{(v'', v)} \text{ and } i' \in vr^*(i) \text{ and } j' \in vr^*(j)) \text{ or } h = 0$  then
4        $I^{(v, v')} \leftarrow I^{(v, v')} \cup \{ (i, j) \}$ 
5       if  $v \in M_r \cup M_v$  and  $v \in M_U$  then
6         return false
7       if  $v \in M_r$  then
8         if a descendant of  $v \in G_U$  also belongs to  $M_U$  then
9           return false
10      if  $v \in M_U$  then
11        if a descendant of  $v \in G_{rv}$  also belongs to  $M_r \cup M_v$  (return-element)
12          then
13            return false
13 return true

```

---

The last algorithm we present, is the one checking if query and update CDAGs have a conflicting chain. The procedure takes as input the set of terminal nodes denoting chains for used and return query chains, as well as for update nodes and for update and all descendant nodes. In particular,  $M_r \subseteq V_{rv}$  is set of nodes that are terminal of query return chains,  $M_v \subseteq V_{rv}$  is set of nodes that are terminal of query used chains,  $M_U \subseteq V_U$  is set of nodes that are in the suffix of an update chain and  $M'_U \subseteq V_U$  is set of nodes that are in the suffix of an update chain and whose whole subtree is impacted by the update.

The procedure iterates over all the levels of the CDAGs and inductively computes a product graph, where edges between nodes are sets of pairs of ids of query and update expressions that simulate the same chain.

At level  $h$ , a pair of codes  $(i, j)$  is added on the edge connecting a node  $v$  with a node  $v'$  (Line 4) only if it is the continuation of some chains computed at level  $h - 1$  (Line 3). Formally, this happens when for some node  $v''$  parent of  $v$  there is a pair of codes  $(i', j')$  such that  $i' \in vr^*(i)$  and  $j' \in vr^*(j)$ . Here  $vr^*(i)$  denotes the reflexive and transitive closure of  $vr(i)$ . This test is trivial for  $h = 0$  since the root has no incoming edge. The set  $I^{(v, v')}$  denotes all pair of edges between  $v$  and  $v'$  in the product graph.

In order to detect conflicts between query and update chains, and hence to implement Definition 4.4.1, we perform the following three tests. If we reached a node  $v$  such that

$v$  is both the terminal node of a used or return chain in  $G_{rv}$ , and an update node of  $G_U$ , then we raise a dependence. Otherwise, if  $v$  is only the terminal of a return chain, we visit  $G_U$  to search for a descendant of  $v$  that also belongs to  $M_U$ . If this node exists, we raise a dependence. Otherwise, if  $v$  is either an update node of  $M_U$  or an update node with all descendants of  $M'_U$  then we visit  $G_{rv}$  and we search for a descendant of  $v$  that also belongs to  $M_r \cup M_v$ . Of course, this is done by properly following edges corresponding to related subexpressions.

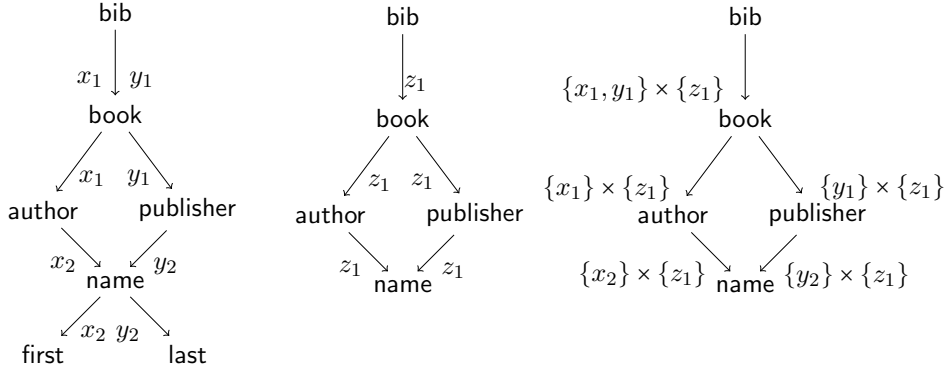
To illustrate, consider DTD  $d_{11}$  and the following query and updates.

$$q = \text{for } x' \text{ in } x//author \quad , \quad \text{for } y' \text{ in } y//publisher \\ \text{return } x'//first \quad \quad \text{return } y'//last$$

$$u_1 = \text{delete } z//name \quad \quad u_2 = \text{for } z' \text{ in } z//publisher \\ \text{return delete } z'//first$$

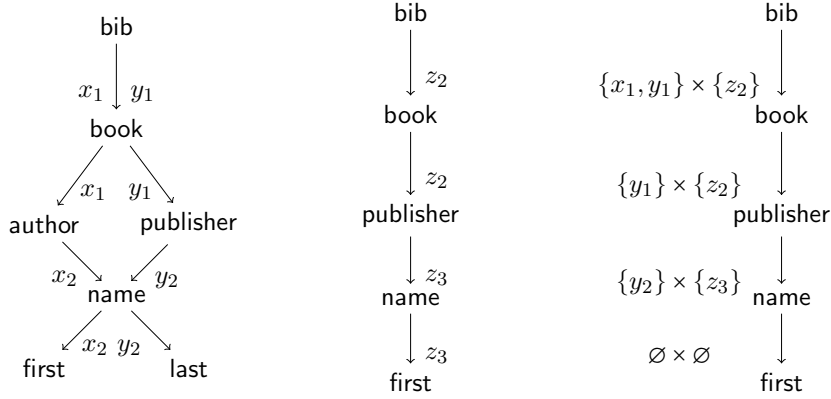
We run Algorithm 11 and show the result of the product graph for the pairs  $q-u_1$  and  $q-u_2$ , respectively. Below, we assume that  $x, y$  and  $z$  are all bound to **bib**. Furthermore,  $id(x//author) = x_1$ ,  $id(x'//first) = x_2$ ,  $id(y//publisher) = y_1$ ,  $id(y'//last) = y_2$ ,  $id(z//name) = z_1$ ,  $id(z//publisher) = z_2$  and  $id(z'//first) = z_3$ .

The product graph for  $q-u_1$  is the following.



In this case there are two conflicting chains because the query and the update access and impacts nodes labeled with *name*. After node  $v_{[name,3]} \in G_U$  has been recognized as a node in the prefix of an update chain, the procedure start searching for a descendant of node  $v_{[name,3]} \in G_{rv}$  that is also the terminal of a used or return chain, by following ids  $x_2, y_2$ . The procedure then reaches both  $v_{[first,4]}$  and  $v_{[last,4]}$  and hence a conflict is correctly raised by the procedure.

The product graph for  $\mathbf{q}\text{-}\mathbf{u}_2$  is the following.



In this case there is no conflicting chain between the two expressions. Notice first that chains for the left subexpression of sequence  $\mathbf{q}$  does not simulate any chain inferred for  $\mathbf{u}$ . This is reflected by the fact that in the product graph there is no edge with id  $x_1$  after node  $v_{[\text{book},1]}$ . Similarly, also the right subexpression of sequence  $\mathbf{q}$  does not simulated after node  $v_{[\text{name},3]}$  in the CDAG, because in the product graph there is no edge connecting  $v_{[\text{name},3]}$  with node  $v_{[\text{first},4]}$ .

## 5.4 Complexity of Chain Inference with CDAGs

In this section we present some results on the complexity of algorithms defined in the previous section. We mainly distinguish between recursive and non-recursive schemas, and outline how complexity can be lowered for some interesting classes of queries and updates.

**Theorem 5.4.1.** *Let  $\mathbf{d}$  be a non recursive DTD with  $n$  type definitions, and  $\mathbf{e}$  a query or update expression of size  $m$ . By using Algorithms 1-10 of Section 5.3 chain inference can be done in  $O(mn^3)$  space and  $O(mn^5)$  time.*

*Proof.* When  $\mathbf{d}$  is not recursive, the  $k$  value stop being determinant for the analysis since no label repeats twice in any schema chain. We first consider the case where  $\mathbf{e}$  is not an insert or replace update. By Lemma 5.2.1, we know that the worst-case size of the CDAG inferred for  $\mathbf{e}$  is  $O(n^3)$ . This, together with the fact that each edge can be labeled with a code, and that the number of codes is upper bounded by the size of the query entails that the CDAG needs at most  $O(mn^3)$  space. Moreover, by using Algorithms 1-11 of Section 5.3 Chain inference for any navigational step, can be computed in the worst case with a traversal of the CDAG, hence in  $O(n^3)$ , for each input nodes. The number of input nodes is upper bounded by the number of nodes in the CDAG. By Lemma 5.2.1, there are at most  $O(n^2)$  nodes in the CDAG. For the whole expression it follows that inference time is in  $O(mn^5)$ . If  $\mathbf{e}$  is an insert or replace update, then by

making a cartesian product of chains for the target and source update expressions with Procedure BUILD UPDATE, we have an overhead which is at most quadratic in the number or terminal nodes of each expression. However, since we assumed that updates preserve the schema, this cannot result in more than  $O(n^2)$  pairs of CDAG nodes corresponding to types satisfying the insert or replace expression. This overhead does not affect the asymptotical complexity.  $\square$

**Corollary 5.4.2.** *Let  $\mathbf{d}$  be a non recursive DTD with  $n$  type definitions, and  $\mathbf{e}$  a query or update expression of size  $m$ . By using Algorithms 1-10 of Section 5.3, the followings holds.*

1. *If  $\mathbf{e}$  does not use step filter `node()` then chain inference can be done in  $O(mn^4)$  time;*
2. *If  $\mathbf{e}$  does not use recursive downward an backward axis then chain inference can be done in  $O(mn^2)$  time;*
3. *If  $\mathbf{e}$  does not use both step filter `node()` and recursive axis then chain inference can be done in  $O(mn)$ ;*
4. *If each navigational step of  $\mathbf{e}$  returns a fixed number of chains, then chain inference can be done in  $O(mn^3)$ .*

*Proof.* Analogous to the proof of Theorem 5.4.1.  $\square$

The above restriction are often met in practice, and in particular by expressions used in our testbed (when the recursive component of the XMark schema is not visited at all by the expression).

**Theorem 5.4.3.** *Let  $\mathbf{d}$  be a recursive DTD with  $n + 1$  type definitions, and  $\mathbf{e}$  a query or update expression of size  $m$ . By using Algorithms 1-10 of Section 5.3 chain inference can be done in  $O(mkn^3)$  space and  $O(mk^2n^5)$  time.*

*Proof.* We first consider the case where  $\mathbf{e}$  is not an insert or replace update. By Lemma 5.2.2 we know that the worst-case size of the CDAG inferred for  $k$ -chains of  $\mathbf{e}$  is  $O(kn^3)$ . Because each edge can be labeled with any code, the worst-case space occupancy for the CDAG is in  $O(mkn^3)$ . By using Algorithms 1-10 of Section 5.3, chain inference for any navigational step can be computed in the worst case with a traversal of the CDAG, hence in  $O(kn^3)$ , for each input nodes. The number of input nodes is upper bounded by the number of nodes in the CDAG. By Lemma 5.2.1, there are at most  $O(kn^2)$  nodes in the CDAG. For the whole expression then inference time is in  $O(mk^2n^5)$ . If  $\mathbf{e}$  is an insert or replace update, Procedure BUILD UPDATE introduce an overhead which is quadratic in the the size of the number or terminal nodes of each expression. At most  $O(kn^2)$  nodes are inferred for each expression. However, since we assumed that updates preserve the schema, this cannot result in more than  $O(mkn^3)$  pairs of CDAG nodes corresponding to types satisfying the insert or replace expression. Hence asymptotic complexity is not affected by this operation.  $\square$



Further restrictions on the the expression would lower again upper bounds.

**Corollary 5.4.4.** *Let  $\mathbf{d}$  be a recursive DTD with  $n + 1$  type definitions, and  $\mathbf{e}$  a query or update expression of size  $m$ . By using Algorithms 1-10 of Section 5.3, the followings holds.*

1. *If  $\mathbf{e}$  does not use step filter `node()` then chain inference can be done in  $O(mk^2n^4)$  time;*
2. *If  $\mathbf{e}$  does not use recursive downward an backward axis then chain inference can be done in  $O(mkn^3)$  time;*
3. *If  $\mathbf{e}$  does not use both step filter `node()` and recursive axis then chain inference can be done in  $O(mkn^2)$  time;*
4. *If each navigational step of  $\mathbf{e}$  returns a fixed number of chains, then chain inference can be done in  $O(mkn^3)$  time.*

*Proof.* Analogous to the proof of Theorem 5.4.1. □

As suggested by the statement complexity of chain inference can be lowered in many cases of practical relevance. In particular, if we assume that during chain inference each XPath step can have a fixed number of CDAG nodes as input, time complexity goes down to  $O(mkn^3)$ . The size of the input is likely to be close to 1 for most XPath steps used in practice. This holds in particular for XMark and XPathMark expressions. Another fact observable from such expressions is that they employ a small number of recursive navigations, thus further lowering the complexity.

**Theorem 5.4.5.** *Checking independence between two CDAGs inferred for a query  $\mathbf{q}$  and an update  $\mathbf{u}$  over a DTD  $\mathbf{d}$ , can be done by Function 11 in  $O(km^2n^3)$  time, where  $k = k_{\mathbf{q}} + k_{\mathbf{u}}$  and  $m = \min\{|\mathbf{q}|, |\mathbf{u}|\}$ .*

*Proof.* Function 11 iterates on the levels of the CDAGs. Each CDAG has at most  $nk$  levels, and there are at most  $n^2$  edges between the two levels. Moreover, each edge is labeled with at most  $m$  codes. Each code for an edge in the query CDAG may be paired with each code for the same edge in the update CDAG, and this can be done in  $O(m^2)$  per edge. Putting everything together we have  $O(km^2n^3)$ . □

## Conclusions

In this section we provided algorithms for efficiently implementing our independence analysis based on schema chains. We showed that chain inference is challenging from the computational point of view, because schema constraints can lead to infer an exponential number of chains. To overcome this limitation, we proposed a conservative implementation of the chain analysis that runs in polynomial space and time. In the next chapter experiments will show that the proposed CDAG algorithm is highly precise.

## Chapter 6

# Experiments

In this chapter we report on experiments made for validating the efficiency of the chain-based independence analysis. We describe the benchmarks we used as well as results obtained witnessing precision and efficiency of our technique.

### 6.1 Benchmarks and Experimental Settings

Our technique has been fully implemented in Java 6. We performed extensive experiments by using our Java implementation, in order to measure *i*) efficiency, *ii*) precision and *iii*) scalability of our static analysis. We used two different benchmarks: a first one based on XMark /XPathMark, and a second one, dubbed R-benchmark, we specifically designed to measure scalability, that are now described.

**Independence benchmark** We used a superset of the independence benchmark adopted by Benedikt and Cheney in [BC09a]. Our benchmark is composed of a set of 36 queries  $v_i$  and a set of 31 updates  $u_i$ . A query belongs to either the XMark query set  $q_1$ – $q_{20}$  [SWK<sup>+</sup>02], or to the XPathMark query set  $A1$ – $A8/B1$ – $B8$  [Fra05]. Queries in the group  $Ai$  of XPathMark only use downward axes, whereas  $Bi$  queries use upward and horizontal axes as well. Concerning updates, a first set corresponds to those used in [BC09a]. These are derived from the XPathMark query set  $A1$ – $A8/B1$ – $B8$  and are of the form  $UAi=\text{delete } Ai$  or  $UBi=\text{delete } Bi$ .

It is worth observing that not all of the delete-updates preserve the schema, but as already seen (Section 4) this is safe because if a chain does not belong to the schema then it is not inferred for these updates. Instead, remaining updates are defined so as to preserve the schema. We added a set of 15 updates formed by insert expressions  $UI1$ – $UI5$ , rename expressions  $UN1$ – $UN5$ , and replace expressions  $UP1$ – $UP5$ . These updates have been defined so as to cover all different types of nodes in XMark documents, and in particular those parts defined by mutually recursive types. It is worth remarking that, even if not all of the delete-updates of the testbed preserve the schema (see  $UA4$ ,  $UA5$ ,  $UA6$ ,  $UA7$ ,  $UA8$ ,  $UB1$ ,  $UB5$ ,  $UB6$ ,  $UB7$ ,  $UB8$ ), the correctness of our technique is still

ensured, since no new chain is created by these expressions. As outlined before, our technique is just unaware of new chains built by breaking schema constraints. In light of this, insert, rename and replace update expressions have been chosen in order to preserve document validity.

Before performing the tests, XMark and XPathMark expressions have been opportunely rewritten into expressions belonging to the XQuery fragment we consider (Section 3), as done in [BC10]. The rewriting essentially consists of: putting predicate conditions in disjunctive form, removing attribute use, and extracting paths from functions calls and arithmetic expressions. Clearly, the rewriting is such that a query and an update are independent if the rewritten query and update are. Updates extending XMark/XPathMark benchmark are reported in Table 6.2.

This benchmark has been used in order to measure precision and efficiency of our technique.

**R-benchmark** This benchmark is formed by schemas and expressions with a massive use of recursion. The benchmark is constituted of five saturated-recursive schemas with 1, 3, 5, 10 and 20 recursive types, denoted by `d1`, `d3`, `d5`, `d10`, and `d20` respectively. In general, schema `dn` is defined as

$$\begin{aligned} a_0 &\leftarrow (a_1 | \dots | a_n)^* \\ a_1 &\leftarrow (a_1 | \dots | a_n)^* \\ &\vdots \\ a_n &\leftarrow (a_1 | \dots | a_n)^* \end{aligned}$$

The benchmark also features a series of XPath expressions consisting on 1, 5 and 10 consecutive `descendant-or-self :: node()` steps, denoted by `e1`, `e5` and `e10` respectively. In general, expression `en` is defined as

$$\text{descendant-or-self :: node}()_1 / \dots / \text{descendant-or-self :: node}()_n$$

This benchmark is designed for understanding the impact of recursion in the performances of our analysis, and draw conclusions on its scalability.

## Configuration

We ran all tests on a desktop 4-core Intel Xeon 2.13 GHz machine with 8 GB RAM (the JVM was given 2 GB) running Linux. To avoid perturbations coming from system activity, we ran each experiment ten times, discarded the best and the worst performance, and computed the average of the remaining times.

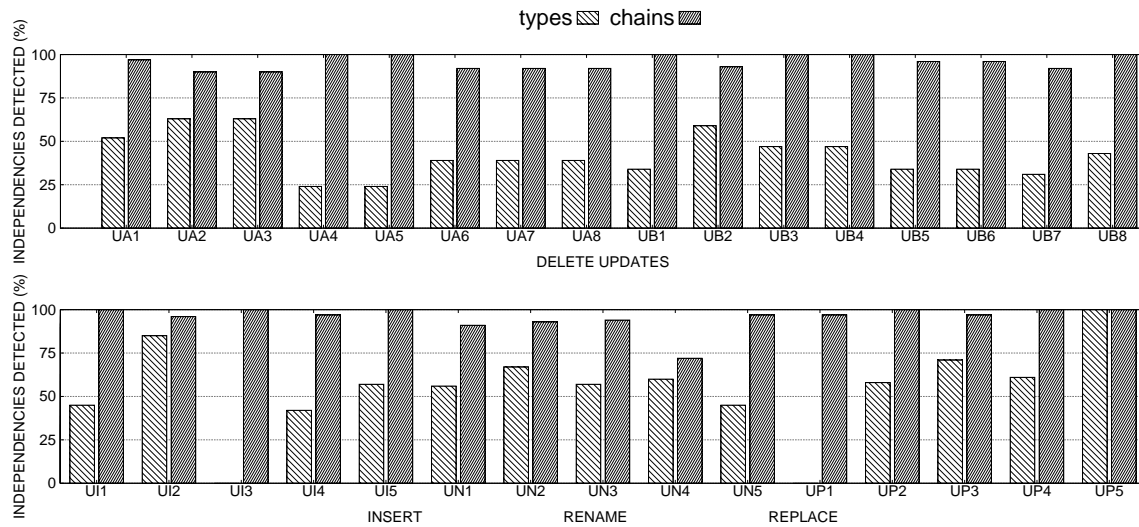


Figure 6.1: Static Analysis Precision: Types vs Chains

## 6.2 Experimental Results

### Precision on XMark

Independence (Definition 2.5.2) is undecidable in general [BC09a], so for the purpose of measuring precision, for each update  $u_i$  we manually determined independent pairs  $(u_i, q_j)$ , and reported in Table 6.2 of the appendix. (note that for most pairs in the considered testbed independence is evident, so this process is much less time consuming than one may guess). We then express precision as the percentage of independent pairs that are deemed independent by our static analysis too. To estimate improvements wrt the alternative schema-based technique [BC09a] we computed the same percentages for that technique by using the public tool [Che09].

Results are reported in Figure 6.2.b. Our chain-based analysis turned out to be precise. Percentages goes from 72% to 100%, while the average precision is 96%. Also, Figure 6.2.b shows that the analysis proposed in [BC09a] (that has an average detection of 49%) is always outperformed in terms of precision by our static analysis, and in some cases improvements are huge. This happens in particular for updates UB1, UB5, UB6, UB8 (employing backward and horizontal axes). For these updates, the over-approximation made by type rules in [BC09a] entails a high number of false negatives. Our chain based inference instead is so precise to avoid most of these false-negatives. In general, improvements in terms of precision go from 8% (UN4) to 96% (UP1), and the average gain is 46%. In particular, precision of our analysis remains high in the presence of views using upward and horizontal axes (XPathMark queries in the group B). These queries are likely to be among the most expensive ones to re-evaluate after document updating.

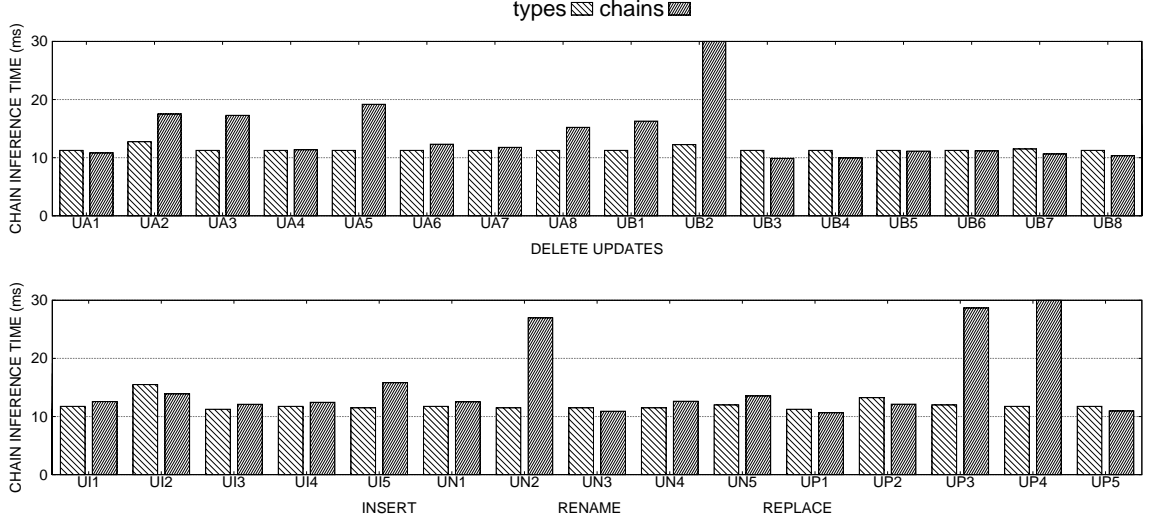


Figure 6.2: Static Analysis Time: Types vs Chains

As expected, we found that in some cases independence is not captured by our static analysis. This happens for pairs that make use of particular node position, identity or value conditions whose semantics is not precisely captured by our static analysis. Also, our system is not able to understand that the renaming performed by UN4 does not change any document tag, thus excluding independence in many cases (the same holds for [BC09a]).

### Inference time on XMark

We measured the time needed by the static analysis to detect independence of each update wrt the whole set of XMark queries. The XMark schema is particularly suitable for testing the performances of our technique since the type dependency graph of this schema contains 5 mutually recursive types that form two cliques of size 2 and 3 respectively. We recall that the execution cost depends on the three parameters  $|d|$ ,  $|e|$  and  $k$ . In this testbed we have  $|d| = 76$ , and  $|e| \leq 20$ , while multiplicity values  $k$  range from 2 to 6. As observed in Section 5.4, in many cases chain inference can be substantially lowered.

Time values include the time for CDAGs inference and comparison, for each pair of expressions. Results are collected in Figure 6.2. It shows that the analysis is quite fast: in the worst case the analysis is performed in less than 40 ms for the whole set of queries, while the average cost is around 15 ms. According to complexity results of Section 5.4, inference time is influenced by i) the  $k$  values needed by a query-update pair and ii) the number of recursive types of the schema effectively unfolded. We see small changes in inference time values according to the  $k$  value (e.g., the pair UB1-UB2). Yet, two expressions having the same  $k$  value may have different time costs for chain inference,

depending on the effective number of recursive types unfolded by the analysis (e.g., the pair UI3-UP3).

Running times obtained from the available OCaml implementation [Che09] of the analysis presented in [BC09a] are rather close to ours: the average time for analyzing an update vs all of the queries is around 10 ms. It is worth observing, that inference time for [BC09a] has no sensible oscillations, while in our case inference time depends on  $k$ , hence on the query and update expressions. The analysis presented in [BC09a] has worst case time complexity  $O((|d|^2+|q|)^2+|u|)$ , and thus is expected to be faster than our analysis in the presence of recursive schemas. Nevertheless, as shown shortly, our running times remain low enough to ensure high time savings in queries maintenance, even when queries are defined on relatively small documents.

### Benefits of the Independence Analysis

Besides experimenting on our prototype, we verified that the precision and the efficiency of our method directly translate in time-savings when our static analysis is used to optimize view maintenance. Along the line of [BC09a], we simulated a view-maintenance scenario, where views are expressed as query, and on updates views are simply refreshed. We measured time savings obtained by avoiding the refreshing of views (queries) which our analysis deem as independent of an update.

We used three XQuery engines: Saxon 9.2EE, BaseX 7.0.1 and QizX 4.4<sup>1</sup>. This, in order to make a faithful comparison of the three systems. We considered a 1MB XMark document and we scaled to 10MB and 100MB, in order to measure time savings in larger XML corpora. Our test results only take into account query answering time, dried of import time and indexing time. This, in order to make a faithful comparison of the three systems. Saxon is a main-memory engine, that uses data summaries for rapidly access data. Differently, BaseX and QizX are persistent system that exploit indexing. We settled the common *full-text* indexes for the latter engines. Document were updated off-line, in order to ensure query answering independent from update executions. For this experiment only, the JVM was given 4GB of RAM, in order to minimize memory swapping.

Results are reported in Figure 6.2.c. Refreshing time has been measured without import and indexing time (a full-text index has been settled for both BaseX and QizX), that we assume asymptotically constants since done only once, and updating and index-maintenance times, that is orthogonal to our goals. As in [BC09a], for each update  $u_i$  we measured the time  $r_i$  needed for refreshing *all* the 36 queries after the update, and the time  $r_i^{type}$  and  $r_i^{chain}$  needed to refresh only queries that are not deemed as independent by the static analysis of [BC09a] and by ours, respectively. In Figure 6.2.c, for each of the three used engines we report the averages of all refreshing times  $r_i, r_i^{type}$

---

<sup>1</sup>[www.saxproject.org](http://www.saxproject.org), [www.basex.org](http://www.basex.org), [www.xmlmind.com](http://www.xmlmind.com)

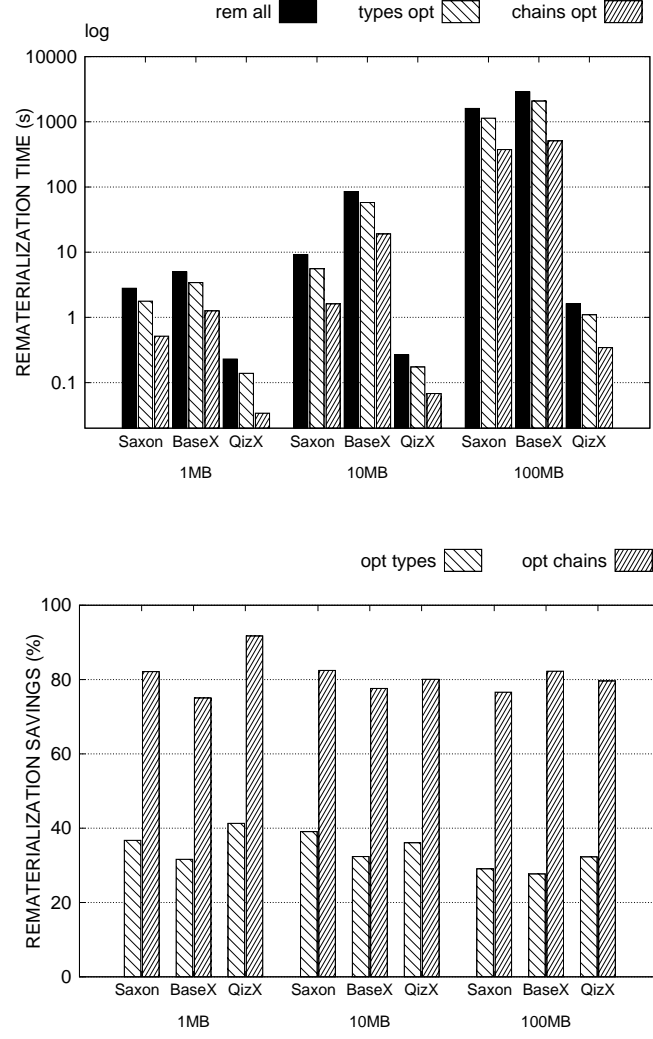


Figure 6.3: View Refreshing: Types vs Chains

$r_i^{chain}$ . As a consequence of time efficiency and precision of our static analysis, even for a relatively small document of 1MB, our independence analysis ensures high time savings for all engines: 82% for Saxon, 75% for BaseX and 85% for QizX. While type based analysis [BC09a] ensures much lower time savings: 36% for Saxon, 31% for BaseX and 37% for QizX. These percentages are essentially the same as those obtained for 10MB and 100MB documents, both for our technique and for that of [BC09a]. This is because in the considered benchmark, queries that are not statically deemed as independent of an update, and hence refreshed, are the most expensive ones to refresh.

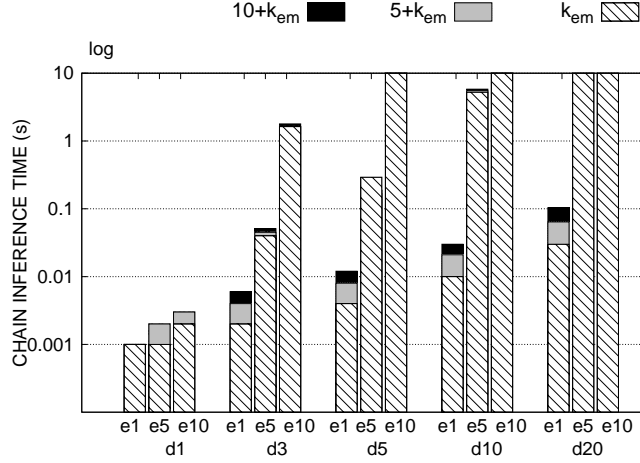


Figure 6.4: Chain Inference Time on R-benchmark

### Scalability on R-benchmark

We ran experiments on the R-benchmark by ranging over each saturated recursive schema  $dn$  ( $n$  ranging over  $\{1, 3, 5, 10, 20\}$ ), and each expression  $em$  ( $m$  ranging over  $\{1, 5, 10\}$ ). Also, for each run we varied the multiplicity values, by considering  $k_{em}$ ,  $k_{em} + 5$  and  $k_{em} + 10$ .

Results are now described. The schema  $d5$  is quite complex, it contains 5 mutually recursive types. We can see from Figure 6.2 that even with such complex form of recursion, for  $e5$ , and for each  $k \in \{5, 10, 15\}$ , chain inference is still fast (inference time is around a decimal of a second). For schema  $d10$ , featuring an extremely complex form of recursion, inference time is around five seconds for  $e5$ , while for  $e10$  the time exceeds ten seconds. The same happens for more complex cases. These test results show that even for forms of recursions that are unlikely to occur in practice (like the  $d5$ - $e5$  case), chain inference is still fast, while it takes more than one second for extremely complex cases.

We see that augmenting  $k$  does not substantially raise the inference time for an expression. This explained by the fact that a larger  $k$  value just enlarges the space of returned chains, without resulting in more context switches (where chains for step  $i$  are passed as input of chains for step  $i+1$ ) that dominate the cost of inference. We see that for  $d10$  and  $d20$ , computing chains for  $e5$  is expensive, while for smaller queries the analysis can be afforded on such schemas. In such a case, even if unlikely to occur in practice, our analysis can still ensure substantial time savings in view maintenance and thus it may be worth running the analysis.

In this chapter we presented a set of experiments validating the precision and the efficiency of our independence analysis. On the XMark benchmark, our method resulted to have double precision as the related approach based on types and comparable inference time, thus validating the interest of using our analysis based on chains. Despite the fact



that our method runs in polynomial space and time, the non-linear complexity required us to perform a worst-case analysis, so as to outline when a massive use of recursion may slow down inference time. This confirmed the analysis is fast for a large class of schemas we find in practice.

A1	A2	A3	A4	A5	A6	A7	A8	B1	B2	B3	B4	B5	B6	B7	B8				
1	2	2	1	2	2	1	1	2	3	2	2	3	3	2	3				
Q1	Q2	Q3	Q4	Q5	Q6	Q7	Q8	Q9	Q10	Q11	Q12	Q13	Q14	Q15	Q16	Q17	Q18	Q19	Q20
1	1	2	1	2	3	3	1	2	1	1	1	1	2	2	2	1	1	2	1

Table 6.1:  $k$  values for XMark/XPathMark query and update expressions

UI1 $k = 3$	for x in <i>//person/people[not(homepage)]</i> return insert <homepage>{"www.myhome.com"}</homepage> into x	Insert a default node for unfilled homepages.
UI2 $k = 3$	let i := <i>//item[1]</i> return for x in <i>//regions</i> return if(empty(x/item)) then insert i into x else ()	Insert the first document item into empty regions.
UI3 $k = 2$	insert <i>site/closed_auctions/closed_auction[1]</i> into <i>/site/closed_auctions</i>	Duplicate the first closed auction.
UI4 $k = 3$	for x in <i>//watches/watch</i> return insert <watch open_auction={"open_auction5"}> after x	Insert a node for each watch element.
UI5 $k = 1$	insert "first_text_element" as first into <i>//text[1]</i>	Mark the first text element with a statement.
UN1 $k = 2$	for x in <i>//item/description/text/keyword</i> return rename x with <i>bold</i>	Rename some keywords elements as bold.
UN2 $k = 2$	for x in <i>//bold</i> return rename x with <i>emph</i>	Rename all bold elements as emph.
UN3 $k = 2$	for x in <i>//profile</i> return rename x/education with <i>interest</i>	Rename all education elements as interest.
UN4 $k = 2$	for x in <i>//item[payment]</i> return rename x with <i>item</i>	Rename some item elements as item.
UN5 $k = 3$	for x in <i>//person[not(address) and phone]</i> return rename x/phone with <i>address</i>	Rename some phone element as address.
UP1 $k = 2$	replace <i>//open_auctions/open_auction[@id="open_auction0"]</i> with ()	Replace an open auction with empty.
UP2 $k = 3$	for x in <i>//open_auction[@id="open_auction0"]</i> return replace x/bidder/personref[@id="person1"] with <personref id={"person89"}>	Replace a bidder in an open auction.
UP3 $k = 2$	for x in <i>//emph</i> return let y:= x/text() return replace x with <bold>{y}</bold>	Replace all emph elements with bold ones, yet preserving content.
UP4 $k = 3$	for x in <i>//text[count(/keyword="fornitures")=0]</i> return replace x with <text>{"element_removed"}</text>	Replace some text elements removing their content.
UP5 $k = 2$	replace () with <i>//keyword//bold</i>	Replace nothing.

Table 6.2: insert, rename and replace update expressions extending benchmark of [BC09a]

	A1	A2	A3	A4	A5	A6	A7	A8	B1	B2	B3	B4	B5	B6	B7	B8	Q1	Q2	Q3	Q4	Q5	Q6	Q7	Q8	Q9	Q10	Q11	Q12	Q13	Q14	Q15	Q16	Q17	Q18	Q19	Q20	Types[BC09a]	Chains	TOT	Types[BC09a] (%)	Chains (%)	
UA1	D	D	D	D	D				I						I							I	I	I	I	I				I		N	N			I		17	30	31	55	97
UA2	D	D	D	D	D					D					I							I	I	I							I	N	N			I		20	27	30	67	90
UA3	D	D	D	D	D					D					I							I	I	I							I	N	N			I		20	27	30	67	90
UA4	I	I	I	D	D						I	I	I	I	I	I	I	I	I	I	I	I	I				I	I	I		I				I	I		8	34	34	24	100
UA5	I	I	I	D	D					I	I	I	I	I	I	I	I	I	I	I	I	I						I	I		I				I	I		8	34	34	24	100
UA6	I	I	I			D	D	D	I	I					N	D					I	I	I	D	D	D	D	I	I	I			D			N		11	25	27	41	93
UA7	I	I	I			D	D	D	I	I					N	D					I	I	I	D	D	D	D	I	I	I			D		I	N		11	25	27	41	93
UA8	I	I	I			D	D	D	I	I					N	D					I	I	I	D	D	D	D	I	I	I			D		I	N		11	25	27	41	93
UB1	I	I	I			I	I	I	D	I				D	D	I					I	I	I	I	I	I	I	I	I	D	I			I	D	I		11	31	31	35	100
UB2	I	D	D	I	N					N		D	D	I	I		D	I	D	D	D	I	I	I				I	I	D	I		I		I	I		20	31	33	61	94
UB3											D	D					D	I	D	D	D	I	I	I			I	I	I		I			I	I		15	30	30	50	100	
UB4											D	D					D	I	D	D	D	I	I	I			I	I	I		I			I	I		15	30	30	50	100	
UB5						I	I	I	D	I				D	D						I	I	I	I	N	I	I	I	I	I	D	I		I	D	I	11	30	31	35	97	
UB6						I	I	I	D	I				D	D						I	I	I	I	N	I	I	I	I	D	I			I	D	I	11	30	31	35	97	
UB7						D	D	D	I	I				I	I	N	D				I	I	I	D	D	D	D	I	I	I	I			D		I	N	11	25	27	41	93
UB8						I	D	D			I	I	I	I	I		D	I	I	I	I	I	I	I	I	I	D	I	I	I	I			D		I	I	15	35	35	43	100
UI1						I	D	D													I	D	D	I	N		D	I	I	I	I		D			D		15	31	31	48	100
UI2									D	D				D	D	I					I	D	D	I	N					D	D				D	I		23	26	27	85	96
UI3	D	D	D	D	D	I	I	I	I	D	I	I	I	I	I	I	I	I	I	I	D		D	D	D	I	I	I	I	I	I	D	D	I	I	I	0	24	24	0	100	
UI4						I	I	I									I	I	I	I	I	I	I	I	I	I	I	I	I	I	I				I	N	15	35	36	42	97	
UI5	I	I	I	I	I																I	I	I							D	I	I	I		I		20	35	35	57	100	
UN1	N	N	N	I	I					I					I	I					I	I	I							D	I	I	I		I		20	32	35	57	91	
UN2	D	D	D	I	D					N					I	I					I	I	I							D	I	I	D	D		I		20	28	30	67	93
UN3			I		I	I															I	I	I			D	I	N	I	I	I				N		20	33	35	57	94	
UN4									N	N				N	N	I					I	N		I	N					N	I				N		23	26	36	64	72	
UN5			I		I	I	D	D		I				I	I						I	I	I	I	D	I	I	I	I	I	I		I	I	N		15	32	33	45	97	
UP1	I	I	I	I	I	I	I	I	I	N	D	D		I	I	D	I	D	D	D	I	I	D	I	I	I	D	D	I	I	I	I	I	I	D	I	I	0	25	26	0	96
UP2											D	D		I	I	I	I	I	I	D	I	I							I	I	I				I		19	33	33	58	100	
UP3	D	D	D	I	D					N				I	I	I					I	I	I							D	I	D	D		I		20	28	29	69	97	
UP4	D	D	D	D	D				I	D				I	I	I					I	I	I	I	I					D	I	D	D		I		17	27	27	63	100	
UP5																																					36	36	36	100	100	
																																						49%		96%		

Table 6.3: Query-Update independent pairs detected. D=dynamic dependence; I=independence detected by chains but not by types; N=independence undetected by both approaches; White space = independence detected by both approaches

# Chapter 7

## Extensions

In this chapter we illustrate how our chain-based type system can be extended in order to deal with features not considered in the previous formal development. We will first focus on additional features concerning queries and updates, and then on those concerning schema mechanisms.

### 7.1 Queries and Updates

The XQuery fragment we have considered, is the same as the one considered in the related approaches [BC09a, BC10], and leaves out several query mechanisms. These can be handled by extending our framework by either by means of query rewriting or by defining new inference rules.

#### 7.1.1 Query rewriting

Query rewriting can be adopted to rewrite a query-update pair  $q-u$  into a pair  $q'-u'$  belonging to the language we presented in Section 2.3, and such that static independence holds for  $q-u$  if it holds for  $q'-u'$ . Thus, here query rewriting is used wrt independence, not wrt to equivalence of XQuery expressions.

As already mentioned, rewriting has already been used in related approaches [BC09a, BCCN06] as well as in our tests for dealing with XMark queries in performed tests. For instance, XMark query  $q_{14}$

```
let $auction := doc("auction.xml") return
for $i in $auction/site/item
where
contains(string(exactly-one($i/description)), "gold")
return $i/name/text()
```

is rewritten into

```

let $auction := doc("auction.xml") return
for $i in $auction/site/item
if ($i/description/descendant-or-self::node())
then $i/name/text()
else ()

```

The rewriting replaces the **where** clause with an **if-then-else** clause; function calls are discarded and only the used paths are retained; note the adding of the step **descendant-or-self::node()** in order to capture the fact that discarded function calls needs descendants of nodes selected by the path  $\$i/description$ . Recall that symbol  $\$$  is used to denote a variable in XQuery.

Query rewriting is an obvious solution to deal with *conditional* queries of the form

**if** ( $q$ ) **then**  $q_1$  **else**  $q_2$

If we do not consider negation, the condition  $q$  is rewritten into a query  $q'$ , by replacing in  $q$  each operator  $\theta \in \{\text{and}, \text{or}, =, \leq\}$  with the query concatenation operator “,”. For instance, let  $x$  and  $y$  be two variables, predicate

$q = (x/a \text{ and } x/b) \text{ or } x/c/text()=y/f/text()$

is rewritten into  $q' = (x/a, x/b, x/c/text(), y/f/text())$ . It is easy to see that if an update does not impact  $q'$  then it does not impact  $q$ .

*Negation* requires to take into account that an update modifying a node which is neither used nor returned by  $q$  can impact **not**( $q$ ). For instance, assume variable  $x$  bound to a node having  $a$  and  $b$  children only, then an update deleting all  $a$  children impacts **not**( $x/b$ ).

Negation is taken into account by rewriting **not**( $q$ ) into  $q'$  by simply replacing each node-test in  $q$  with the node-test **node()**. This rewriting ensures a sound analysis with a loss of precision in some cases. For the previous example, **not**( $x/b$ ) is rewritten into  $x/node()$ , enabling the chain analysis to exclude independence with the aforementioned deletion.

The above described rewriting is essentially the same used in [BCCN06] in the context of type-based projection for query optimization. According to the above illustrated rewriting, a projection for  $q'$  is indeed a sound projection for  $q$  as well.

### 7.1.2 New Inference Rules

Query mechanisms left out can be handled by introducing new rules for inferring chains and by determining the  $k$  multiplicity value. For instance we can add new rules to handle

**where** and **order – by** clauses and builtin functions. These new rules are quite similar to the existing ones, and presented in Table 7.1 whereas rules for inferring the multiplicity  $k$  value are presented in Table 7.2.

For the **where** clause the extension consists of adding two new rules for **for** and **let** expressions containing a **where** clause. For the **order – by** clause a similar reasoning applies. For conciseness, we restrict to descending ordering (the parameter does not affect the analysis), and assume **where** always co-occurs with **order – by**.

Concerning *types switch* expressions like

```
typeswitch(q)
  case T1 return q1
    ⋮
  case Tn-1 return qn-1
default return qn
```

a sound and precise analysis can be ensured without involving types  $T_i$  in the analysis. The chain inference for these expressions is similar to the chain inference for **if** expressions.

During chain inference for each query  $q_i$ , the type  $T_i$  could be taken into account to improve precision, by filtering out those chains inferred from  $q$  that can not be generated by any of the  $T_i$ 's in the spirit of the precise type analysis provided in [BCF03] for pattern-matching in the CDuce language.

For *functions call*  $f(q_1, \dots, q_n)$  new rules needs to take into account the nature of the function itself, and more specifically the way input parameters are used. Table 7.1 collects rules for several XQuery built-in functions used in the XMark query set that show how to deal with possible different cases.

For instance, the rule for **count**( $q$ ) is simple: it states that *used* chains of **count**( $q$ ) are those of  $q$  plus the return chains of  $q$ , and also that **count**( $q$ ) has no return chains. The return-to-used conversion is illustrated by the following example. Consider the query **count**( $x//employee$ ), for the argument  $x//employee$  we infer return chains pointing to *employee* elements selected by  $x//employee$ .

Since counting these nodes does not depend on their descendants (e.g., *employee* has *salary* element as child), for a precise-analysis the return-to-used conversion is needed. Concerning element chains eventually inferred from the input  $q$ , these do not have to be considered as chains for **count**( $q$ ). Recall that element chains are inferred for **return** clauses of queries inside updates, so that precise chain inference for insert/replace updates can be done.

For the function **string** (that converts a value or a node to a string) the rule is simpler: used and return chains of **string**( $q$ ) are those of  $q$ . In this case the conversion is

not needed due to the semantics of the `string` function: the descendants of return nodes of `q` are needed to compute the function output. Aggregation functions `sum`, `min`, `max` and `avg`, as well as conversion functions such as `data` and `number` are treated in the same way as `string`. Finally, `contains` function is defined similarly too, by collecting used and return chains of the queries passed as parameters.

For *user-defined functions calls* the rewriting method can be used as follows. Assuming that the definition of the function `f` is *non recursive* and given by

$$(x_1, \dots, x_n) = q$$

the function call `f(q1, ..., qn)` is rewritten into

```
let x1 := q1
  ⋮
let xn := qn
return q
```

for chain inference and independence analysis (of course the rewriting is recursively done inside `q` eventually). Dealing with recursive functions raises problems wrt to determining the multiplicity value  $k$ , and is out of the scope of this work.

Concerning *transform* expressions of the form

```
copy x := q1
modify u1
return q2
```

we have to notice that they do not change the input document, but rather a copy of some of its fragments.

In order for a transform query to be independent of an update `u`, we need that `q1`, `q2` and the 'query component' of `u1` (the part of `u1` that queries the schema instance) be independent of `u`. So chain inference for transform queries can be specified by relying on the existing chain inference for `q1` and `q2`, while for inferring chains from `u1` some modifications on update chain inference need to be done so that chains coming from the query component of `u` be taken into account. We leave this extension as future work. Note that the same kind of extension is needed for the related approach [BC09a].

Concerning the inference of the  $k$  multiplicity value for added constructs this is done along the lines of definitions for the core considered in Section 2.3. The extended definition is reported in Table 7.2.

$$\begin{array}{c}
\frac{
\begin{array}{c}
\Gamma \vdash_C q_1 : (r_1, v_1, e_1) \\
\Gamma[x \mapsto c] \vdash_C q^w : (r_c^w, v_c^w, e_c^w) \\
\Gamma[x \mapsto c] \vdash_C q^o : (r_c^o, v_c^o, e_c^o) \\
\Gamma[x \mapsto c] \vdash_C q_2 : (r_c, v_c, e_c) \text{ for any } c \in r_1
\end{array}
}{
\begin{array}{c}
\text{for } x \text{ in } q_1 \\
\text{where } q^w : (\bigcup_{c \in r_1} r_c, v_1 \cup \bigcup_{\substack{c \in r_1 \\ \lambda \in \{w, o\}}} (r_c^\lambda \cup v_c^\lambda) \cup \bigcup_{\substack{c \in r_1 \\ r_c \cup e_c \neq \emptyset}} (v_c \cup \{c\}), \bigcup_{c \in r_1} e_c) \\
\text{order-by } q^o \\
\text{return } q_2
\end{array}
} \text{(FORWHEREOB)}
\\[20pt]
\frac{
\begin{array}{c}
\Gamma \vdash_C q_1 : (r_1, v_1, e_1) \\
\Gamma[x \mapsto r_1] \vdash_C q^w : (r^w, v^w, e^w) \\
\Gamma[x \mapsto r_1] \vdash_C q^o : (r^o, v^o, e^o) \\
\Gamma[x \mapsto r_1] \vdash_C q_2 : (r, v, e)
\end{array}
}{
\begin{array}{c}
\text{let } x := q_1 \\
\text{where } q^w : (r, r_1 \cup \bigcup_{\lambda \in \{w, o\}} (r^\lambda \cup v^\lambda) \cup v_2, e_2) \\
\text{order-by } q^o \\
\text{return } q_2
\end{array}
} \text{(LETWHEREOB)}
\\[20pt]
\frac{
\begin{array}{c}
\Gamma \vdash_C q : (r_0, v_0, e_0) \\
\Gamma[x \mapsto r] \vdash_C q_i : (r_i, v_i, e_i) \text{ for any } i \in \{1, \dots, n\}
\end{array}
}{
\begin{array}{c}
\text{typeswitch } q \\
\text{case } T_1 \text{ return } q_1 \\
\Gamma \vdash_C : : (\bigcup_{i \in \{0 \dots n\}} r_i, \bigcup_{i \in \{0 \dots n\}} v_i, \bigcup_{i \in \{1 \dots n\}} e_i) \\
\text{case } T_{n-1} \text{ return } q_{n-1} \\
\text{default return } q_n
\end{array}
} \text{(TYPESWITCH)}
\\[20pt]
\frac{\Gamma \vdash_C q : (r, v, e)}{\Gamma \vdash_C \text{count}(q) : (\emptyset, r \cup v, \emptyset)} \text{(COUNT)} \qquad \frac{\Gamma \vdash_C q : (r, v, e)}{\Gamma \vdash_C \text{string}(q) : (r, v, \emptyset)} \text{(STRING)}
\\[20pt]
\frac{\Gamma \vdash_C q_i : (r_i, v_i, e_i) \text{ for } i = 1, 2}{\Gamma \vdash_C \text{contains}(q_1, q_2) : (r_1 \cup r_2, v_1 \cup v_2, \emptyset)} \text{(CONTAINS)}
\end{array}$$

Table 7.1: Chain Inference Rules for Extensions



$$\frac{f = \mathcal{F}(a, q_1) + \max\{\mathcal{F}(a, q_2), \mathcal{F}(a, q^w), \mathcal{F}(a, q^o)\}}{\mathcal{F}(a, \text{for/let } x \ q_1 \text{ where } q^w \text{ order-by } q^o \text{ return } q_2) = f} (\mathcal{F}\text{-ORDER-BY})$$

$$\frac{f = \mathcal{F}(a, q) + \max\{\mathcal{F}(a, q_i)\}}{\mathcal{F}(a, \text{typeswitch } q \text{ case } T_i \text{ return } q_i \text{ default return } q_n) = f} (\mathcal{F}\text{-TYPESWITCH})$$

$$\frac{f = \max\{\mathcal{F}(a, q_1), \mathcal{F}(a, q_2)\}}{\mathcal{F}(a, \text{contains}(q_1, q_2)) = f} (\mathcal{F}\text{-CONTAINS})$$

$$\frac{f = \mathcal{F}(a, q)}{\mathcal{F}(a, \text{string}(q)) = f} (\mathcal{F}\text{-STRING}) \quad \frac{f = \mathcal{F}(a, q)}{\mathcal{F}(a, \text{count}(q)) = f} (\mathcal{F}\text{-COUNT})$$

$$\frac{r = \mathcal{R}(q_1) + \max\{\mathcal{R}(q_2), \mathcal{R}(q^w), \mathcal{R}(q^o)\}}{\mathcal{R}(\text{for/let } x \ q_1 \text{ where } q^w \text{ order-by } q^o \text{ return } q_2) = r} (\mathcal{R}\text{-ORDER-BY})$$

$$\frac{r = \mathcal{R}(q) + \max\{\mathcal{R}(q_i)\}}{\mathcal{R}(\text{typeswitch } q \text{ case } T_i \text{ return } q_i \text{ default return } q_n) = r} (\mathcal{R}\text{-TYPESWITCH})$$

$$\frac{r = \max\{\mathcal{R}(q_1), \mathcal{R}(q_2)\}}{\mathcal{R}(\text{contains}(q_1, q_2)) = r} (\mathcal{R}\text{-CONTAINS})$$

$$\frac{r = \mathcal{R}(q)}{\mathcal{R}(\text{string}(q)) = r} (\mathcal{R}\text{-STRING}) \quad \frac{r = \mathcal{R}(q)}{\mathcal{R}(\text{counts}(q)) = r} (\mathcal{R}\text{-COUNT})$$

Table 7.2:  $\mathcal{F}(,)$  and  $\mathcal{R}(,)$  definition for Extensions

## 7.2 Schemas

### 7.2.1 Attributes

As most of the studies on XQuery type systems [BCCN06, BC09a], we have excluded attribute types in our presentation, mainly because these does not introduce any particular issue. In any case, our implementation and our tests take attribute types into account.

To handle attribute types, the implementation considers chains possibly having a final label of the form  $@\alpha$ . A chain  $c.@\alpha$  is generated by a schema if an element pointed by  $c$  in a schema instance may have an attribute labeled as  $\alpha$  (w.l.o.g, we assume the symbol  $@$  not used in element tags defined by the schema). To infer such chains, the following new rule handling the **attribute** axis has been implemented:

$$\mathbf{A_C}(c, \text{attribute}) \stackrel{\text{def}}{=} \{ c.@\alpha \mid c.@\alpha \in \mathbf{C} \}$$

Chain based independence is straightforwardly extended to handle the attribute extension. Furthermore, the inference of the  $k$  multiplicity value is not impacted by attribute types, as they cannot occur more than once in a chain.

### 7.2.2 Keys and foreign-keys integrity constraints

Concerning ID/IDREF constraints in DTDs, and key/keyref constraints in XSDs (studied in [AFL02]), we assume they are preserved by updates, as we assume that validity is preserved (Section 2.4). So, in order to ensure precise and sound independence analysis, chain inference does not need to consider these constraints. Our notion of static independence only concerns the type component of the schema, while these constraints pose restrictions on the values of attributes and elements in a document, and do not impact its structure.

To better illustrate, consider the following XMark  $q_8$  [SWK<sup>+</sup>02] query

```
let $auction := doc("auction.xml")
for $p in $auction/site/people/person
let $a :=
  for $t in $auction/site/closed_auctions/closed_auction
  where $t/buyer/@person = $p/@id
  return $t
return <item person = "$p/name/text()" > count($a) </item >
```

Here ID/IDREF attributes are used, and selected by the following path expressions.

```
p1 = doc("auction.xml")/site/people/person/@id
p2 = doc("auction.xml")/site/site/closed_auctions/
      closed_auction/buyer/@person
```

Hence our chain inference system generates chains

```
site.people.person.@id
site.site.closed_auctions.closed_auction.buyer.@person
```

If an update *u* changes either an *id* or *person* attributes, then are generated update chains like

```
site.people.person : @id
site.site.closed_auctions.closed_auction.buyer : @person
```

These chains conflict with query chains, hence static independence is excluded.

**Keyref constraints in XML Schema** XML Schema provides powerful mechanisms to define keys and references formed by multiple attribute/element nodes. To this end, a fragment of XPath including downward axes and no predicates is used to specify the position and structure of keys and references. The example below is borrowed from the W3C specification [TBMM04].

```
<xs:key name = "regKey" >
<!-- vehicles are keyed by a pair of state and plate -->
  <xs:selector xpath = "../vehicle" />
  <xs:field xpath = "@state" />
  <xs:field xpath = "@plateNumber" />
</xs:key >

<xs:keyref name = "carRef" refer = "regKey" >
<!-- people's cars are a reference -->
  <xs:selector xpath = "../car" />
  <xs:field xpath = "@regState" />
  <xs:field xpath = "@regPlate" />
</xs:keyref >
```

Both declarations are made inside the definition of a *root* element (whose definition can be found in [TBMM04]). So the two selectors *../vehicle* and *../car* navigates starting from a *root* element, and respectively select the parents of the key *regKey* and foreign key *carRef*. Both key and foreign keys consist of a pair of attributes selected by relative path expressions in the *xs:field* elements. Field paths navigate starting from nodes selected by respective selectors. An important distinction wrt DTDs is that a field path may also select subelements of elements selected by the selector.

It is worth observing that the specification of key and keyref constraints and the specification of the type of the documents are different in nature. The purpose of key and keyref is to constraint the values of the elements of documents valid wrt a schema and have no impact on their structure. These constraints are defined using paths navigating valid documents. As a consequence, as for DTDs, key and keyref constraints can be totally

ignored by the static analysis in order to ensure a sound and precise static detection of independence.

One of the nice features of our static analysis is that it can be directly extended so as to capture Extended DTDs, and thus it is suitable also for XML Schema and a core of RelaxNG type languages.

### 7.2.3 Extended DTDs

As outlined in Chapter 1, EDTDs differ from DTDs mainly because they allow to write several type definitions for the same element. In order to perform our chain analysis, we have to reflect this aspect while doing chain inference.

Recall that in an EDTD [PV00] we have a “type alphabet”  $\Sigma' = \{a_i | a \in \Sigma\}$  and a mapping to a “tag alphabet”  $\mu$  such that  $\mu(a_i) = a$  for all  $a_i \in \Sigma'$ . This captures the fact that two types differently indexed produce the same label but possibly different content models.

That said, it is sufficient to change the definition of reachability Definition 4.1.2, and of some step chain inference rules. Both changes are straightforward. Reachability should be redefined so as chains contain indexed symbols  $a_i$  instead of  $a$ , referring to different definitions of the same element tags.

Then, chain inference rule for step filtering has to be modified in order to deal with indexed symbols, as follows.

$$T_C(c.a_i, b) = \{c.a_i \mid \mu(a_i) = b\}$$

Note that inference rules for axis chain inference remains unchanged.

Once chains have been inferred for a query and an update, to check query-update independence we use a slight variation of Definition 4.4.1, where we look for overlapping chains up to indexes  $i$  in chains symbols  $a_i$ ’s.

Notice that precision of the inference as well as complexity results remain unchanged for the EDTD case, since this can be simply seen as a schema with more definitions. Also  $<_r$  definition can be easily adapter for dealing with interleaving operator of Relax NG.

For EDTDs, the CDAG compression needs to be a bit more careful. In Figure 7.2.3 we write an EDTD  $d_{10}$ , featuring two definitions of element **name**, that are used in different contexts. Consider CDAG storing chains of  $d_{10}$  made by considering only the tag-name of types, depicted in (a). This way of representing chains would include some artifacts (chains that do not belong to the schema), such as the following ones.

`book.author.name.String`   `book.publisher.name.first.String`   `book.publisher.name.last.String`

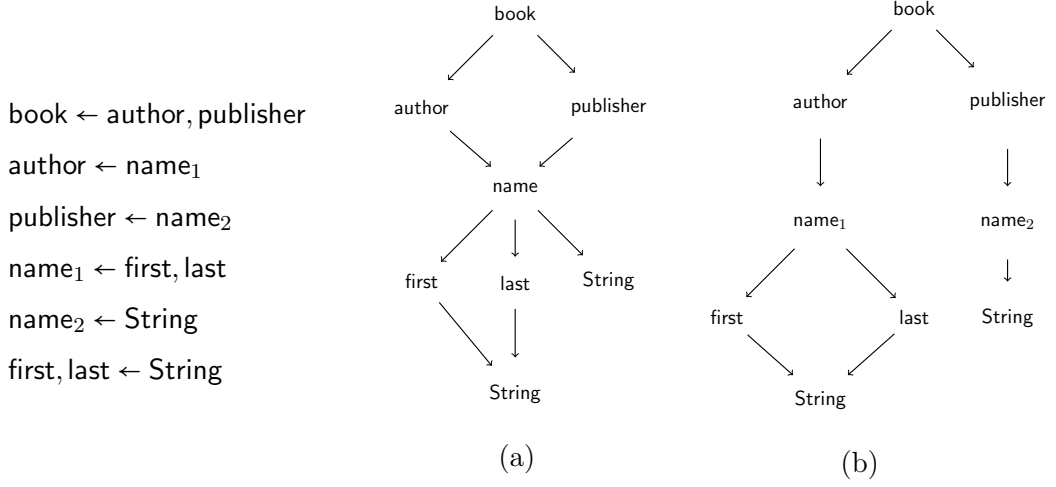


Figure 7.1: Example of CDAGs for EDTD

However, this situation can be simply avoided by considering that type `name1` and `name2` come from different definitions. The result is depicted in (b). We stress that types `name1` and `name2` are considered as different when we store chains, but as the same type `name` when we do independence detection. Notice that the CDAG inferred from an EDTD is again polynomial in the size of the EDTD.

## Conclusions

In this chapter we discussed some extension of our static analysis based on schema chains. We would like to stress that dealing with extensions is easy because our chain inference and independence analysis are based on the notions of used, return and element nodes (Chapter 4). These are *universal* and *essential* notions, in the sense that each processed node of whatever kind of query one wants to add to the framework falls in one of these three categories. Thus, generalizing our framework for a new query construct mainly consists of identifying how used, return and element nodes are determined. This simply requires the understanding of the semantics of query constructs. In the next chapter we draw the conclusions of our work and we outline possible future perspectives.

## Chapter 8

# Conclusions and Future Perspectives

In this thesis we have tackled the problem of statically detecting query-update independence in the context of typed XML databases. After having outlined the important issues raised by this problem in the context of XML data management, we provided an exhaustive overview of related existing techniques, by including both typed and untyped approaches. As discussed, existing techniques still suffer of several limitations.

To overcome these limitations, in this thesis we made the following contributions:

- We presented a type system able to statically detect XML query-update independence. One of the main features of the type system is the chain inference component, allowing to infer information at the basis of an highly precise analysis.
- We proposed a method to restrict the analysis to a finite set of chains in the presence of recursive schemas.
- We proposed techniques for a sound and efficient implementation of the above mentioned finite analysis.
- We conducted extensive experiments in order to validate precision and efficiency of our static analysis. As shown by test results, our technique succeeded from both aspects.

We also illustrated how the static analysis can be extended in order to deal with several features not considered in the core query and update languages, in order to show that the analysis is amenable to be adopted in most of the XQuery query and update languages, by preserving precision and efficiency.

## 8.1 Future Perspectives

We discuss here some future directions that we believe are worth being investigated.

### $\alpha$ -types

One of the main goals of this thesis is to show that schema information can be an extremely powerful tool for optimizing XML queries and updates. We presented a type system for inferring static projections of queries and updates (Theorems 4.3.4 and 4.3.6), used for the purpose of checking independence, that takes deeply into account informations concerning the hierarchical structure of data. We believe that a schema-based analysis, that is precise, cannot abstract away from such structural constraints. At the same time, when designing a method, there is a quest for tractability and efficiency of static analysis. Our method can run in polynomial time by introducing some approximations and has been proved fast by experiments. However particular applications may be interested in further lowering the computational time needed by the analysis. As shown in Chapters 3 and 6, one way to do this is to simplify the system and infer simple types [BC09a] rather than chains of types. To mitigate the lack of precision that this approach comports, it would be interesting to study inference for a *parametric* notion of type, called  $\alpha$ -type. An  $\alpha$ -type is a chain (not necessarily rooted) whose length is given by  $\alpha$  (a positive natural value) and whose first  $\alpha-1$  type labels determine the *context* of the type indicated by the  $\alpha^{th}$  symbol, typing a set of nodes accessed by the query or update. The parameter  $\alpha$  determines the amount of contextual information used to disambiguate types during the inference, and it can be opportunely tuned so as to vary the precision of the independence analysis wrt the computational resources that are available, because also the cost of the inference would result to be mainly determined by  $\alpha$ . Note that with such an approach we would have that  $\alpha = 1$  would enable to capture the type system in [BC09a], while unbounded  $\alpha$  would enable to capture the infinite analysis here proposed. One of the main interesting aspects of this research consists of finding methods for computing optimal  $\alpha$  values, by taking into account queries, updates, schemas, and required constraints about execution time.

### Probabilistic Independence Analysis

The decision procedure that we developed to check static independence gives a binary answer yes/no. However, because our method performs a conservative analysis, a query and an update may still be independent over a schema instance, even if this has been excluded by the static analysis. This may happen for instance, because a conflict is detected on an optional element tag that occurs quite rarely in practice. This said, it seems interesting to refine our method so as to take into account schemas with probabilities. Probabilistic schemas are useful in many contexts. As discussed in [AAD<sup>+</sup>12], they can be used for instance to describe the distribution of the data, telling how many element nodes of a certain type are found in the database. By assigning to each chain an uncertainty value one could also estimate the probability that a conflict happened. This for

instance can be used to design more sophisticated view-maintenance algorithms that do maintenance only when the estimated probability of conflict is higher than a threshold.

### **Schema-less Analysis**

Designed for the Web, XML documents do not always come with a schema. Besides relying on schema inference techniques such as the one proposed by Bex et al. [BNSV10], this case can be handled by performing chain inference starting from a data summary (e.g., DataGuide) in place of a schema. A DataGuide essentially is a superset of the chains belonging to the data, and can be efficiently inferred [GW97]. Our system rules can be easily adapted because of the chain based nature of a DataGuide.

### **Independence Analysis for the Semantic Web**

Query-update independence for Semantic Web databases is an open problem. Semantic web databases are graphs written in RDF. Despite designed for tree-shaped data, we believe that our method can be applied also for RDF graphs by making an analysis of accessed classes of data, provided that an ontology or a data summary is available. Also, our method could be applied for n-SPARQL, the navigational language for RDF [PAG10].





## Chapter 9

# Proofs

In this chapter we prove the statements concerning the correctness of our chain analysis. The whole proof is organized in two main parts. In the first part we prove that the chain-based independence analysis is sound, in the case where a possibly infinite set of chains is inferred for an expression over a schema. In the second part we prove that the finite restriction to the subset of  $k$ -chains yields an analysis that is equivalent to the infinite one.

Proving that the possibly infinite chain analysis is sound mainly consists of showing that the sets of return, used, element and update chains inferred for an expression correctly abstract the data that is accessed/updated during the evaluation of the expression, over a valid schema instance. As a corollary of this we conclude the correctness of our notion of chain-based independence.

Proving that the finite restriction to the subset of  $k$ -chains yields an analysis that is equivalent to the infinite one. Consists of two main tasks. We first show that all chains can be reduced to  $k$ -chains. Then we show that any pair of conflicting chains can be reduced to a pair of conflicting  $k$ -chains. These proofs are developed by first introducing *linear-path queries* (see Section 9.2).

### 9.1 Soundness of Chain Inference in the Possibly Infinite Case

In this section we prove that our analysis makes a sound abstraction of accessed/updated data in the possibly infinite case. Because updates are defined in terms of queries, and queries in turn are defined by using navigational capabilities, the proof is structured as follows. We first prove correctness for XPath navigational steps. We then prove correctness for XQuery expressions. We finally prove correctness for XQuery Update expressions. Furthermore, we show that for XPath navigational steps chain inference is also complete.

### 9.1.1 Soundness and Completeness of Schema Type Relations

Correctness (and also completeness) of chain inference for XPath navigational steps relies on the correctness of the auxiliary relations defined to capture with types the parent-child and sibling relations among documents nodes, namely  $\Rightarrow_d$  and  $<_r$ .

We show that  $\Rightarrow_d$  and  $<_r$  are both sound and complete. Soundness means that for any pair of parent-child (respectively sibling) locations belonging to a valid tree, there exists a correspondent pair of types in  $\Rightarrow_d$  (respectively in  $<_r$ ). Completeness means that these relations are also minimal, thus removing one pair from  $\Rightarrow_d$  or  $<_r$  we break soundness.

**Lemma 4.1.3** (Soundness of  $\Rightarrow_d$ ) *Let  $d$  be a DTD and  $t = (\sigma, l_t) \in d$  a valid tree. If  $(l, l') \in \text{Child}_t$  then  $\text{type}(l) \Rightarrow_d \text{type}(l')$ .*

*Proof.* Since  $t \in d$ , by definition of validity (Section 2.2) it follows that  $\text{type}(l_2) \in d(\text{type}(l_1))$  and thus  $\text{type}(l_1) \Rightarrow_d \text{type}(l_2)$ .  $\square$

**Lemma 4.2.2** (Soundness of  $<_r$ ) *Let  $d$  be a DTD and  $t = (\sigma, l_t) \in d$  a valid tree. If  $(l, l') \in \text{FollowingSibling}_t$  and  $c_l^t = c_{l'}^t$  then  $\text{type}(l) <_{d(c)} \text{type}(l')$ .*

*Proof.* If  $(l_1, l_2) \in \text{FollowingSibling}_t$  and  $(l, l_1) \in \text{Child}_t$  and  $r = d(c_l^t)$  then, by definition of validity, there exists a word  $w \in \text{Lang}(r)$  where  $\text{type}(l_1)$  occurs before  $\text{type}(l_2)$ . We show that this in turn implies that  $\text{type}(l_1) <_r \text{type}(l_2)$ . By cases on  $r$ .

*Base.*

(Cases  $r = \epsilon$  and  $r = a$ ) Trivial since such  $w$  cannot exist.

*Induction.*

(Case  $r = r_1, r_2$ ) Because  $\text{Lang}(r_1, r_2) = \text{Lang}(r_1) \cup \text{Lang}(r_2) \cup \text{Lang}(r_1) \times \text{Lang}(r_2)$  then either  $w \in \text{Lang}(r_1) \cup \text{Lang}(r_2)$  or  $w \in \text{Lang}(r_1) \times \text{Lang}(r_2)$ . In the first case, we conclude by induction on  $r_1, r_2$ . In the second case, we conclude since  $(\text{type}(l_1), \text{type}(l_2)) \in \text{Sym}(r_1) \times \text{Sym}(r_2) \subseteq <_r$ .

(Case  $r = r_1 | r_2$ ) In this case we have that  $\text{Lang}(r) = \text{Lang}(r_1) \cup \text{Lang}(r_2)$ . Hence at least one of  $\text{type}(l_1) <_{r_1} \text{type}(l_2)$  and  $\text{type}(l_1) <_{r_2} \text{type}(l_2)$  hold, and we conclude by induction.

(Cases  $r = r^+$  and  $r = r^*$  and  $r = r^?$ ) Immediate since  $<_{r_1^+} = <_{r_1, r_1}$  and  $r^* = r^+ | \epsilon$  and  $r^? = r | \epsilon$ .

$\square$

Let  $L = (l_1, \dots, l_n)$  be a sequence of location, in the following, with a little abuse of notation, we write  $\text{type}(L)$  for denoting the word  $\text{type}(l_1), \dots, \text{type}(l_n)$ .

**Lemma 4.1.4** (Completeness of  $\Rightarrow_d$ ) *For all DTD  $d$ , if  $a \Rightarrow_d b$  then there exists a valid tree  $t = (\sigma, l_t) \in d$ , and a pair of locations  $(l, l') \in \text{Child}_t$ , such that  $\text{type}(l) = a$  and  $\text{type}(l') = b$ .*

*Proof.* *i)* Recall that the DTD  $d$  does not contain any type definition which is non-terminating or unreachable from the root. If  $a \Rightarrow_d b$  then there exists a chain  $c \in C_d$  of the form

$$a_1 \dots a_{(n-1)} \cdot a_n$$

where  $a_1 = s_d$ ,  $a_{n-1} = a$ , and  $a_n = b$ . In order to conclude the thesis, we want to show that there exists a valid tree  $t \in d$  and a pair of locations  $(l_1, l_2) \in \text{Child}_t$  such that

$$c_{l_2}^t = a_1 \dots a_n$$

Such  $t$  can be built with the following recursive procedure.

Start with  $i = 1$ .

1. (Generate an element of the desired chain)

Let  $l$  be a location such that  $\text{type}(l) = a_i$ . Attach to  $l$  a finite sequence of locations  $L$  such that the following hold.

- $\text{type}(L) \in \text{Lang}(d(a_i))$ ;
- $\text{type}(l') = a_{(i+1)}$  for some  $l' \in L$ ;

2. (Recursively generate the desired chain)

For all  $l' \in L$  such that  $\text{type}(l') = a_{(i+1)}$  do the following. If  $i < n$  then increment  $i$  by 1 and apply recursively step 1 on  $l'$ . Else, the desired chain of nodes is obtained. Go to the next step.

3. (Complete the tree so as to make it valid)

For all  $l' \in \text{dom}(t)$ , such that  $\text{Child}_t(l') = \emptyset$  and that  $\epsilon \notin \text{Lang}(\text{type}(l'))$  do the following. Attach to  $l'$  a finite and valid sequence of locations  $L'$  such that let  $r' = d(\text{type}(l'))$

- $r''$  is the regular expression obtained by replacing each  $*$ -guarded or optional subterm of  $r'$  with  $\epsilon$  (e.g.  $a^?|b^+, a^*$  becomes  $\epsilon|b^+$ );
- $\text{type}(L') \in \text{Lang}(r'')$

Repeat this step until obtaining a valid tree.

This procedure always terminates producing a finite tree since the all definitions in  $d$  terminates. Furthermore, this tree contains a pair of locations  $(l_1, l_2) \in \text{Child}_t$  such that  $c_{l_2}^t = a_1 \dots a_n$ , where of course  $\text{type}(l_1) = a$  and  $\text{type}(l_2) = b$ .  $\square$

**Lemma 4.2.3** (Completeness of  $<_r$ ) *For all DTD  $d$ , If  $a <_{d(c)} b$  then there exists a valid tree  $t = (\sigma, l_t) \in d$ , and a pair of locations  $(l, l') \in \text{FollowingSibling}_t$ , such that  $\text{type}(l) = a$  and  $\text{type}(l') = b$ , with  $c_l^t = c.a$  and  $c_{l'}^t = c.b$ .*

*Proof.* By Lemma 4.1.4 we know that we can build a tree instance with a location  $l''$  typed by a chain  $c \in C_d$ . If  $a <_{d(c)} b$  we claim that we can build a sequence of locations  $l_1, \dots, l_n$  such that the word  $\text{type}(l_1), \dots, \text{type}(l_n)$  belongs to the language of  $d(c)$ , and for some locations  $l_i$  and  $l_j$  with  $1 \leq i < j \leq n$  we have that  $\text{type}(l_i) = a$  and  $\text{type}(l_j) = b$ . We show this by induction on the structure of  $r = d(c)$ .

*Base.*

(Cases  $r = \epsilon$  and  $d(c) = a$ ) Trivial since they do not satisfy the hypothesis.

*Induction.*

(Case  $r = r_1 r_2$ ) If it holds that  $a <_{r_1} b \cup a <_{r_2} b$  then we conclude by induction. Otherwise, because  $a <_{d(c)} b$  we have that  $a \in \text{Sym}(r_1)$  and  $b \in \text{Sym}(r_2)$ . Therefore there exist two words  $w_1 \in \text{Lang}(r_1)$  and  $w_2 \in \text{Lang}(r_2)$  such that  $a \in w_1$  and  $b \in w_2$ . At this point take a sequence of locations  $L$  such that  $\text{type}(L) = w_1 w_2$ .

(Case  $r = r_1 | r_2$ ) In this case it holds that  $a <_{r_1} b \cup a <_{r_2} b$  and we conclude by induction.

(Cases  $r = r^+$  and  $r = r^*$  and  $r = r^?$ ) Immediate since  $<_{r_1^+} = <_{r_1, r_1}$  and  $r^* = (r^+ | \epsilon)$  and  $r^? = (r | \epsilon)$ .

To conclude the thesis it suffices to complete the tree so as to obtain a valid instance. This can be done by applying the procedure described at point 3 of Lemma 4.1.4.  $\square$

### 9.1.2 Axis chain inference

In this section we prove that chain inference for XPath navigational is sound.

**Lemma 9.1.1** (Soundness of Step Chain Inference). *Let  $d$  be a DTD,  $t = (\sigma, l_t) \in d$  be a valid tree, and  $l \in \text{dom}(\sigma)$  be a location of the tree. If*

$$\sigma, (x := l_x) \models x/\text{axis} :: \phi \Rightarrow \sigma, L$$

*then for each location  $l' \in L$  we have that*

$$c_l^\sigma \in T_C(A_C(c_{l_x}^\sigma, \text{axis}), \phi)$$

*Proof.* By case analysis on **axis**. Since **node()** is used as test filtering, we consider the correspondent axis evaluation and we leave as implicit the set-to-sequence conversion of rule by mean of document order.

[**axis = self**] We conclude immediately since by definition

$$[[\text{self}]]_{l_x}^\sigma = \{ l_x \} \quad \text{and} \quad A_C(c_{l_x}, \text{self}) = \{ c_{l_x} \}$$

[**axis = child**] By definition

$$[[\mathbf{child}]]_{l_x}^\sigma = \mathit{Child}_\sigma(l_x)$$

Since the document is valid wrt the DTD  $\mathbf{d}$ , for all  $l \in \mathit{Child}_\sigma(l_x)$ , let  $\mathbf{a}_x = \mathbf{type}(l_x)$  and  $\mathbf{a} = \mathbf{type}(l)$  we have that  $\mathbf{a}_x \Rightarrow_{\mathbf{d}} \mathbf{a}$ , with  $\mathbf{a}_x \in \Sigma$  and  $\mathbf{a} \in \Sigma_{\mathbf{S}}$ . Therefore the type chain of the location  $l$  is  $\mathbf{c}_l = \mathbf{c}_{l_x}.\mathbf{a}$  and by definition  $\mathbf{c}_{l_x}.\mathbf{a} \in \mathbf{C}_{\mathbf{d}}$ . We conclude since

$$\mathbf{A}_{\mathbf{C}_{\mathbf{d}}}(\mathbf{c}_{l_x}, \mathbf{child}) = \{ \mathbf{c}_{l_x}.\mathbf{a} \mid \mathbf{c}_{l_x}.\mathbf{a} \in \mathbf{C}_{\mathbf{d}} \}$$

[**axis = descendant**] By definition we have that

$$[[\mathbf{descendant}]]_{l_x}^\sigma = \mathit{Descendant}_\sigma(l_x)$$

For all  $l \in \mathit{Descendant}_\sigma(l_x)$  let  $L = (l_1, l_2, \dots, l_n)$  be the sequence of nodes such that  $(l_i, l_{i+1}) \in \mathit{Child}_\sigma$  for  $i = 1..n-1$ , with  $l_1 = l_x$  and  $l_n = l$ . Since the document is valid wrt the DTD  $\mathbf{d}$ , for all pair of locations  $l_i, l_{i+1} \in L$ , let  $\mathbf{a}_i = \mathbf{type}(l_i)$ , we have that  $\mathbf{a}_i \Rightarrow_{\mathbf{d}} \mathbf{a}_{i+1}$ . Therefore  $\mathbf{a}_{l_x} \Rightarrow_{\mathbf{d}}^+ \mathbf{a}_l$ , and the type chain of the location  $l$  is  $\mathbf{c}_l = \mathbf{c}_{l_x}.\mathbf{a}_2.\mathbf{a}_3 \dots \mathbf{a} \in \mathbf{C}_{\mathbf{d}}$ . We conclude since

$$\mathbf{A}_{\mathbf{C}_{\mathbf{d}}}(\mathbf{c}_{l_x}, \mathbf{descendant}) = \{ \mathbf{c}_{l_x}.\mathbf{c} \mid \mathbf{c}_{l_x}.\mathbf{c} \in \mathbf{C}_{\mathbf{d}}, \mathbf{c} \neq \epsilon \}$$

Notice that case **axis = descendant-or-self** is analogous.

[**axis = parent**] By definition

$$[[\mathbf{parent}]]_{l_x}^\sigma = \mathit{Parent}_\sigma(l_x)$$

Since the document is valid wrt the DTD  $\mathbf{d}$ , for  $l \in \mathit{Parent}_\sigma(l_x)$ , let  $\mathbf{a}_x = \mathbf{type}(l_x)$  and  $\mathbf{a} = \mathbf{type}(l)$  we have  $\mathbf{a} \Rightarrow_{\mathbf{d}} \mathbf{a}_x$ . Therefore  $\mathbf{c}_{l_x} = \mathbf{c}.\mathbf{a}_x$ , where the type chain of the location  $l$  is  $\mathbf{c}_l = \mathbf{c}$ . We conclude since

$$\mathbf{A}_{\mathbf{C}_{\mathbf{d}}}(\mathbf{c}_l.\mathbf{a}_x, \mathbf{parent}) = \{ \mathbf{c}_l \}$$

[**axis = ancestor**] This case is similar to the case of **descendant-or-self** axis.

By definition we have that

$$[[\mathbf{ancestor}]]_{l_x}^\sigma = \mathit{Ancestor}_\sigma(l_x)$$

Given  $l \in \mathit{Ancestor}_\sigma(l_x)$ , let  $L = (l_1, l_2, \dots, l_n)$  be the sequence of nodes such that  $(l_i, l_{i+1}) \in \mathit{Child}_\sigma$  for  $i = 1..n-1$  with  $l_1 = l$  and  $l_n = l_x$ . Since the document is valid wrt the DTD  $\mathbf{d}$ , for all pair of locations  $l_i, l_{i+1} \in L$ , let  $\mathbf{a}_i = \mathbf{type}(l_i)$  we have that  $\mathbf{a}_i \Rightarrow_{\mathbf{d}} \mathbf{a}_{i+1}$  and  $\mathbf{a}_i \in \Sigma_{\mathbf{S}}$ . Therefore  $\mathbf{a}_l \Rightarrow_{\mathbf{d}}^+ \mathbf{a}_{l_x}$  and the type chain of location  $l$  is  $\mathbf{c}_{l_x} = \mathbf{c}_l.\mathbf{a}_2 \dots \mathbf{a} \in \mathbf{C}_{\mathbf{d}}$ . We conclude since

$$\mathbf{A}_{\mathbf{C}_{\mathbf{d}}}(\mathbf{c}_{l_x}, \mathbf{ancestor}) = \{ \mathbf{c} \mid \mathbf{c}_{l_x} = \mathbf{c}.\mathbf{c}', \mathbf{c}', \mathbf{c}'' \neq \epsilon \}$$

Notice that case **axis = ancestor-or-self** is analogous.

[**axis = following-sibling**] By definition we have that

$$[[\text{following-sibling}]]_{l_x}^\sigma = \text{FollowingSibling}_\sigma(l_x)$$

Let  $l_p$  be the parent of  $l_x$ , i.e.,  $l_p \in \text{Parent}_\sigma(l_x)$ . Since the document is valid wrt the DTD  $d$ , let  $a_x = \text{type}(l_x)$  and  $a_p = \text{type}(l_p)$  we have that  $a_p \Rightarrow_d a_x$ . Therefore  $c_{l_x} = c_{l_p}.a_x \in C_d$ . Since the document is valid wrt the DTD  $d$ , for  $l_f \in \text{FollowingSibling}_\sigma(l_x)$ , let  $a_f = \text{type}(l_f)$  we have that  $a_p \Rightarrow_d a_f$  and  $a <_{d(c_{l_p})} b$ . Therefore  $c_{l_f} = c_p.a_f \in C_d$ . We conclude since by definition

$$\mathbf{A}_{C_d}(c_{l_x}, \text{following-sibling}) = \{ c_p.a_f \in C_d \mid c_{l_x} = c_p.a_x \wedge a_x <_{d(c_{l_p})} a_f \}$$

Notice that case **axis = preceding-sibling** is analogous.

□

### Test chain inference

Soundness of chain inference for test filtering is stated by the following theorem. For simplicity, we present it together with completeness.

**Lemma 9.1.2** (Test Filtering Soundness and Completeness). *If  $\sigma, (x := l_x) \models x/\text{self} :: \phi \Rightarrow \sigma, L$  then*

$$\mathbf{T}_{C_d}(c_{l_x}^\sigma, \phi) = \{ c_{l_x}^\sigma \mid l_x \in L \}$$

*Proof.* By case analysis on  $\phi$ . Since **self** is used as navigational axis, we directly translate the judgment into the proper test filtering  $[[\phi]]_\sigma^{l_x}$ . Also note that the sequence  $L$  is equal to  $l_x$  alone when the test is evaluated to true and  $L$  is empty otherwise.

[ $\phi = a \in \Sigma$ ] By definition of node test filtering we have that  $[[a]]_\sigma^{l_x} = \{ l_x \mid l_x \leftarrow a[L] \in \sigma \}$  and we conclude immediately completeness since  $\mathbf{T}_{C_d}(c_{l_x}, a) = \{ c_{l_x} \mid c_{l_x} = c.a \}$ .

[ $\phi = \text{text}()$ ] As before, completeness follows by definition since  $[[a]]_\sigma^{l_x} = \{ l_x \mid l_x \leftarrow \text{"txt"} \}$  and  $\mathbf{T}_{C_d}(c_{l_x}, a) = \{ c_{l_x} \mid c_{l_x} = c.\text{String} \}$ .

[ $\phi = \text{node}()$ ] By definition  $[[\phi]]_\sigma^{l_x} = l_x$  and  $\mathbf{T}_{C_d}(c_{l_x}, \text{node}()) = \{ c_{l_x} \}$ .

□

### 9.1.3 Completeness of step chain inference

**Lemma 9.1.3** (Completeness of Step Chain Inference). *Let  $\mathbf{d}$  be a DTD and  $\mathbf{c} \in \mathbf{C}_{\mathbf{d}}$  a chain. If*

$$\mathbf{c}' \in \mathbf{T}_{\mathbf{C}_{\mathbf{d}}}(\mathbf{A}_{\mathbf{C}_{\mathbf{d}}}(\mathbf{c}, \mathbf{axis}), \mathbf{node}())$$

*then there exists a valid tree  $t = (\sigma, l_t) \in \mathbf{d}$  and a location  $l_x \in \text{dom}(t)$  such that  $\mathbf{c}_{l_x} = \mathbf{c}$  and, assuming that  $\sigma, (\mathbf{x} := l_x) \models \mathbf{x}/\mathbf{axis} :: \mathbf{node}() \Rightarrow \sigma, L$  we have that  $\mathbf{c}_{l'}^\sigma = \mathbf{c}'$ , for some  $l' \in L$ .*

*Proof.* As before, since  $\mathbf{node}()$  is used as test filtering, we directly translate the judgment into the proper axis evaluation and we leave as implicit the set-to-sequence conversion by mean of document order. Then we conclude by a case analysis on  $\mathbf{axis}$ .

[ $\mathbf{axis} = \mathbf{self}$ ] By Lemma 4.1.4, for all input chains  $\mathbf{c}$  then there exists a valid tree  $t = (\sigma, l_t) \in \mathbf{d}$  and a location  $l_x \in \text{dom}(t)$  such that  $\mathbf{c}_{l_x} = \mathbf{c}$ . We conclude immediately since by definition we have that  $\mathbf{A}_{\mathbf{C}_{\mathbf{d}}}(\mathbf{c}_{l_x}, \mathbf{self}) = \{ \mathbf{c}_{l_x} \}$  and  $\llbracket \mathbf{self} \rrbracket_{l_x}^\sigma = \{ l_x \}$ .

[ $\mathbf{axis} = \mathbf{child}$ ] By definition

$$\mathbf{A}_{\mathbf{C}_{\mathbf{d}}}(\mathbf{c}, \mathbf{child}) = \{ \mathbf{c.a} \mid \mathbf{c.a} \in \mathbf{C}_{\mathbf{d}} \}$$

By Lemma 4.1.4, for all input rooted chains  $\mathbf{c.a} \in \mathbf{C}_{\mathbf{d}}$  then there exists a valid tree  $t = (\sigma, l_t) \in \mathbf{d}$  and a location  $l \in \text{dom}(t)$  such that  $\mathbf{c}_l = \mathbf{c.a}$ . To conclude it suffices to fix  $\mathbf{c}_{l_x} = \mathbf{c}$ .

[ $\mathbf{axis} = \mathbf{descendant}$ ] By definition we have that

$$\mathbf{A}_{\mathbf{C}_{\mathbf{d}}}(\mathbf{c}, \mathbf{descendant}) = \{ \mathbf{c.c}' \mid \mathbf{c.c}' \in \mathbf{C}_{\mathbf{d}}, \mathbf{c}' \neq \epsilon \}$$

By Lemma 4.1.4, for all input rooted chains  $\mathbf{c.c}' \in \mathbf{C}_{\mathbf{d}}$  then there exists a valid tree  $t = (\sigma, l_t) \in \mathbf{d}$  and a location  $l \in \text{dom}(t)$  such that  $\mathbf{c}_l = \mathbf{c.c}'$ . To conclude it suffices to fix  $\mathbf{c}_{l_x} = \mathbf{c}$ .

Case [ $\mathbf{axis} = \mathbf{descendant-or-self}$ ] is analogous.

[ $\mathbf{axis} = \mathbf{parent}$ ] By definition we have that

$$\mathbf{A}_{\mathbf{C}_{\mathbf{d}}}(\mathbf{c.a}, \mathbf{parent}) = \{ \mathbf{c} \}$$

By Lemma 4.1.4, for all input rooted chains  $\mathbf{c} \in \mathbf{C}_{\mathbf{d}}$  then there exists a valid tree  $t = (\sigma, l_t) \in \mathbf{d}$  and a location  $l \in \text{dom}(t)$  such that  $\mathbf{c}_l = \mathbf{c}$ . To conclude it suffices to fix  $\mathbf{c}_{l_x} = \mathbf{c.a}$ .

[ $\mathbf{axis} = \mathbf{ancestor}$ ] By definition we have that

$$\mathbf{A}_{\mathbf{C}_{\mathbf{d}}}(\mathbf{c}, \mathbf{ancestor}) = \{ \mathbf{c}' \mid \mathbf{c}' \leq \mathbf{c}, \mathbf{c} \neq \epsilon \}$$

By Lemma 4.1.4, for all input rooted chains  $\mathbf{c.c}' \in \mathbf{C}_{\mathbf{d}}$  then there exists a valid tree  $t = (\sigma, l_t) \in \mathbf{d}$  and a location  $l \in \text{dom}(t)$  such that  $\mathbf{c}_l = \mathbf{c}$ . To conclude it suffices to fix  $\mathbf{c}_{l_x} = \mathbf{c.c}'$ .

Case  $\mathbf{ancestor-or-self}$  is analogous.



[**axis = following-sibling**] By definition we have that

$$\mathbf{A}_{C_d}(\mathbf{c}, \mathbf{a}, \text{following-sibling}) = \{ \mathbf{c}.b \in C_d \mid \mathbf{a} <_{d(\mathbf{c})} b \}$$

By Lemma 4.1.4, for all input rooted chains  $\mathbf{c} \in C_d$  then there exists a valid tree  $t = (\sigma, l_t) \in \mathbf{d}$  and a location  $l_p \in \text{dom}(t)$  such that  $\mathbf{c}_{l_p} = \mathbf{c}$ , and such that the sequence of children of  $l_p$  contains two locations  $l, l'$  such that  $\text{type}(l) = \mathbf{a}$  and  $\text{type}(l') = \mathbf{b}$  and  $\mathbf{a} <_{d(\mathbf{c}_{l_p})} \mathbf{b}$ .

Case **preceding-sibling** is analogous.

□

Putting everything together, we can prove the theorems of soundness and completeness for chain inference of XPath Step expressions.

LEMMA 4.3.1 (STEP CHAIN SOUNDNESS). Let  $\mathbf{d}$  be a DTD,  $t = (\sigma, l_t) \in \mathbf{d}$  be a valid tree, and  $l \in \text{dom}(\sigma)$  be a location of the tree. If

$$\sigma, (\mathbf{x} := l) \models \mathbf{x}/\mathbf{axis} :: \phi \Rightarrow \sigma, L$$

then for each location  $l' \in L$  we have that

$$\mathbf{c}_{l'}^\sigma \in \mathbf{T}_C(\mathbf{A}_C(\mathbf{c}_l^\sigma, \mathbf{axis}), \phi)$$

*Proof.* By Lemma 9.1.1 and Lemma 9.1.2. □

LEMMA 4.3.1 (STEP CHAIN COMPLETENESS). Let  $\mathbf{d}$  be a DTD and  $\mathbf{c} \in C_d$  a chain. If

$$\mathbf{c}' \in \mathbf{T}_{C_d}(\mathbf{A}_{C_d}(\mathbf{c}, \mathbf{axis}), \phi)$$

then there exists a valid tree  $t = (\sigma, l_t) \in \mathbf{d}$  and a location  $l \in \text{dom}(t)$  such that  $\mathbf{c}_l = \mathbf{c}$  and, assuming that  $\sigma, (\mathbf{x} := l) \models \mathbf{x}/\mathbf{axis} :: \phi \Rightarrow \sigma, L$  we have that  $\mathbf{c}_{l'}^\sigma = \mathbf{c}'$ , for some  $l' \in L$ .

*Proof.* By Lemma 9.1.3 and Lemma 9.1.2. □

#### 9.1.4 Soundness of Query Chain Inference

In this section we prove the correctness of chain inference for query expressions. We introduce first some auxiliary definitions and properties that we will use later in our formal development.

### Auxiliary definitions

The notion of correctness we stated (Theorem 4.3.4) relies on tree projections (Section 4.3.1 ).

Let  $\sigma$  be a store and  $L_1, L_2$  be two sequences of locations belonging to  $\text{dom}(\sigma)$ . The merge of  $L_1, L_2$ , denoted by  $L_1 \uplus L_2$ , is defined as the sequence  $L$  such that

- $l \in L_1 \cdot L_2$  iff  $l \in L$ ;
- for all  $l, l' \in L$  we have that  $l \neq l'$ ;
- $(l_i, l_{i+1}) \in L$  implies  $(l_i, l_{i+1}) \in \text{docOrder}_\sigma$ .

Let  $\sigma$  be a store and  $\sigma_1, \sigma_2$  be two of its projections. The stores obtained by only common and uncommon locations between  $\sigma_1$  and  $\sigma_2$  are defined as follows.

$$\text{common}(\sigma_1, \sigma_2) \stackrel{\text{def}}{=} \{ l \leftarrow \mathbf{a}[L_1 \uplus L_2] \mid l \leftarrow \mathbf{a}[L_1] \in \sigma_1, l \leftarrow \mathbf{a}[L_2] \in \sigma_2 \}$$

$$\text{uncommon}(\sigma_1, \sigma_2) \stackrel{\text{def}}{=} \{ l \leftarrow \mathbf{a}[L] \mid l \notin \text{dom}(\sigma_1) \cap \text{dom}(\sigma_2) \}$$

**Definition 9.1.4** (Merge of Projections). *Let  $\sigma$  be a store and  $\sigma_1, \sigma_2$  be two of its projections. The merge of  $\sigma_1$  and  $\sigma_2$ , denoted by  $\sigma_1 \uplus \sigma_2$ , is defined as*

$$\sigma_1 \uplus \sigma_2 \stackrel{\text{def}}{=} \text{common}(\sigma_1, \sigma_2) \cup \text{uncommon}(\sigma_1, \sigma_2)$$

*Let  $l \in \text{dom}(\sigma_1) \cap \text{dom}(\sigma_2)$  be a location, the merge of the trees  $t_1 = (\sigma_1, l)$  and  $t_2 = (\sigma_2, l)$ , denoted by  $t_1 \uplus t_2$ , is defined as  $t_1 \uplus t_2 \stackrel{\text{def}}{=} (\sigma_1 \uplus \sigma_2, l)$ .*

We outline some simple properties of the merge operator.

**Proposition 9.1.5.** *Given a tree  $t$  and two of its projections,  $t_1, t_2 \leq t$ , the following properties hold.*

1.  $t_1 \leq t_1 \uplus t_2$
2.  $t_1 \uplus t_2 \leq t$
3.  $t'_1 \leq t_1$  and  $t'_2 \leq t_2$  implies  $t'_1 \uplus t'_2 \leq t_1 \uplus t_2$

*Proof.* 1) Follows from the definition of merge.

2) Follows from the definition of merge and from  $t_i \leq t$ .

3) Follows since, by the point 1), we have that  $t'_1 \leq t_1 \uplus t_2$  and  $t'_2 \leq t_1 \uplus t_2$ . Therefore  $t'_1 \uplus t'_2 \leq t_1 \uplus t_2$ .  $\square$

We prove now an auxiliary lemma concerning return chains.

**Proposition 9.1.6.** *Let  $\mathbf{d}$  be a DTD,  $t = (\sigma, l_t) \in \mathbf{d}$  a valid tree,  $\mathbf{q}$  a query,  $\gamma$  a dynamic environment and  $\Gamma$  a static environment such that  $\sigma \models_{\mathbf{d}} \gamma : \Gamma$ . Provided  $\sigma, \gamma \models \mathbf{q} \Rightarrow \sigma_{new}, L$  and  $\Gamma \vdash_{C_d} \mathbf{q} : (r, v, e)$ , if  $l \in L \cap \text{dom}(t)$  then  $c_l^\sigma \in r_{\mathbf{q}}$ .*

*Proof.* By induction on the structure of  $\mathbf{q}$ .

*Base.*

$[\mathbf{q} = ()]$  Immediate since  $r_{\mathbf{q}} = \emptyset$ .

$[\mathbf{q} = \text{"txt"}]$  Immediate since  $l \notin \text{dom}(t)$ .

*Induction*

$[\mathbf{q} = \mathbf{q}_1, \mathbf{q}_2]$  By induction on  $\mathbf{q}_1, \mathbf{q}_2$  and by noticing that  $r_{\mathbf{q}} = r_{\mathbf{q}_1} \cup r_{\mathbf{q}_2}$ .

$[\mathbf{q} = \langle a \rangle \mathbf{q}' \langle /a \rangle]$  Immediate since  $l \notin \text{dom}(t)$ .

$[\mathbf{q} = \mathbf{x}/\text{step}]$  By correctness of step chain inference (Theorem 4.3.1), and since resulting chains are in  $r_{\mathbf{q}}$ .

$[\mathbf{q} = \text{for } \mathbf{x} \text{ in } \mathbf{q}_1 \text{ return } \mathbf{q}_2]$  First we notice that  $L_{\mathbf{q}} = L_{\mathbf{q}_2}$ . We conclude by inductive hypothesis on  $\mathbf{q}_2$  and because  $r_{\mathbf{q}} = r_{\mathbf{q}_2}$ .

$[\mathbf{q} = \text{let } \mathbf{x} := \mathbf{q}_1 \text{ return } \mathbf{q}_2]$  Analogous to the **for** case.

$[\mathbf{q} = \text{if } (\mathbf{q}_0) \text{ then } \mathbf{q}_1 \text{ else } \mathbf{q}_2]$  First we notice that either  $L_{\mathbf{q}} = L_{\mathbf{q}_1}$  or  $L_{\mathbf{q}} = L_{\mathbf{q}_2}$ . We conclude by inductive hypothesis on  $\mathbf{q}_1, \mathbf{q}_2$  and because  $r_{\mathbf{q}} = r_{\mathbf{q}_1} \cup r_{\mathbf{q}_2}$ .

□

**Theorem 4.3.4** ((Soundness of Query Chains)) *Let  $\mathbf{d}$  be a DTD,  $t = (\sigma, l_t) \in \mathbf{d}$  a valid tree,  $\mathbf{q}$  a query,  $\gamma$  a dynamic environment and  $\Gamma$  a static environment such that  $\sigma \models_{\mathbf{d}} \gamma : \Gamma$ . Provided  $\sigma, \gamma \models \mathbf{q} \Rightarrow \sigma_{new}, L$  and  $\Gamma \vdash_{\mathbf{C}_d} \mathbf{q} : (r, v, e)$  the following holds.*

1. *if  $t^{min}$  is a minimal  $\mathbf{q}$ -projection of  $t$  then  $t^{min} \leq t|_{\mathcal{L}_{r \cup v}^t} \leq t$*
2. *if  $t^{new} = (\sigma_{new}, l_{new})$  with  $l_{new} \in L \setminus \text{dom}(\sigma)$  is a fresh subtree then  $\text{dom}(t^{new}) = \mathcal{L}_e^{t^{new}}$*

For the sake of conciseness we adopt the following abbreviations. The sequence  $L_{\mathbf{q}}$  denotes the result  $\mathbf{q}$  evaluation, that is  $\sigma, \gamma \models \mathbf{q} \Rightarrow \sigma_{\mathbf{q}}, L_{\mathbf{q}}$ . The sets of chains  $r_{\mathbf{q}}, v_{\mathbf{q}}, e_{\mathbf{q}}$  denote the sets of return, used and element chains inferred for  $\mathbf{q}$ , that is  $\Gamma \vdash_{\mathbf{C}_d} \mathbf{q} : (r_{\mathbf{q}}, v_{\mathbf{q}}, e_{\mathbf{q}})$ .

We prove the point 1.

*Proof.* By structural induction on  $\mathbf{q}$ .

*Base.*

$[\mathbf{q} = (), \mathbf{q} = \text{“txt”}]$  We conclude immediately since the  $\mathbf{q}$ -projection of  $t$  is empty.

*Induction.*

$[\mathbf{q} = \mathbf{q}_1, \mathbf{q}_2]$  A minimal  $\mathbf{q}$ -projection  $t'$  is the merge of a minimal  $\mathbf{q}_1$ -projection  $t'_1$  with a minimal  $\mathbf{q}_2$ -projection  $t'_2$ . This follows from the semantics rule (SQ-CONCAT). By inductive hypothesis we have that

$$t'_1 \leq t|_{\bar{r}_{\mathbf{q}_1} \cup v_{\mathbf{q}_1}} \quad \text{and} \quad t'_2 \leq t|_{\bar{r}_{\mathbf{q}_2} \cup v_{\mathbf{q}_2}}$$

By Proposition 9.1.5

$$t' = t'_1 \uplus t'_2 \leq t|_{\bar{r}_{\mathbf{q}_1} \cup v_{\mathbf{q}_1}} \uplus t|_{\bar{r}_{\mathbf{q}_2} \cup v_{\mathbf{q}_2}}$$

By chain inference rule (CONCAT) we have that  $r_{\mathbf{q}} = r_{\mathbf{q}_1} \cup r_{\mathbf{q}_2}$  and  $v_{\mathbf{q}} = v_{\mathbf{q}_1} \cup v_{\mathbf{q}_2}$ . Therefore we conclude

$$t' \leq t'_1 \uplus t'_2 \leq t|_{\bar{r}_{\mathbf{q}} \cup v_{\mathbf{q}}}$$

$[\mathbf{q} = \langle a \rangle \mathbf{q}' \langle /a \rangle]$  In this case, a tree projection is minimal for  $\mathbf{q}$  if and only if it is minimal for  $\mathbf{q}'$ . This follows from the fact that element construction always produces fresh locations that do not belong to the original store. This is specified by semantics rules (SQ-ELT) and (SQ-COPY). By inductive hypothesis, if  $t'$  is a minimal projection for  $\mathbf{q}'$  then  $t' \leq t|_{\bar{r}_{\mathbf{q}'} \cup v_{\mathbf{q}'}}$ . Following rule (ELT) we have that  $r_{\mathbf{q}} = \emptyset$  and

$$v_{\mathbf{q}} = \bar{r}_{\mathbf{q}'} \cup v_{\mathbf{q}'}$$

We conclude since  $t' \leq t|_{v_{\mathbf{q}}}$ .

[ $q = \mathbf{x}/\text{step}$ ] Semantics of step navigation is defined by rule (SQ-STEP). We assume that  $\mathbf{x}$  is bound to a single location  $l_x$  in  $\gamma$ . The general case where  $\mathbf{x}$  is bound to a sequence of locations goes by induction and concludes this case.

We distinguish between two main cases depending on the axis steps being forward or not.

**(axis is forward)** Since **axis** is forward, for all location  $l \in L_q$  we have that either  $l$  is a descendant of  $l_x$ , or  $l = l_x$  if **axis** = **self**. Let  $t'$  be a minimal  $q$ -projection of  $t$ . Projection  $t'$  allows to correctly evaluate  $q$ , hence it is given by all locations in  $L_q$  together with their ancestors and their descendants. Formally,  $t' = (\sigma|_{\mathcal{L}}, l_t)$  where

$$\mathcal{L} = \bigcup_{l \in L_q} \{l\} \cup \llbracket \text{ancestor} \rrbracket_l^\sigma \cup \llbracket \text{descendant} \rrbracket_l^\sigma$$

Notice that  $t'$  is both minimal and unique. Notice also that it may be the case that  $L_q = ()$  and thus  $t'$  is the empty forest. Because  $\sigma \models_d \gamma : \Gamma$  we have that  $c_{l_x} \in \Gamma(\mathbf{x})$ . By Theorem 4.3.1 step chain inference is correct. Therefore, following rule (STEPF) we can conclude that if  $l \in L_q$  then  $c_l \in r_q$ . Notice also that  $c_l$  types any location  $l'$  that is an ancestors of  $l$ . The descendants of  $l$  are typed by the prolongation of  $c_l$ , that belongs to  $\bar{r}_q$ . By rule (STEPF) no used chain is inferred. Therefore we conclude

$$t' \leq t|_{\bar{r}_q}$$

**(axis is backward or horizontal)** In this case, since the axis is either backward or horizontal, all resulting locations  $l \in L_q$  are not descendants of  $l_x$ . Notice however that it may be the case that  $l = l_x$  if **axis** = **ancestor-or-self**. Let  $t'$  be a minimal  $q$ -projection of  $t$ . Because  $t'$  allows to correctly evaluate  $q$ , it must contain: *i*) all locations in the result sequence  $L_q$ , together with their respective ancestors and descendants and *ii*) the input location  $l_x$  together with its ancestors, but only in the case that  $L_q \neq ()$ . We stress that we may not need to include the descendants of  $l_x$  in the projection, for instance if we are evaluating an horizontal axis. Formally,  $t' = (\sigma|_{\mathcal{L}}, l_t)$  where  $\mathcal{L}$  is

$$\bigcup_{l \in L_q} \{l\} \cup \llbracket \text{ancestor} \rrbracket_l^\sigma \cup \llbracket \text{descendant} \rrbracket_l^\sigma \cup \{l_x \mid L_q \neq ()\}$$

Notice that for backward and horizontal axis we have  $\llbracket \text{ancestor} \rrbracket_l^\sigma \subseteq \llbracket \text{ancestor} \rrbracket_{l_x}^\sigma$ . Also,  $t'$  is unique. Because  $\sigma \models_d \gamma : \Gamma$  we have that  $c_{l_x} \in \Gamma(\mathbf{x})$ . By Theorem 4.3.1 step chain inference is correct. Therefore, for all  $l \in L_q$  we have that  $c_l \in r_q$ . Furthermore, if  $L \neq ()$  then  $c_{l_x} \in v_q$ . We conclude that

$$t' \leq t|_{\bar{r}_q \cup v_q}$$

When  $\gamma(x) = L$  query evaluation is performed by means of iteration, as defined by rules (SQ-ITER) and (SQ-ITERBASE). Let  $t'$  be the minimal  $\mathbf{q}$ -projection for  $\gamma = \{ l \in L \}$ . By inductive hypothesis, chain inference is correct for each  $l \in L$ . Let  $t'_l$  be the minimal query  $\mathbf{q}$ -projection for  $\gamma = \{ l \}$ . We conclude since by Lemma 9.1.5

$$t' = \biguplus_{l \in L} t'_l \leq t_{\bar{r}_q \cup v_q}$$

$[\mathbf{q} = \text{for } \mathbf{x} \text{ in } \mathbf{q}_1 \text{ return } \mathbf{q}_2]$  Semantics of iteration is defined by rule (SQ-FOR). Query  $\mathbf{q}_2$  is evaluated once for each location  $l \in L_{\mathbf{q}_1}$  binding  $\mathbf{x}$  with  $l$  in  $\gamma$ . We denote by  $t'_2$  a minimal  $\mathbf{q}_2$ -projection of  $t$  obtained by assuming that  $\mathbf{x} \mapsto l \in \gamma$ . We denote by  $\mathcal{L}_1$  the set of result locations of  $\mathbf{q}_1$  that are also productive for  $\mathbf{q}_2$ , that is  $\mathcal{L}_1 = \{ l \in L_1 \mid t_2^l \neq () \}$ .

A minimal  $\mathbf{q}$ -projection of  $t$  is defined as  $t' = t_1 \uplus t_2$  where  $t_1$  is a minimal  $\mathbf{q}_1$ -projection of  $t$  that ensures only locations in  $\mathcal{L}_1$  in the result sequence of  $\mathbf{q}_1$ , and  $t_2$  is a minimal  $\mathbf{q}_2$ -projection of  $t$  obtained by the merge of all  $\mathbf{q}_2$ -minimal projections of  $t$  where  $\mathbf{x}$  is bound to a location  $l \in \mathcal{L}_1$ .

We want to show that

$$t' = t_1 \uplus t_2 \leq t_{\bar{r}_q \cup v_q}$$

We show first  $t_2 \leq t_{\bar{r}_q \cup v_q}$ . By inductive hypothesis, provided that  $l_i \in \mathcal{L}_1$  and  $\Gamma[\mathbf{x} \mapsto \{c_{l_i}\}] \vdash_C \mathbf{q}_2 : (r_i, v_i, e_i)$  it follows that

$$t_2^l \leq t_{\bar{r}_i \cup v_i}$$

By Lemma 9.1.5 we have also that

$$t_2 = \biguplus_{l \in \mathcal{L}_1} t_2^l$$

In general, if the result of  $\mathbf{q}_2$  is non-empty then  $r_i \cup e_i \neq \emptyset$ .

By rule (FOR)

$$r_q = \bigcup_{c_{l_i} \in r_{q_1}} r_i \quad \text{and} \quad v_q \supseteq \bigcup_{r_i \cup e_i \neq \emptyset} \{v_i\}$$

from which it follows

$$t_2 \leq t_{\bar{r}_q \cup v_q}$$

We show  $t_1 \leq t_{v_q}$ . By inductive hypothesis,

$$t_1 \leq t_{\bar{r}_{q_1} \cup v_{q_1}}$$

however since  $t'$  does not need the descendants of result locations of  $\mathbf{q}_1$ , we have also that  $t_1 \leq t_{r_{q_1} \cup v_{q_1}}$ . Then, because we need only result locations in  $\mathcal{L}_1$  we have that  $t_1 \leq t_{r'_1 \cup v_{q_1}}$  where

$$r'_1 = \bigcup_{\substack{c_{l_i} \in r_{q_1} \\ r_i \cup e_i \neq \emptyset}} \{c_{l_i}\}$$

By rule (FOR) we have that

$$r_{q_1} \cup v_{q_1} \subseteq v_q$$

and therefore we conclude

$$t_1 \leq t_{|v_q}$$

Putting everything together, by Lemma 9.1.5

$$t_1 \uplus t_2 \leq t_{|\bar{r}_q \cup v_q}$$

[ $q = \text{let } x := q_1 \text{ return } q_2$ ] The semantics of the let binding is defined by rule (SQ-LET). Let  $t_i$  be a minimal  $q_i$ -projection for  $q_i$ . Notice that in the evaluation of  $q_2$ ,  $x$  is bound to the resulting sequence of  $q_1$  in  $\gamma$  (no iteration is performed by the rule).

By inductive hypothesis

$$t_1 \leq t_{|\bar{r}_{q_1} \cup v_{q_1}} \quad \text{and} \quad t_2 \leq t_{|\bar{r}_{q_2} \cup v_{q_2}}$$

Following rule (LET) we have that

$$r_{q_1} \cup v_{q_1} \subseteq v_q$$

Therefore

$$t_1 \leq t_{|v_q}$$

Furthermore, by rule (LET) we have that

$$\bar{r}_{q_2} \cup v_{q_2} \subseteq \bar{r}_q \cup v_q$$

Therefore

$$t_2 \leq t_{|\bar{r}_q \cup v_q}$$

Now  $t' = t_1 \uplus t_2$  is a minimal  $q$ -projection of  $t$ . By Proposition 9.1.5 we conclude

$$t_1 \uplus t_2 \leq t_{|\bar{r}_q \cup v_q}$$

[ $q = \text{if } (q_0) \text{ then } q_1 \text{ else } q_2$ ] The semantics of the conditional expression is defined by rules (SQ-IF). Let  $t', t_0, t_1$  and  $t_2$  minimal projection of  $t$  for  $q, q_0, q_1$  and  $q_2$ , respectively. By inductive hypothesis

$$t_i \leq t_{|\bar{r}_{q_i} \cup v_{q_i}} \quad \text{for } i = 0..2$$

According to the semantics of the conditional expression

$$t' \leq \biguplus_{i=0..2} t_i$$

and then we conclude by rule (IF) since

$$\biguplus_{i=0..2} t_i \leq t_{|\bar{r}_q \cup v_q}$$

□

We prove the point 2.

*Proof.* By structural induction on  $q$ .

*Base.*

$[q = ()]$  The empty query does not construct any element.

$[q = \text{"txt"}]$  The query constructs a fresh text-node and by rule (TEXT) we have that  $\text{String} \in e$ .

*Induction*

$[q = q_1, q_2]$  The constructed elements returned by  $q$  are obtained by the concatenation of those returned by  $q_1$  and  $q_2$ , respectively. By inductive hypothesis the statement holds for  $q_1, q_2$  and since  $e_q = e_{q_1} \cup e_{q_2}$  we conclude.

$[q = \langle a \rangle q' \langle /a \rangle]$  Following the rule for query semantics (SQ-CONSTR), the query  $q$  adds to the input store  $\sigma$  a new location, say  $l_a$ , such that  $(l_a \leftarrow a[L_{q'}^{\text{copy}}])$ , where  $L_{q'}^{\text{copy}}$  denotes the copy of  $L_{q'}$ . The resulting store is denoted by  $\sigma_q$ . We distinguish four cases depending on whether the subexpression  $q'$  returns *i*) only locations belonging to the input tree, *ii*) only constructed elements, *iii*) both constructed and input elements, or *iv*) the empty sequence.

*i*) We assume that  $l \in L_{q'}$  implies  $l \in \text{dom}(t)$ , where  $t = (\sigma, l_t)$  denotes the input tree. By Proposition 9.1.6 we know that for all  $l \in L_{q'}$  we have  $c_l \in r_{q'}$ . By rule (ELT), if  $c_l = c.b$  and  $c.b.c' \in \bar{r}_{q'}$  then  $a.b.c' \in e_q$ . Therefore let  $t' = \sigma_q @ l_a$  we conclude

$$t' \leq t'_{|e_q}$$

*ii*) We assume that  $l \in L_{q'}$  implies  $l \notin \text{dom}(t)$ , where  $t = (\sigma, l_t)$  denotes the input tree. By inductive hypothesis, for all tree  $t'' = \sigma_q @ l_a$ , with  $l \in L_{q'}$ , we have that

$$t'' \leq t''_{|e_{q'}}$$

By rule (ELT), if  $c \in e_{q'}$  then  $a.c \in e_q$ . Let  $t' = \sigma_q @ l_a$ , we conclude that all locations of  $t'$  are typed by a prefix of a chain in  $e_q$ . Formally we have

$$t' \leq t'_{|e_q}$$

*iii*) In the case where  $q'$  provides both input and fresh data, we conclude by reasoning as in *i*) and *ii*).

*iv*) We assume  $L_{q'} = ()$ . In this case, the result of  $q$  is a tree composed by a single node labeled  $a$ . By inductive hypothesis we assume the system to be *correct* on  $q'$ , but it may not be complete, thus we do not know if  $r_{q'} \cup e_{q'} = \emptyset$  or not. So we distinguish between the two cases. If  $r_{q'} \cup e_{q'} = \emptyset$  (hence the system is



also complete for  $q'$ ) we conclude since by rule (ELT)  $a \in e_q$ . If  $r_{q'} \cup e_{q'} \neq \emptyset$  (hence the system is correct but not complete for  $q'$ ) then rule (ELT) produces element chains as described in *i*) and *ii*). Since all these chains will start with label  $a$ , we can conclude that the resulting element of  $q$  is always typed.

[ $q = x/\text{step}$ ] No element is constructed by the query.

[ $q = \text{for } x \text{ in } q_1 \text{ return } q_2$ ] The elements constructed and returned by  $q$  are those constructed and returned by  $q_2$ . We conclude by inductive hypothesis on  $q_2$ .

[ $q = \text{let } x := q_1 \text{ return } q_2$ ] Analogous to the **for** case.

[ $q = \text{if } (q_0) \text{ then } q_1 \text{ else } q_2$ ] The elements constructed and returned by  $q$  are those constructed and returned either by  $q_1$  or  $q_2$ . We conclude by inductive hypothesis on  $q_1$  and  $q_2$ .

□

Summing up, we proved that query chain inference is sound. This means two things. First, that locations in the input document accessed by the query are typed by some inferred chains. Second, that locations created and returned by the query are typed by some inferred chains. Completeness (whether inference yields to minimal sets of chains) is still open, and we leave it as a future work.

### 9.1.5 Soundness of Update chain Inference

In this section we prove that chain inference for update expressions is correct. As updates are defined in terms of queries, the proof relies on correctness of chain inference for query expressions (Theorem 4.3.4). Chain inference for updates is sound in the sense that all chains required to type locations involved by an update are inferred.

**THEOREM 4.3.6 (SOUNDNESS OF UPDATE CHAIN INFERENCE)** *Let  $\mathbf{d}$  be a DTD,  $t$  a valid tree  $t = (\sigma, l_t) \in \mathbf{d}$ , and  $\mathbf{u}$  an update with at most one free variable  $\mathbf{x}$ ,  $\gamma = \{ \mathbf{x} \mapsto l_t \}$  a dynamic environment, and  $\Gamma = \{ \mathbf{x} \mapsto \mathbf{s}_d \}$  a static environment. Provided that*

$$\sigma, \gamma \models \mathbf{u} \Rightarrow \sigma_\omega, \omega \quad \sigma_\omega \models \omega \rightsquigarrow \sigma_{\mathbf{u}} \quad \Gamma \vdash_{C_d} \mathbf{u} : \mathbf{U}$$

*it holds that*

- *if  $l \in \text{dom}(\sigma)$  is a location in  $t$ , and the update  $\mathbf{u}$  involves  $l$  then there exists  $\mathbf{c} : \mathbf{c}' \in \mathbf{U}$  such that  $\mathbf{c}_l^\sigma = \mathbf{c} : \mathbf{c}'$ , where  $\mathbf{c}' \neq \epsilon$ .*
- *if  $l$  is a location in  $\mathbf{u}(t)$ , i.e.  $l \in \text{dom}(\sigma_{\mathbf{u}} @ l_t)$ , and the update  $\mathbf{u}$  involves  $l$  then there exists  $\mathbf{c} : \mathbf{c}' \in \mathbf{U}$  such that  $\mathbf{c}_l^{\sigma_{\mathbf{u}}} = \mathbf{c} : \mathbf{c}'$ , where  $\mathbf{c}' \neq \epsilon$ .*

*Proof.* By structural induction on  $\mathbf{u}$ .

*Base.*

$[\mathbf{u} = ()]$  No location is involved by the empty update. As defined by rules (SU-EMPTY) and (SAU-EMPTY),  $\mathbf{u}$  yields an empty update pending list, that performs no operation. Notice also that following rule (EMPTY) no chain is inferred.

$[\mathbf{u} = \text{delete } \mathbf{q}_0]$  Update  $\mathbf{u}$  involves all locations in  $L_{\mathbf{q}_0}$  (that is the resulting sequence of  $\mathbf{q}_0$ ). As defined by rule (SU-DELETE), there should be exactly one location in  $L_{\mathbf{q}_0}$ . We assume  $L_{\mathbf{q}_0} = (l_0)$ , with  $l_0 \in \text{dom}(t)$ . By rule (SAU-DELETE)  $l_0$  is detached from its parent. By rule (SANITYCHECK) this operation is correct since  $l_0$  is mutable and thus can be updated. By Proposition 9.1.6  $\mathbf{c}_{l_0}^\sigma \in \mathbf{r}_{\mathbf{q}_0}$ . The statement follows by definition since by following rule (DELETE) we have that given  $\mathbf{c}_{l_0}^\sigma = \mathbf{c}.\mathbf{a}_0$  then  $\mathbf{c}:\mathbf{a}_0 \in \mathbf{U}$ .

$[\mathbf{u} = \text{rename } \mathbf{q}_0 \text{ as } \mathbf{b}]$  Update  $\mathbf{u}$  involves all locations in  $L_{\mathbf{q}_0}$ . As defined by rule (SU-RENAME) there is at most one location  $l_0 \in \text{dom}(t)$  in the resulting sequence of  $\mathbf{q}_0$ . We assume  $L_{\mathbf{q}_0} = (l_0)$ . By rule (SAU-RENAME) the label of  $l_0$  is renamed in label  $\mathbf{b}$ . This is correct since by rule (SANITYCHECK)  $l_0$  is mutable, and thus can be updated, and the location is not the target of several **rename** expressions. By Proposition 9.1.6  $\mathbf{c}_{l_0}^\sigma \in \mathbf{r}_{\mathbf{q}_0}$ . In this case, we have that the renamed location belongs both to the original and to the updated store. We conclude now that in both cases the location is typed. Let  $\mathbf{c}_{l_0}^\sigma = \mathbf{c}.\mathbf{a}_0$  by rule (RENAME) we have  $\mathbf{c}:\mathbf{a}_0 \in \mathbf{U}$ . Let  $\mathbf{c}_{l_0}^{\sigma_{\mathbf{u}}} = \mathbf{c}.\mathbf{b}$  by rule (RENAME) we have that  $\mathbf{c}:\mathbf{b} \in \mathbf{U}$ .

$[u = \text{insert } q \text{ (into|into as first|into as last) } q_0]$  The locations involved by  $u$  are copies of source locations in  $L_q$ , together with their descendants. As defined by rules (SAU-INSERTINTO), (SAU-INSERTFIRST) and (SAU-INSERTLAST), these locations are in turn inserted as children of location  $L_{q_0}$ . By rule (SU-INSERT1) there is at most one target location in the resulting sequence of  $q_0$ , hence  $L_{q_0} = (l_0)$ . Source locations in  $L_q$  are typed by inferred chains, as stated by Theorem 4.3.4. We recall this fact, distinguishing between fresh and input source locations.

- If  $l \in L_q$  belongs to the input document then by Proposition 9.1.6 we have that  $c_l^\sigma \in r_q$ . Furthermore, for all  $l'$  descendant of  $l$  we have  $c_{l'}^\sigma \in \bar{r}_q$ .
- If  $l \in L_q$  is a fresh location, by Theorem 4.3.4, we have that  $l$  is typed by some chains in  $e_q$ , formally  $c_l^{\sigma_w}.e \in e_q$ . Furthermore, for all  $l'$  descendant of  $l$  we have  $c_{l'}^{\sigma_w}.e' \in e_q$ . Recall that  $\sigma_w$  is an intermediate store containing source locations that are arguments of elementary update commands, to be applied further.

Locations involved by  $u$  belongs to the updated store  $\sigma_u$  only. We conclude showing that chains typing these locations are in  $U$ . By Proposition 9.1.6  $c_{l_0}^\sigma \in r_{q_0}$ . By rule (INSERT1) we have that

- if  $c.a \in r_q$  and  $c.a.c' \in \bar{r}_q$  then  $c_{l_0}^\sigma : a.c' \in U$
- if  $e \in e_q$  then  $c_{l_0}^\sigma : e \in U$

$[u = \text{insert } q \text{ (after|before) } q_0]$  This case is similar to the previous one, except for the fact that chains for after/before expressions have a different definition. The locations involved by  $u$  are copies of source locations in  $L_q$ , together with their descendants. As defined by rules (SAU-INSERTAFTER), and (SAU-INSERTBEFORE), these locations are in turn inserted as children of location  $L_{q_0}$ . By rule (SU-INSERT2) there is at most one target location in the resulting sequence of  $q_0$ , hence  $L_{q_0} = (l_0)$ . Source locations in  $L_q$  are typed as stated by Theorem 4.3.4. We recall this fact, distinguishing between fresh and input source locations.

- If  $l \in L_q$  belongs to the input document then By Proposition 9.1.6 we have that  $c_l^\sigma \in r_q$ . Furthermore, for all  $l'$  descendant of  $l$  we have  $c_{l'}^\sigma \in \bar{r}_q$ .
- If  $l \in L_q$  is a fresh location, by Theorem 4.3.4 we have that  $l$  is typed by some chains in  $e_q$ , formally  $c_l^{\sigma_w}.e \in e_q$ . Furthermore, for all  $l'$  descendant of  $l$  we have  $c_{l'}^{\sigma_w}.e' \in e_q$ . Recall that  $\sigma_w$  is an intermediate store containing source locations that are arguments of elementary update commands not yet applied.

Locations involved by  $u$  belongs to the updated store  $\sigma_u$ . We conclude showing that chains for those locations are in  $U$ . By Proposition 9.1.6  $c_{l_0}^\sigma \in r_{q_0}$ . Now we can conclude since by rule (INSERT1), if  $c_{l_0}^\sigma = c.a$  we have that

- if  $c.a \in r_q$  and  $c.a.c' \in \bar{r}_q$  then  $c : c' \in U$

- if  $e \in e_q$  then  $c : e \in U$

$[u = \text{replace } q_0 \text{ with } q]$  In this case soundness follows as for delete and insert cases. The locations involved by  $u$  are copies of source locations in  $L_q$  together with their descendants, and target locations returned by  $q_0$ . By rule (SU-REPLACE) there is only one target location  $l_0$  in the resulting sequence of  $q_0$ , hence  $L_{q_0} = (l_0)$ . We first treat the typing of source locations and the typing of target locations. Source locations in  $L_q$  are typed as stated by Theorem 4.3.4. We recall this fact, distinguishing between fresh and input source locations.

- If  $l \in L_q$  belongs to the input document then By Proposition 9.1.6 we have that  $c_l^\sigma \in r_q$ . Furthermore, for all  $l'$  descendant of  $l$  we have  $c_{l'}^\sigma \in \bar{r}_q$ .
- If  $l \in L_q$  is a fresh location, by Theorem 4.3.4 we have that  $l$  is typed by some chains in  $e_q$ , formally  $c_l^{\sigma_w}.e \in e_q$ . Furthermore, for all  $l'$  descendant of  $l$  we have  $c_{l'}^{\sigma_w}.e' \in e_q$ . Recall that  $\sigma_w$  is an intermediate store containing source locations that are arguments of elementary update commands not yet applied.

By Proposition 9.1.6  $c_{l_0}^\sigma \in r_{q_0}$ . Now we can conclude that source locations are typed since by rule (REPLACE) let  $c_{l_0}^\sigma = c.a$  we have that

- if  $c.a \in r_q$  and  $c.a.c' \in \bar{r}_q$  then  $c : a.c' \in U$
- if  $e \in e_q$  then  $c : e \in U$

Concerning target locations, the statement follows since let  $c_{l_0}^\sigma = c.a$  we have that  $c : a \in U$ .

*Induction.*

$[u = u_1, u_2]$  The locations involved by the concatenation of two updates is the union of the locations involved by each expression. We conclude by inductive hypothesis on  $u_1, u_2$ .

$[u = \text{for } x \text{ in } q \text{ return } u_1]$  The locations involved by the update are those involved by  $u_1$  when  $x$  is bound to each node in the resulting sequence of  $q$ . By Proposition 9.1.6 the chain corresponding to all return location of  $q$  are inferred in  $r_q$ . By inductive hypothesis the statement holds for the update  $u_1$  when  $x$  is bound to a return chain inferred for  $q$ . This is described by rule (FOR), and therefore we conclude that the statement holds for  $u$ .

$[u = \text{let } x := q \text{ return } u_1]$  The locations involved by the update are those involved by  $u_1$  when applied with  $x$  bound to the resulting sequence of  $q$ . By Proposition 9.1.6 the chain corresponding to all return location of  $q$  are inferred in  $r_q$ . By inductive hypothesis the statement holds for the update  $u_1$  when  $x$  is bound to the resulting sequence of  $q$ .

This is described by rule (LET), and therefore we conclude that the statement holds for  $u_1$ .

$[u = \text{if } (q) \text{ then } u_1 \text{ else } u_2]$  The chains involved by  $u$  are those involved by  $u_1$  or  $u_2$ . By inductive hypothesis the statement holds for  $u_1, u_2$  and we conclude that it holds also for  $u$ .

□

In this section we proved the soundness of update chain inference. We showed that involved locations are always typed by some update chains, and that these update chains are correctly inferred. As for queries, completeness remains open and it is left as a future work. It is worth noticing however that completeness of chain inference for updates requires (as a necessary condition) chain inference for queries to be complete.

### 9.1.6 Soundness of Chain-based Independence

**Proposition 9.1.7.** *If  $q \perp_{C_d} u$ , then we have:*

$$\mathcal{I}_U^t \cap \mathcal{L}_{\bar{r} \cup v}^t = \mathcal{I}_U^{u(t)} \cap \mathcal{L}_{\bar{r} \cup v}^{u(t)} = \emptyset$$

*Proof.* We treat separately parts of the statement concerning  $t$  and  $u(t)$ .

We show  $\mathcal{I}_U^t \cap \mathcal{L}_{\bar{r} \cup v}^t = \emptyset$ , aiming at obtaining a contradiction. Assume there is an involved location  $l_0$  belonging to a minimal  $q$ -projection, hence  $l \in \mathcal{I}_U^t \cap \mathcal{L}_{\bar{r} \cup v}^t$ . We recall that the set of nodes of  $t$  typed by update chains in  $U$  is defined as

$$\mathcal{I}_U^t = \{ l \in \text{dom}(t) \mid c_l = c : c' \in U \quad c' \neq \epsilon \}$$

and the set of nodes in a tree  $t$  typed by chains in  $\bar{r} \cup v$  is defined as

$$\mathcal{L}_{\bar{r} \cup v}^t = \{ l \in \text{dom}(t) \mid c_{l_0}^\sigma \cdot c \in \bar{r} \cup v \}$$

From these definitions we conclude that  $c_{l_0}^\sigma \in U$  and  $c_{l_0}^\sigma \in \bar{r} \cup v$ . Hence,

$$U \cap (\bar{r} \cup v) \neq \emptyset$$

By hypothesis  $q \perp_{C_d} u$  and then

$$\text{confl}(r, U) = \text{confl}(U, r) = \text{confl}(U, v) = \emptyset$$

Assume  $c : \alpha.c' \in U$ . Since  $\text{confl}(r, U) = \emptyset$  there cannot exist  $c_q \in r$  such that  $c_q \leq c : \alpha.c'$ . Since  $\text{confl}(U, r) = \emptyset$  there cannot exist  $c_q \in r$  such that  $c.\alpha.c' \leq c_q$ . Hence we conclude that  $c.\alpha.c' \notin \bar{r}$ . Finally, since  $\text{confl}(U, v) = \emptyset$  there cannot exist  $c_q \in v$  such that  $c \leq c_q$ . Hence we conclude  $c.\alpha.c' \notin v$ . Therefore from  $q \perp_{C_d} u$  follows

$$U \cap (\bar{r} \cup v) = \emptyset$$

a contradiction.

The same reasoning holds for proving  $\mathcal{I}_U^{u(t)} \cap \mathcal{L}_{\bar{r} \cup v}^{u(t)} = \emptyset$ .

□

**Theorem 4.4.2** (SOUNDNESS OF  $C_d$  INDEPENDENCE) *Let  $d$  be a DTD,  $t = (\sigma, l_t) \in d$  a valid tree,  $q$  a query,  $u$  an update,  $\gamma$  a dynamic environment and  $\Gamma$  a static environment such that  $\sigma \models_d \gamma : \Gamma$ . Then,*

$$q \perp_{C_d} u \quad \text{implies} \quad q \perp_d u$$

*Proof.* If  $q \perp_{C_d} u$  then by Proposition 9.1.7 we have that locations involved by the update are disjoint from location in a query projection. Since a query projection ensures preservation of query semantics under evaluation we can conclude  $q \perp_d u$ .

□

## 9.2 Soundness of the Chain Analysis for the Finite Case

In this section we prove that the independence analysis based on chains in possibly infinite case can be equivalently performed on a finite subset of  $k$ -chains. The main statement we prove is that a query and an update are conflicting if and only if they are conflicting on a finite set of  $k$ -chains. As a consequence of this, two expressions are independent if they do not conflict on a specific set of  $k$ -chains.

The proof is constituted of two main parts. We first prove that each chain inferred from an expression can be *folded* to a  $k$ -chain inferred for the same expression, yet preserving all informations concerning the behavior of the query. Then we show that this folding can be done also for a pair of conflicting expressions thus preserving all interactions between a query and update also in the finite analysis.

We develop our proof with the aid of *linear-path queries*, a subclass of queries defined in Chapter 2 that simplify the formal development yet allowing us to represent all query and update chains.

### 9.2.1 Folding Lemma for a Single Expression

In this section we prove Lemma 4.4.4 that states the existence of a folding relation between query chains and  $k$ -query-chains. We will reduce the problem to that of showing a folding relation for a subclass of expressions considered in this paper, called linear-path queries (*lp queries*).

Linear-path queries are defined in order to represent all navigational paths of a query that are superposed at the syntactic level, and are defined as follows.

**Definition 9.2.1** (Linear-path queries). *A linear-path query, denoted by  $q$ , is an expression matching the following grammar*

$$\begin{aligned} q &::= \langle a \rangle q \langle /a \rangle \mid \bar{q} \\ \bar{q} &::= \text{for } x \text{ in } x/\text{axis}::\phi \text{ return } \bar{q} \mid x/\text{axis}::\phi \mid \text{"txt"} \end{aligned}$$

Linear-path queries are the equivalent of navigational XPath expression without branching, plus both element and text node construction. These expressions are called linear because the grammar allows only for linear recursion. In fact, an iteration possibly nests only another iteration (on its right branch), and an element construction possibly nests either another element construction or an iteration. Also, element constructions are allowed only at the outermost level of the query.

A query  $q$  is always reformulated in a set of  $lp$  queries  $\tau_q$ . To illustrate the process, consider the query  $q$  defined as

```
for  $x_1 := x//author, x//editor$ 
return
 $x_1/self :: author$ 
```

is translated into

$$\tau_q = \left\{ \begin{array}{l} \text{for } x_1 \text{ in } x//author \\ \text{return} \\ x_1/self :: author \end{array}, \begin{array}{l} \text{for } x_1 \text{ in } x//editor \\ \text{return} \\ x_1/self :: author \end{array}, x//author, x//editor \right\}$$

It is easy to see that this reformulation provides an over-approximation of all navigational paths embodied in  $q$ , and hence all chains inferred for the original query are contained in the set of chains inferred for the reformulated ones. At the same time, we have that  $k_q = 2$  and for all  $q \in \tau_q$  it holds that  $k_q \leq k_q$ , hence all chains inferred for a reformulated expression can be folded to  $k_q$ -chains.

These two properties together makes sufficient to show a folding on the set of  $lp$ -queries, to scale on the whole language considered in this work. By showing a  $k_q$ -folding for a superset of  $q$  chains, we immediately conclude the  $k_q$ -folding for all  $q$  chains. This simplification is without loss of generality. Queries  $lp$  are a sort of normalized form for the problem of chain inference that simplify the formal development. In previous versions of this work, we developed proofs directly dealing with the general case. These resulted to be prolix, because they required involved properties, without clarifying the crux of the question. Folding for updates expressions will be given later, by means of linear-path queries.

### Reformulating a query in a set of $lp$ -queries

The set of linear-paths induced by a query is defined in Table 9.1. We describe the reformulation process. Given a query  $q$ , the set of  $lp$  queries induced by  $q$  is obtained by splitting paths in concatenations (“,”) and conditional expressions, by approximating let-bindings, by linearizing all iterations of the query and, finally, by rewriting element constructions.

Let-bindings are approximated by means of iteration. By doing this in principle we obtain a correct multiplicity values for the expression, however there is a gap to fill in terms of inferred chains. In fact, iterations discard chains bound to variables if they do not satisfy navigational specifications, while rule (LET) retains all chains bound to a variable  $x$  no regardless. For this reason, we put the iteration in union with all expression bound with a new let-variable. These are return chains instead of used ones, but the nature of the chains is irrelevant for the end of showing a folding relation.



To illustrate path splitting and let-binding approximation, consider the following example.

$$\text{LPQ} \left( \begin{array}{l} \text{let } x_1 := x//\text{author}, x//\text{editor} \\ \text{return} \\ x_1/\text{self} :: \text{author} \end{array} \right) = \left\{ \begin{array}{l} \text{for } x_1 \text{ in } x//\text{author} \\ \text{return} \\ x_1/\text{self} :: \text{author} \end{array}, \begin{array}{l} \text{for } x_1 \text{ in } x//\text{editor} \\ \text{return} \\ x_1/\text{self} :: \text{author} \end{array}, x//\text{author}, x//\text{editor} \right\}$$

Here the chains pointing to *editor* elements are discarded by the iteration. This is correct from the semantics point of view but it is not compliant with the (LET) rule that also produces chains pointing to editor as used chains.

The functions  $R_{\text{for}}^\bullet$  and  $R_e^\bullet$  are introduced for the linearization of the iterations and for lifting element constructions.

We describe now the linearization of iterations. Given a **for** construct, we denote by  $R_{\text{for}}$  the function that linearizes expressions of the form **for**  $x$  **in**  $q_1$  **return**  $q_2$ .  $R_{\text{for}}$  is defined as the identity on any input query, except for the ones matching the following case.

$$R_{\text{for}} \left( \begin{array}{l} \text{for } x_1 \text{ in} \\ (\text{for } x_2 \text{ in } q_2 \text{ return } q_1) \\ \text{return } q_0 \end{array} \right) \stackrel{\text{def}}{=} \begin{array}{l} \text{for } x_2 \text{ in } R_{\text{for}}(q_2) \\ \text{for } x_1 \text{ in } R_{\text{for}}(q_1) \\ \text{return } R_{\text{for}}(q_0) \end{array}$$

Notice that the rule applies locally on a query  $q$  and it has to be iterated until fixpoint for linearizing all subexpressions, until  $R_{\text{for}}(q) = q$ . We denote this by  $R_{\text{for}}^\bullet(q)$ . Also, we denote by  $R_{\text{for}}^\bullet(Q)$  the rewriting of a set of queries  $Q$ .

The fixpoint rewriting works as in the following example.

$$\begin{array}{ccc} \begin{array}{l} \text{for } x_1 \text{ in} \\ (\text{for } x_2 \text{ in} \\ \text{for } x_3 \text{ in } x//\text{book} \\ \text{return } x_3//\text{author} \\ \text{return } x_2//\text{name}) \\ \text{return } x_1//\text{first} \end{array} & \xrightarrow{R_{\text{for}}^\bullet} & \begin{array}{l} \text{for } x_2 \text{ in} \\ (\text{for } x_3 \text{ in } x//\text{book} \\ \text{return } x_3//\text{author}) \\ \text{return} \\ \text{for } x_1 \text{ in } x_2//\text{name} \\ \text{return } x_1//\text{first} \end{array} \\ & & \\ & \xrightarrow{R_{\text{for}}^\bullet} & \begin{array}{l} \text{for } x_3 \text{ in } x//\text{book} \\ \text{for } x_2 \text{ in } x_3//\text{author} \\ \text{for } x_1 \text{ in } x_2//\text{name} \\ \text{return } x_1//\text{first} \end{array} \end{array}$$

$\text{LPQ}(() ) \stackrel{\text{def}}{=} \emptyset$	$\text{LPQ}(\text{"txt"}) \stackrel{\text{def}}{=} \{ \text{"txt"} \}$
$\text{LPQ}(q_1, q_2) \stackrel{\text{def}}{=} \text{LPQ}(q_1) \cup \text{LPQ}(q_2)$	
$\text{LPQ}(x/\text{step}) \stackrel{\text{def}}{=} \{x/\text{step}\}$	
$\text{LPQ}(q) \stackrel{\text{def}}{=} R_e^\bullet(\text{LPQ}(q))$	
$\text{LPQ}(\text{for } x \text{ in } q_1 \text{ return } q_2) \stackrel{\text{def}}{=} \bigcup_{q_i \in \text{LPQ}(q_1)} R_{\text{for}}^\bullet(\text{for } x \text{ in } q_1 \text{ return } q_2)$	
$\text{LPQ}(\text{let } x := q_1 \text{ return } q_2) \stackrel{\text{def}}{=} \text{LPQ}(\text{for } x \text{ in } q_1 \text{ return } q_2) \cup \text{LPQ}(q_1)$	
$\text{LPQ}(\text{if } (q_0) \text{ then } q_1 \text{ else } q_2) \stackrel{\text{def}}{=} \bigcup_i \text{LPQ}(q_i)$	
$\text{LPQ}(\langle a \rangle q \langle /a \rangle) \stackrel{\text{def}}{=} \bigcup_{q \in \text{LPQ}(q)} \langle a \rangle q \langle /a \rangle$	

Table 9.1: Linear-path query extraction

For *lp* queries, element construction is allowed only at the outermost level of the query.<sup>1</sup> Expressions that do not satisfy this property are rewritten in proper ones by lifting element construction. For example, the query

`for x in y//b return <new>x</new>`

is rewritten in

`<new>for x in y//b return x</new>`

First notice that the semantics of the two expressions is different: the former returns a forest of  $n$  trees, where  $n$  is the number of  $b$  elements in a tree, while the latter returns a single tree. Also, that there may be no type definition in the schema employing  $b$  as a label. In this case the chain inference process outputs the empty set for the former expression, while it outputs the element chain *new* for the latter. However, while the semantics of the expressions may differ, the rewriting is still sound for the purpose of chain inference. In this case, the set of used, return and element chains inferred for the original and the rewritten expression coincide. In general, it is easy to show that the set of *element* chains inferred for a *lp* query set a sound approximation of the element chains of the original query.

We denote by  $R_e(q)$  the function that rewrites a query by lifting element constructions.  $R_e(q)$  is defined as the identity for all input query, except for the ones matching

---

<sup>1</sup>It is worth recalling that we assume also that query do not construct elements in the left branch of an iteration or let binding.

one of the following two cases.

$$\begin{aligned}
 & R_e(\text{for } x \text{ in } q_1 \text{ return } \langle a \rangle q_2 \langle /a \rangle) \\
 & \quad \quad \quad \stackrel{def}{=} \\
 & \langle a \rangle \text{for } x \text{ in } q_1 \text{ return } R_e(q_2) \langle /a \rangle \\
 \\
 & R_e(\langle a \rangle q \langle /a \rangle) \stackrel{def}{=} \langle a \rangle R_e(q) \langle /a \rangle
 \end{aligned}$$

The function is intentionally defined to work only on that specific form of iteration. It is meant to be applied over queries that are (almost) *lp* queries, up to element construction. Notice that the rules have to be applied until fixpoint is reached, i.e.  $R_e(q) = q$ . We denote this by  $R_e^\bullet(q)$ . We will use the same notation also when speaking about set of queries. The fixpoint is unique and always reached by the translation.

As an example consider the previous query

$$\text{for } x \text{ in } y//b \text{ return } \langle new \rangle q \langle /new \rangle$$

where  $q = \text{for } x_2 \text{ in } y//a \text{ return } \langle sub \rangle$ . We have that

$$\begin{aligned}
 & R_e^\bullet(\text{for } x \text{ in } y//b \text{ return } \langle new \rangle q \langle /new \rangle) \\
 = & R_e^\bullet(\langle new \rangle \text{for } x \text{ in } y//b \text{ return } R_e(q) \langle /new \rangle) \\
 = & R_e^\bullet(\langle new \rangle \text{for } x \text{ in } y//b \text{ return} \\
 & \quad \quad \quad \langle sub \rangle \\
 & \quad \quad \quad \text{for } x_2 \text{ in } y//a \text{ return } () \\
 & \quad \quad \quad \langle /sub \rangle \\
 & \quad \quad \quad \langle /new \rangle) \\
 = & \langle new \rangle \\
 & \quad \langle sub \rangle \\
 & \quad \text{for } x \text{ in } y//b \text{ return} \\
 & \quad \quad \text{for } x_2 \text{ in } y//a \text{ return } () \\
 & \quad \quad \langle /sub \rangle \\
 & \quad \langle /new \rangle
 \end{aligned}$$

This translation is reminiscent of path extraction, however there are two difference. The first is that in the definition of the set of *lp* queries associated to an expression, we do not take into account variable binding. We avoid this process because it is not needed in order to have a sound *k*-folding. The second difference, and most important, is that this process is described formally, but actually never implemented. It is just used to show that the *k*-analysis is sound in principle.

Correctness wrt chains of the  $LPQ()$  transformation of queries is stated as follows.

**Proposition 9.2.2.** *Let  $d$  be a DTD,  $q$  a query, and  $\Gamma$  a static environment. Let us denote by  $\tau_q$  the set of used, return and element chains inferred for  $q$  wrt  $\Gamma$  and by  $\tau_{LPQ(q)}$  the union of used, return, and element chains inferred for each expression in  $LPQ(q)$ . Then,  $\tau_q \subseteq \tau_{LPQ(q)}$ .*

*Proof.* By induction on the structure of  $q$ . □

By over approximating the set of chains inferred for an expression, we have that also the result of a query is always over approximated by its associated set of linear-path queries. However, what is important for us is not only the resulting locations of a query, but also the whole set of locations accessed. This can be showed as a corollary of Proposition 9.2.2.

### Computing the Multiplicity Value for $lp$ -queries

In the general case (Section 2.3), the definition of the  $k$  value is (necessarily) structured, because several navigational paths are superposed at the syntactic level. Differently, linear-path allow us to directly compute the multiplicity value for a given query expression defined in Chapter 4, without inductive definition. This will be useful in the proof in order to show that the finite on  $k$ -chains is equivalent to the possibly infinite analysis.

**Proposition 9.2.3.** *Let  $q$  be an  $lp$  query without element construction then*

$$k_q = \max_{a \in \Sigma} \left\{ \sum_{\text{step} \in q} \mathcal{F}(a, \text{step}) + \mathcal{R}(\text{step}) \right\}$$

*Furthermore, in the case that  $q$  is nested inside an element construction, we have*

$$k_{\langle a \rangle q \langle /a \rangle} = \begin{cases} 1 + k_q & \text{if } a \text{ is the most frequent label in } q \\ k_q & \text{otherwise} \end{cases}$$

*Where  $a$  is the most frequent label in  $q$  if, for all  $b \in \Sigma$ , if we have that*

$$\sum_{\text{step} \in q} \mathcal{F}(b, \text{step}) \leq \sum_{\text{step} \in q} \mathcal{F}(a, \text{step})$$

*Proof.* Straightforward by induction on  $q$ . □

The following lemma shows that all  $lp$ -queries inferred for an expression have a multiplicity value bounded by that of the original query.

**Proposition 9.2.4.** *Let  $d$  be a DTD,  $q$  a query, and  $\Gamma$  a static environment. Let us denote by  $\tau_q$  and  $\tau_{LPQ(q)}$  the union of the set of used, return, and element chains inferred for  $q$  wrt  $\Gamma$ , and for each expression in  $LPQ(q)$ , respectively. Then, for all query  $q \in LPQ(q)$  we have that  $k_q \leq k_q$ .*

*Proof.* By induction on the structure of  $q$ . □

The existence of a folding relation between query chains and  $k$  query chains is stated by the following theorem.

**Lemma 9.2.5** (Folding of linear-path queries). *Let  $d$  be a DTD,  $q$  a query with at most only one free variable  $x$ , and  $\Gamma = \{ x \mapsto s_d \}$  a static environment. For all  $q \in \text{LPQ}(q)$ , provided that  $\Gamma \vdash_{C_d} q : (r, v, e)$  and  $\tau = r \cup v \cup e$ , for each chain  $c \in \tau$  there exists a chain  $c' \in \tau$  such that  $c \hookrightarrow_d^* c'$  and  $c'$  is a  $k_q$ -chain.*

Thanks to the above lemma, we can finally conclude that the folding relation exists for chains inferred for all queries.

**LEMMA 4.4.4 (FOLDING)** *Let  $d$  be a DTD,  $q$  a query with at most only one free variable  $x$ , and  $\Gamma = \{ x \mapsto s_d \}$  a static environment. Assume that  $\Gamma \vdash_{C_d} q : (r, v, e)$  and  $\tau = r \cup v \cup e$ . For each chain  $c \in \tau$  there exists a chain  $c' \in \tau$  such that  $c \hookrightarrow_d^* c'$  and  $c'$  is a  $k_q$ -chain.*

*Proof.* Let  $c \in \tau_q$  be a chain. By Lemma 9.2.2 we have that  $c \in \tau_{\text{LPQ}(q)}$ . We conclude since by Lemma 9.2.5  $c$  can be folded to a  $k_q$ -chains.  $\square$

In the the remainder of this section we prove Lemma 9.2.5.

### Folding Lemma for $lp$ -queries

The height of an  $lp$ -query  $q$ , denoted by  $h(q)$ , is defined as

$$\begin{aligned} h(()) &= 0 \\ h(\text{"txt"}) &= 0 \\ h(x/\text{axis} :: \phi) &= 0 \\ h(\text{for } x \text{ in } q_1 \text{ return } q_2) &= 1 + \max\{h(q_1), h(q_2)\} \\ h(\langle a \rangle q \langle /a \rangle) &= 1 + h(q) \end{aligned}$$

The following lemma states that any chain can be folded to a 1-chain.

**Lemma 9.2.6.** *Let  $d$  be a DTD. For all chain  $c \in C_d$  there exists a 1-chain  $c' \in C_d$  such that  $c \hookrightarrow_d^* c'$ .*

*Proof.* Given  $c \in C_d$ , if  $c$  is a 1-chain then we are done. Otherwise  $c$  is of the form  $c = c_1.w.c_2$ , where  $w$  is a cycle  $w = a.a_1 \dots a_n.a$ . By definition of the folding relation we have that

$$c_1.w.c_2 \hookrightarrow_d c_1.a.c_2$$

Now if  $c_1.a.c_2$  is a 1-chain then we are done, otherwise we repeat the folding process until we obtain a chain  $c'$  with no cycles, such that  $c \hookrightarrow_d^* c'$ . This process terminates because the number of symbols in  $c$  is finite.  $\square$

**Lemma 9.2.5** (Folding of linear-path queries) *Let  $\mathbf{d}$  be a DTD,  $\mathbf{q}$  a query with at most only one free variable  $\mathbf{x}$ , and  $\Gamma = \{ \mathbf{x} \mapsto \mathbf{s}_d \}$  a static environment. For all  $q \in \text{LPQ}(\mathbf{q})$ , provided that  $\Gamma \vdash_{C_d} q : (r, v, e)$  and  $\tau = r \cup v \cup e$ , for each chain  $c \in \tau$  there exists a chain  $c' \in \tau$  such that  $c \mapsto_d^* c'$  and  $c'$  is a  $k_q$ -chain.*

PROOF. By induction on  $h(q)$ .

*Base* ( $h = 0$ ) In this case  $q$  has one of the following three forms

$$() \quad | \quad \text{"txt"} \quad | \quad \mathbf{x}/\text{axis} :: \phi$$

We conclude the proof by a case analysis.

$[q = ()]$  Straightforward since  $\tau_q = \emptyset$ .

$[q = \text{"txt"}]$  Recall that because the query “txt” is independent of any update, in Chapter 2 we assumed that “txt” is nested into a for expression. Furthermore, because we do not navigate on constructed sequences, we have that “txt” should be in the right branch of the for expression. Then, in order to do chain inference for “txt” we must have a navigational step in the left branch of  $q$  that makes  $k_q > 0$ . In the case where  $k_q = 0$  it means that no operation is performed by the left branch of the query and thus chains for “txt” are not inferred by the type system rules. We conclude since  $\tau_q = \{\text{String}\}$  and String is a 1-chain.

$[q = \mathbf{x}/\text{axis} :: \phi]$  Because  $\mathbf{x} \mapsto \mathbf{s}_d$  we have that  $\tau_q \neq \emptyset$  only if **axis** is one of **self**, **child**, **descendant**, or **descendant-or-self**. We assume  $\phi = \mathbf{a} \in \Sigma$  and we proceed with a case analysis on **axis**. The proof for the case  $\phi \in \{\text{node}(), \text{text}()\}$  is analogous, since **node**() can be simulated by ranging over all labels in  $\Sigma$ , and **text**() by considering only type String.

$[\mathbf{axis} = \text{child}]$  In this case  $k_q = 1$  and we conclude since either  $\tau_q = \emptyset$  or  $\tau_q = \{ \mathbf{s}_d.\mathbf{a} \}$ , where  $\mathbf{s}_d \neq \mathbf{a}$  by hypothesis of non-recursive schema root (see Chapter 2).

$[\mathbf{axis} = \text{self}]$  In this case  $k_q = 1$  and we conclude since either  $\tau_q = \emptyset$  or  $\tau_q = \{ \mathbf{s}_d \}$ .

$[\mathbf{axis} = \text{descendant}]$  In this case  $k_q = 1$  since  $\mathcal{R}(\text{descendant} :: \mathbf{a}) = 1$ . We conclude since by Lemma 9.2.6 any chain can be folded to a 1-chain.

$[\mathbf{axis} = \text{descendant-or-self}]$  Analogous to the **self** and **descendant** cases.

*Induction* ( $h > 0$ )

In this In this case  $q$  has one of two forms.

$$\langle \mathbf{a} \rangle q' \langle / \mathbf{a} \rangle \quad | \quad \text{for } \mathbf{x} \text{ in } q_1 \text{ return } q_2$$

$[q = \langle \mathbf{a} \rangle q' \langle / \mathbf{a} \rangle]$  We prove this case by making a case analysis wrt the set of used, return and element chains inferred for  $q$ .

**used chains** By rule (ELT),  $v_q = \bar{r}_{q'} \cup v_{q'}$ . All chains inferred for the inner query  $q'$  are obtained by a derivation tree of height  $h(q') < h(q)$ . Then by inductive hypothesis on  $q'$  we conclude the thesis for all chains in  $r_{q'} \cup v_{q'}$ , since they are  $k_{q'}$ -chains and provided that  $k_{q'} \leq k_q$  they are  $k_q$ -chains.

**return chains** By rule (ELT),  $r_q = \emptyset$  hence no return chain is inferred.

**element chains** We have two main cases to consider.

- If  $r_{q'} \cup v_{q'} = \emptyset$  then  $e_q = \{a\}$ . We conclude since  $k_q = 1$ .
- If  $r_{q'} \cup v_{q'} \neq \emptyset$  then  $e_q = e_q^e \cup e_q^r$  where

$$e_q^e = \{ a.c \mid c \in e_{q'} \} \quad e_q^r = \{ a.a.c' \mid c.a \in r_{q'}, c.a.c' \in \bar{r}_{q'} \}$$

The two sets are built from element and return chains of  $q'$ , respectively. Let us consider them separately.

- If  $a.c$  belongs to  $e_q^e$  then  $c \in e_{q'}$ . By inductive hypothesis  $c$  is a  $(k_{q'})$ -chain. If  $k_q = 1 + k_{q'}$  then  $a$  is the most frequent label in  $c$  with frequency upper bounded by  $k_{q'}$ . Therefore  $a.c$  is a  $(1 + k_{q'})$ -chain and hence a  $k_q$ -chain. In the case that  $k_q = k_{q'}$  then  $a$  is not the most frequent label in  $c$ . Thus the frequency of  $a$  in  $c$  is strictly upper bounded by  $k_{q'}$ . Therefore  $a.c$  is still a  $k_{q'}$ -chain and thus a  $k_q$ -chain.
- If  $a.a.c'$  belongs to  $e_q^r$  then  $c.a \in r_{q'}$  and  $c.a.c' \in \bar{r}_{q'}$ . By inductive hypothesis  $c.a.c'$  can be folded to a  $(1 + k_{q'})$ -chain  $c_1.a.c'_1$ . Thus  $a.a.c'$  can be folded to  $a.a.c'_1$ . Provided that  $a.c'_1$  is a  $(1 + k_{q'})$ -chain, we have that  $a.a.c'$  is again a  $(1 + k_{q'})$ -chain or the prefix  $a.c'$  can be folded to a  $(k_{q'})$ -chain, thus obtaining a  $(1 + k_{q'})$ -chain.

$[q = \text{for } x \text{ in } q_1 \text{ return } q_2]$  Here  $q_2$  possibly nests an arbitrary number  $n$  of iterations. In particular, we assume that the rightmost subexpression of  $q$  is a navigational step. It may also be that the rightmost subexpression is a text-element construction, but in that case the proof is a particular instance of the one given below where one does not make any use of the static variable environment. Any derivation tree of chain inference for  $q$  must end by applying the (FOR) rule, i.e., it looks like this (for simplicity we do not detail all rule operations)

$$\begin{array}{c}
\Gamma \vdash_C \mathbf{step}_1 : \tau_1 \\
\hline
\mathbf{x}_1 \mapsto \tau_1, \Gamma \vdash_C \mathbf{step}_2 : \tau_2 \\
\hline
\vdots \\
\hline
\{\mathbf{x}_i \mapsto \tau_i\}_{i \leq n-1}, \Gamma \vdash_C \mathbf{step}_n : \tau_n \quad \{\mathbf{x}_i \mapsto \tau_i\}_{i \leq n}, \Gamma \vdash_C \mathbf{step}_{n+1} : \tau_{n+1} \\
\hline
\Gamma \vdash_C \text{for } \mathbf{x}_1 \text{ in } \mathbf{step}_1 \dots \mathbf{x}_n \text{ in } \mathbf{step}_n \text{ return } \mathbf{step}_{n+1} : \tau_q
\end{array}$$

where  $\mathbf{step}_i$  is of the form  $\mathbf{x}_i/\mathbf{axis}_i :: \phi_i$ , with  $\mathbf{x}_i$  either a free-variable of  $q$  bound to the root-type in  $\Gamma$ , or a bound-variable of  $q$  previously defined in  $\{\mathbf{x}_1, \dots, \mathbf{x}_{i-1}\}$ . Also, the resulting set of chains  $\tau_q \subseteq \tau_n \cup \tau_{n+1}$ , as specified by rule (FOR).

We denote by  $h$  the height of the derivation tree for  $q$ . For each variable  $\mathbf{x}_i$ , chains in  $\tau_i$  are obtained by a derivation of height  $h_i < h$ . By inductive hypothesis, the chains bound to  $\mathbf{x}_i$  can be folded to  $k_i$ -chains, where  $k_i$  is

$$k_i = \max_{a \in \Sigma} \left\{ \sum_{\mathbf{step}=\mathbf{step}_1}^{\mathbf{step}_i} \mathcal{F}(a, \mathbf{step}) + \mathcal{R}(\mathbf{step}) \right\}$$

Since  $\tau_q \subseteq \tau_n \cup \tau_{n+1}$ , and all chains in  $\tau_n$  are  $k_q$ -chains, what is left to show is that all chains inferred for  $\mathbf{step}_{n+1}$  can be folded to  $k_{n+1}$  chains, provided that the static environment  $\Gamma$  binds variables with chains that can be folded to  $k_n$  chains. This allows us to conclude the proof since  $k_{n+1} = k_q$ .

We proceed with a case analysis on  $\mathbf{step}_i$ , assuming as before that  $\phi_i = \mathbf{a} \in \Sigma$ . We focus on the cases where  $\tau_q \neq \emptyset$ , i.e., the inference yields a non-empty result for some chains bound to  $\mathbf{x}$  in  $\Gamma$ . When  $\tau_q = \emptyset$  the thesis follows straightforwardly.

[ $\mathbf{step}_{n+1} = \mathbf{x}/\mathbf{child} :: \mathbf{a}$ ] In this case

$$\begin{aligned}
\tau_{n+1} &= \bigcup_{c \in \Gamma(\mathbf{x})} \mathbf{T}_{C_d}(\mathbf{A}_{C_d}(c, \mathbf{child}), \mathbf{a}) \\
&= \{ c.\mathbf{a} \mid c \in \Gamma(\mathbf{x}) \}
\end{aligned}$$

By inductive hypothesis, for each  $c \in \Gamma(\mathbf{x})$  there exists a  $k_n$ -chain  $c' \in \Gamma(\mathbf{x})$  such that  $c \hookrightarrow_d^* c'$ . Provided this, it follows  $c.\mathbf{a} \hookrightarrow_d^* c'.\mathbf{a}$ . We need to show that  $c'.\mathbf{a}$  is a  $k_{n+1}$ -chain. We do this by distinguishing two cases, depending on whether  $k_n < k_{n+1}$  or  $k_n = k_{n+1}$ .



If  $k_n < k_{n+1}$ ,  $c'$  being a  $k_n$ -chain, it follows that  $c'.a$  is a  $k_{n+1}$ -chain.

If  $k_n = k_{n+1}$ , then  $c'.a$  is a  $k_n$ -chain because  $a$  appears strictly less than  $k_n$  times in  $c'$ . In other words,  $a$  is not the most frequent label in  $\text{step}_1, \dots, \text{step}_n$ . This can be shown as follows by considering that

$$\begin{aligned} \sum_{i=1}^n \mathcal{F}(a, \text{step}_i) + \mathcal{R}(\text{step}_i) \\ < \\ \sum_{i=1}^{n+1} \mathcal{F}(a, \text{step}_i) + \mathcal{R}(\text{step}_i) \end{aligned}$$

and provided that

$$\begin{aligned} \max_{a \in \Sigma} \left\{ \sum_{i=1}^n \mathcal{F}(a, \text{step}_i) + \mathcal{R}(\text{step}_i) \right\} \\ = \\ \max_{a \in \Sigma} \left\{ \sum_{i=1}^{n+1} \mathcal{F}(a, \text{step}_i) + \mathcal{R}(\text{step}_i) \right\} \end{aligned}$$

we conclude the case because

$$\sum_{i=1}^n \mathcal{F}(a, \text{step}_i) + \mathcal{R}(\text{step}_i) < k_n$$

In the following we denote by  $c_a$  the fact that  $a$  is the last symbol of the chain  $c$ , that is  $c = c'.a$ .

[ $\text{step}_{n+1} = \text{descendant} :: a$ ] In this case

$$\begin{aligned} \tau_{n+1} &= \bigcup_{c \in \Gamma(x)} \mathbf{T}_{C_d}(\mathbf{A}_{C_d}(c, \text{descendant}), a) \\ &= \{ c.c_a \in C_d \mid c \in \Gamma(x) \} \end{aligned}$$

Given  $c.c_a \in \tau_{n+1}$ , by inductive hypothesis there exists a  $k_n$ -chain  $c' \in \Gamma(x)$  such that  $c \hookrightarrow_d^* c'$ . By Lemma 9.2.6 the suffix  $c_a$  can be folded to a 1-chain.<sup>2</sup> At this point we have that  $c.c_a \hookrightarrow_d^* c'.c'_a$ . Provided that  $c'$  is a  $k_n$ -chain and  $c'_a$  is a 1-chain,  $c'.c'_a$  is  $(1+k_n)$ -chain. We conclude since  $\mathcal{R}(\text{descendant}) = 1$  and thus  $k_{n+1} = k_n + 1$ .

[ $\text{step}_{n+1} = \text{descendant-or-self} :: a$ ] Analogous to the case **descendant**.

[ $\text{step}_{n+1} = \text{parent} :: a$ ] In this case

$$\begin{aligned} \tau_{n+1} &= \bigcup_{c \in \Gamma(x)} \mathbf{T}_{C_d}(\mathbf{A}_{C_d}(c, \text{parent}), a) \\ &= \{ c_a \mid c = c_a.b \in \Gamma(x) \} \end{aligned}$$

---

<sup>2</sup>Notice however that  $c_2.a \hookrightarrow_d^* c'_2$  does not imply  $c_2 \hookrightarrow_d^* c'_2$ . For instance,  $b.a.b \hookrightarrow_d b$ , yet  $b.a \not\hookrightarrow_d b$ .

where  $\mathbf{b} \in \Sigma \cup \{\text{String}\}$ .

We want to show a folding for  $\mathbf{c}_a.\mathbf{b}$  and for  $\mathbf{c}_a$ . By inductive hypothesis, for each  $\mathbf{c}_b = \mathbf{c}_a.\mathbf{b} \in \Gamma(\mathbf{x})$  there exists a  $k_n$ -chain  $\mathbf{c}'_b$  in  $\Gamma(\mathbf{x})$  such that  $\mathbf{c}_b \hookrightarrow_d^* \mathbf{c}'_b$ . We have two cases to consider, depending on whether  $\mathbf{a}$  and  $\mathbf{b}$  are mutually recursive or not.

If  $\mathbf{a}$  and  $\mathbf{b}$  are not mutually recursive ( $\mathbf{a} \Rightarrow_d \mathbf{b}$  but  $\mathbf{b} \not\Rightarrow_d^* \mathbf{a}$ ) then none of the chains belonging to  $\mathbf{C}_d$  has a cycle involving the two types. So by folding  $\mathbf{c}_b$  to  $\mathbf{c}'_b$  the suffix  $\mathbf{a}.\mathbf{b}$  is preserved. Let  $\mathbf{c}'_b = \mathbf{c}'_a.\mathbf{b}$ , we conclude since  $\mathbf{c}_a \hookrightarrow_d^* \mathbf{c}'_a$ .

If  $\mathbf{a}$  and  $\mathbf{b}$  are mutually recursive ( $\mathbf{a} \Rightarrow_d^* \mathbf{b}$  and  $\mathbf{b} \Rightarrow_d^* \mathbf{a}$ ) then there are two cases to consider. If the prefix  $\mathbf{a}.\mathbf{b}$  is preserved by some foldings then we conclude. We discuss the case where the prefix  $\mathbf{a}.\mathbf{b}$  is never preserved by any folding to  $\mathbf{c}'_b$ . This happens if the schema enforces that  $\mathbf{b}$  is always defined in terms of  $\mathbf{a}$  (i.e.,  $\mathbf{a}$  is always a descendant of  $\mathbf{b}$ ) not the other way.

Notice that in this case  $\mathbf{b} \neq \text{String}$  and necessarily  $\mathbf{b} \in \Sigma$ . This implies that the first occurrence of  $\mathbf{b}$  precedes any occurrence of  $\mathbf{b}$  in the input chain  $\mathbf{c}_b$ , hence  $\mathbf{c}_b$  is of the form

$$\mathbf{c}_b = \mathbf{c}_1. \underbrace{\mathbf{b}.\mathbf{b}_1 \dots \mathbf{b}_n.\mathbf{a}.\mathbf{b}}_w$$

and the cycle  $w$  is folded in order to respect the multiplicity value  $k_n$ , thus obtaining  $\mathbf{c}'_b = \mathbf{c}'_1.\mathbf{b}$ .<sup>3</sup>

Then we do the following construction. Take  $w' = \mathbf{b}.w_a.\mathbf{b}$  as the folding of  $w$ , where  $\text{last}(w_a) = \mathbf{a}$  and  $w_a.\mathbf{b}$  is a 1-chain. Concatenate  $\mathbf{c}'_1$  with  $w'$ , thus obtaining a  $(1+k_n)$ -chain

$$\mathbf{c}' = \mathbf{c}'_1.w'$$

Now it remains to show that  $k_{n+1} = 1 + k_n$ . We show it by contradiction. Assume  $k_{n+1} = k_n$ . Then, as shown for the case  $\text{axis} = \text{child}$ ,  $\mathbf{a}$  is not the most frequent label in  $\text{step}_1, \dots, \text{step}_n$ . Since  $\mathbf{a}$  belongs to  $w_a$ , then also this symbols are not the most frequent in  $\text{step}_1, \dots, \text{step}_n$ . Thus there exists a chain  $\mathbf{c}'_a$  preserving the suffix  $\mathbf{b}.\mathbf{a}$ . A contradiction.

[ $\text{step}_{n+1} = \text{ancestor} :: \mathbf{a}$ ] In this case

$$\begin{aligned} \tau_{n+1} &= \bigcup_{\mathbf{c} \in \Gamma(\mathbf{x})} \mathbf{T}_{\mathbf{C}_d}(\mathbf{A}_{\mathbf{C}_d}(\mathbf{c}, \text{ancestor}), \mathbf{a}) \\ &= \{ \mathbf{c}_a \mid \mathbf{c} = \mathbf{c}_a.\mathbf{c}' \in \Gamma(\mathbf{x}) \} \end{aligned}$$

with  $\mathbf{c}' \neq \epsilon$ . We want to show a folding for  $\mathbf{c}_a$ .

We have two cases to consider, depending on whether  $\mathbf{a}$  is recursive or not.

---

<sup>3</sup>This happens for instance for schemas such as  $\mathbf{d}_1 = \{\mathbf{s} \leftarrow \mathbf{a}; \mathbf{a} \leftarrow \mathbf{b}?\}$  or  $\mathbf{d}_2 = \{\mathbf{s} \leftarrow (\mathbf{a}, \mathbf{b}); \mathbf{a} \leftarrow \mathbf{b}?\}; \mathbf{b} \leftarrow \mathbf{a}^*\}$  or  $\mathbf{d}_3 = \{\mathbf{s} \leftarrow \mathbf{a}; \mathbf{a} \leftarrow \mathbf{b}_1?; \dots \mathbf{b}_{n-1} \leftarrow \mathbf{b}_n; \mathbf{b}_n \leftarrow \mathbf{a}\}$  and a query such as  $q = //a/\text{parent} :: \mathbf{a}$

If  $\mathbf{a}$  is not recursive, then no chain in  $\mathbf{C}_d$  has a cycle involving  $\mathbf{a}$ . The type  $\mathbf{a}$  is in a path between two distinct strongly-connected components of the dependency graph induced by the schema. Therefore  $\mathbf{c} = \mathbf{c}_a.\mathbf{c}_2$  is of the form

$$\mathbf{c} = \mathbf{c}_1.\mathbf{a}.\mathbf{c}_2$$

where  $\mathbf{c}_1$  and  $\mathbf{c}_2$  do not share any label. By inductive hypothesis, there exists a  $k_n$ -chain  $\mathbf{c}' \in \Gamma(\mathbf{x})$  such that  $\mathbf{c} \rightarrow_d^* \mathbf{c}'$ . Since  $\mathbf{c}_1$  and  $\mathbf{c}_2$  do not share any label, it is easy to see that

$$\mathbf{c}' = \mathbf{c}'_1.\mathbf{a}.\mathbf{c}'_2$$

where  $\mathbf{c}_i \rightarrow_d^* \mathbf{c}'_i$ , and  $\mathbf{c}'_i$  is a  $k_n$ -chain, ( $i = 1, 2$ ). Then  $\mathbf{c}' = \mathbf{c}'_1.\mathbf{a}$  is a  $k_n$ -chain folding of  $\mathbf{c}_a$ , and a  $k_{n+1}$ -chain as well.

If  $\mathbf{a}$  is recursive, then we distinguish two sub cases. If  $\mathbf{step}_1, \dots, \mathbf{step}_n$  do not employ any recursive axis, then no  $\mathbf{c} \in \Gamma(\mathbf{x})$  need to be folded. Therefore the proof is concluded straightforwardly by inductive hypothesis on the derivation of the expression that generated the chains for  $\mathbf{x}$ .

Otherwise, if some steps in  $\mathbf{step}_1, \dots, \mathbf{step}_n$  employ a recursive axis then, analogously to the case of  $\mathbf{axis} = \mathbf{parent}$ , the label  $\mathbf{a}$  may not be preserved by the folding relation. If  $\mathbf{a}$  is preserved then the conclusion is immediate. We assume this is not the case.

By inductive hypothesis, there exists a  $k_n$ -chain  $\mathbf{c}' \in \Gamma(\mathbf{x})$  such that  $\mathbf{c} \rightarrow_d^* \mathbf{c}'$ , and  $\mathbf{c}'$  is of the form

$$\mathbf{c}' = \mathbf{b}_1.\mathbf{b}_2 \dots \mathbf{b}_m$$

and  $\mathbf{b}_i \neq \mathbf{a}$  for all  $i \in \{1, \dots, m\}$ .

We describe a procedure for constructing the desired  $k_{n+1}$ -chain.

For all  $i = 1..m$ , if

$$\mathbf{b}_i \Rightarrow_d^* \mathbf{a} \quad \text{and} \quad \mathbf{a} \Rightarrow_d^* \mathbf{b}_{i+1}$$

then construct a new chain from  $\mathbf{c}'$ , by injecting a word  $u$  between  $\mathbf{b}_i$ , such as to obtain a chain of the form

$$\mathbf{c}''_i = \mathbf{b}_1 \dots \mathbf{b}_i . u . \mathbf{b}_{i+1} \dots \mathbf{b}_m$$

such that  $u = \mathbf{a}_1 \dots \mathbf{a}_k$  is a 1-chain with  $\mathbf{a} \in u$ , and  $\mathbf{b}_i \Rightarrow_d \mathbf{a}_1$  and  $\mathbf{a}_n \Rightarrow_d \mathbf{b}_{i+1}$ .

Now, to check that  $\mathbf{c}''_i$  we obtained from the construction is a chain matching the semantics of  $\mathbf{step}_n$  it suffices to check that  $\mathbf{c}''_i \in \Gamma(\mathbf{x})$ . If this is not the case, then we chose another  $\mathbf{b}_i$  and we repeat the process. If no other  $\mathbf{b}_i$  exists than the step is unsatisfiable. If  $\mathbf{c}''_i \in \Gamma(\mathbf{x})$  then we define  $\mathbf{c}''_i = \mathbf{c}''_a.\mathbf{c}''$ , where  $\mathbf{c}''_a$  is the prefix of  $\mathbf{c}''_i$  terminating with  $\mathbf{a}$ . According to chain inference for the ancestor axis, we have that  $\mathbf{c}''_a \in \tau_{n+1}$  with  $\mathbf{c}''_a$  is a  $(1+k_n)$ -chain. Because  $\mathcal{R}(\mathbf{ancestor}) = 1$ ,  $k_{n+1} = 1 + k_n$  and we conclude. It remains to show that some  $\mathbf{c}''_i$  constructed as above belong

to  $\Gamma(\mathbf{x})$ . Assume that for all  $\mathbf{c}_i''$  we have that  $\mathbf{c}_a'' \notin \Gamma(\mathbf{x})$  then we conclude that  $\text{step}_1, \dots, \text{step}_n$  do not employ recursive navigations. A contradiction.

[ $\text{step}_{n+1} = \text{following-sibling} :: \mathbf{a}$ ] In this case

$$\begin{aligned} \tau_{n+1} &= \bigcup_{\mathbf{c} \in \Gamma(\mathbf{x})} \mathbf{T}_{\mathbf{C}_d}(\mathbf{A}_{\mathbf{C}_d}(\mathbf{c}, \text{following-sibling}), \mathbf{a}) \\ &= \{ \mathbf{c}_p \cdot \mathbf{a} \mid \mathbf{c} = \mathbf{c}_p \cdot \mathbf{b} \in \Gamma(\mathbf{x}), \mathbf{b} <_{d(\mathbf{c}_p)} \mathbf{a} \} \end{aligned}$$

We want to show a folding for  $\mathbf{c}_p \cdot \mathbf{a}$ . By case  $\text{axis} = \text{parent}$ , we know that there exists a  $(1 + k_n)$ -chain  $\mathbf{c}'_p \cdot \mathbf{a} \in \tau_n$  such that  $\mathbf{c}_p \cdot \mathbf{a} \hookrightarrow_d^* \mathbf{c}'_p \cdot \mathbf{a}$ . Also, in this case the last label of  $\mathbf{c}'_p$  may be  $\mathbf{a}$  as well. Now independently from whether  $\mathbf{a} = \mathbf{a}$  or not  $\mathbf{c}'_p \cdot \mathbf{a}$  is a  $1 + k_n$ -chain, as well as  $\mathbf{c}'_p \cdot \mathbf{a}$ . If  $k_{n+1} = 1 + k_n$ , then we conclude. Otherwise, if  $k_{n+1} = k_n$  then, analogously to the case  $\text{axis} = \text{child}$ , we can show that  $\mathbf{a}$  is not the most frequent label in  $\text{step}_1, \dots, \text{step}_n$ , and thus that  $\mathbf{c}'_p \cdot \mathbf{a}$  is a  $k_n$ -chain.

□

**Remark 9.2.7.** *As a direct corollary of the lemma we can approximate the problem of query unsatisfiability in presence of a schema just looking on query  $k$ -chains. Indeed, if no  $k$ -chain is inferred for the input query expression, then no data match query needs. Moreover, for the subset of queries and schemas that make our chain inference system enjoying completeness, we can exactly determine whether a query is satisfiable or not in presence of a schema.*

### 9.2.2 Folding Lemma with conflict preservation

In this section we prove Lemma 4.4.5, by showing the existence of a folding relation that preserves also conflicts among chains. We show that such a folding exists for a pair of expressions if the multiplicity value we take is the sum of the multiplicity value of each expression.

**Lemma 9.2.8.** *Let  $d$  be a DTD,  $q_1, q_2$  two queries with at most only one free variable  $\mathbf{x}$ , and  $\Gamma = \{ \mathbf{x} \mapsto \mathbf{s}_d \}$  a static environment. For all  $\mathbf{q}'_i \in \text{LPQ}(q_i)$ , provided that  $\Gamma \vdash_{\mathbf{C}_d} \mathbf{q}'_i : (r_i, v_i, e_i)$  and  $\tau_i = r_i \cup v_i \cup e_i$ , for each chain  $\mathbf{c} \in \tau$  there exists a chain  $\mathbf{c}' \in \tau$  such that  $\mathbf{c} \hookrightarrow_d^* \mathbf{c}'$  and  $\mathbf{c}'$  is a  $k_{\mathbf{q}}$ -chain. For each pair of chains  $(\mathbf{c}_1, \mathbf{c}_2) \in \tau_1 \times \tau_2$  such that  $\mathbf{c}_1 \leq \mathbf{c}_2$ , there exists  $(\mathbf{c}'_1, \mathbf{c}'_2) \in \tau_1 \times \tau_2$  such that  $\mathbf{c}'_1 \leq \mathbf{c}'_2$  and  $\mathbf{c}_i \hookrightarrow_d^* \mathbf{c}'_i$  with  $\mathbf{c}'_i$  a  $(k_{\mathbf{q}'_1} + k_{\mathbf{q}'_2})$ -chain ( $i=1, 2$ ).*

*Proof.* We consider first the case when  $\mathbf{c}_1, \mathbf{c}_2$  are either used or return chains. Since  $\mathbf{c}_1 \leq \mathbf{c}_2$ , we have that  $\mathbf{c}_2$  is the longest chain and we make a case analysis on  $q_2$ . We distinguish three main cases depending on whether  $q_2$  contains some recursive axis or not. If  $q_2$  does not contains (forward or backward) recursive axis,  $\tau_{q_2}$  is finite. Moreover,

all chains in  $\tau_{q_2}$  have height upper bounded by  $|q_2|$  and label frequency upper bounded by

$$k_{q_2} = \max_{a \in \Sigma} \left\{ \sum_{\text{step} \in q_2} \mathcal{F}(a, \text{step}) \right\}$$

Since  $c_1 \leq c_2$ , it follows that  $c_1$  is a  $k_{q_2}$ -chain as well. Let  $k = k_{q_1} + k_{q_2}$  we conclude that all pairs of conflicting chains are  $k$ -chains, as desired. Notice that if  $q_2$  also uses recursive backward navigation the statement still hold. As a matter of facts, a backward navigation would just backtrack on the input chain, without affecting the property that chains inferred for  $q_2$  are  $k_{q_2}$  chains.

If  $q_2$  uses at least one recursive *forward* axis, then an upper bound of the label frequency of chains witnessing the conflict is

$$k = k_{q_1} + k_{q_2}$$

The case requiring the largest  $k$  value occur when each prefix of  $q_2$  navigates some descendant of nodes selected by  $q_1$ , in the sense that  $c_2 \in \tau_{q_2}$ , when  $q_2$  has only one free variable  $x$ , and the chains are derived with the judgment

$$\Gamma[x \mapsto c_1] \vdash_{C_d} q_2 : \tau_{q_2}$$

In this case let  $q_1 \circ q_2$  be the composition of the two expressions, that is <sup>4</sup>

$$q_1 \circ q_2 \stackrel{\text{def}}{=} \text{for } \bar{x} \text{ in } q_1 \text{ return } q_2$$

By definition we have that  $k_{q_1 \circ q_2} = k_{q_1} + k_{q_2}$ . By Lemma 4.4.4 we can fold  $c \in \tau_{q_1 \circ q_2}$  to a  $k_{q_1+q_2}$ -chain  $c'$ , such that  $c'$  is of the form  $c'_1.c''$ , with  $c'_1$  a  $k_{q_1}$ -chain for  $q_1$  and  $c'$  a  $(k_{q_1}+k_{q_2})$ -chain for  $q_2$ , thus exhibiting a pair of folded chains as desired.

In the remaining case where some  $q'_2$  prefix of  $q_2$  does not navigate descendants of nodes returned by  $q_1$ , then either  $q'_2$  is navigating some ancestors of types returned by  $q_1$ , or  $q'_2$  is navigating in parallel a node which is also navigated by  $q_1$  (and of course in order to have a conflict, a sibling navigation should follow in the sequel). In both cases we need a multiplicity value upper bounded by the one given by the case shown above, when each  $q'_2$  prefix of path  $q_2$  navigates some descendant of  $c_1$ .

Correctness when either  $c_1$  or  $c_2$  are element chains follows, since element chains are either completely specified, and then folding is not needed (seen that label frequency is precisely kept into account by  $\mathcal{F}(,)$ ), or they have an head completely specified and tail taken from a return chain, for which we assume the thesis holds by induction.  $\square$

---

<sup>4</sup>Notice that  $q_1 \circ q_2$  may not be linear anymore, thus with a little abuse of notation we denote by  $q_1 \circ q_2$  also  $R_{\text{for}}^\bullet(q_1 \circ q_2)$

Notice that the lemma states something stronger than what is needed for query-update independence purpose. In fact, it does not only concern used and return chains, but it also takes into account element chains, even if they are not involved in the definition of independence.

As before, we generalize from  $lp$  queries to the whole language.

**Lemma 4.4.5** ( $k$ -FOLDING AND CONFLICT PRESERVATION) *Let  $d$  be a DTD,  $q_1, q_2$  two queries with at most only one free variable  $x$ , and  $\Gamma = \{x \mapsto s_d\}$  a static environment. Provided that  $\Gamma \vdash_{C_d} q_i : (r_i, v_i, e_i)$  and  $\tau_i = r_i \cup v_i \cup e_i$ , for each pair of chains  $(c_1, c_2) \in \tau_{q_1} \times \tau_{q_2}$  such that  $c_1 \leq c_2$ , there exists  $(c'_1, c'_2) \in \tau_{q_1} \times \tau_{q_2}$  such that  $c'_1 \leq c'_2$  and  $c_i \hookrightarrow_d^* c'_i$  with  $c'_i$  a  $(k_{q_1} + k_{q_2})$ -chain ( $i=1, 2$ ).*

*Proof.* Let  $LPQ(q_i)$  be the set of linear-path queries associated to  $q_i$ , and  $k'_i = \max\{k_{q'_i} \mid q'_i \in LPQ(q_i)\}$ . By Lemma 9.2.2 we have that  $\tau_{LPQ(q_i)} \supseteq \tau_{q_i}$ . Moreover, by Lemma 9.2.8 all pair of chains  $(c_1, c_2) \in \tau_{LPQ(q_1)} \times \tau_{LPQ(q_2)}$  can be folded to a pair of  $(k'_1 + k'_2)$ -chains preserving the conflict, and by Proposition 9.2.4  $k'_i \leq k_{q_i}$ . Thus all chains in  $\tau_{q_1} \times \tau_{q_2}$  can be folded into  $k_{q_1} + k_{q_2}$ -chains.  $\square$

**Remark 9.2.9.** *As a direct corollary of Lemma 4.4.5 we can approximate the problem of queries disjunction in presence of a schema just looking on query  $k$ -chains. Indeed, if no pair of conflicting  $k$ -chain is found for the input query expressions, then no nodes are selected by both input query expressions on any valid document instance. Furthermore, for the subset of queries and schemas that make our chain inference system enjoying completeness, we can exactly solve the query intersection problem in presence of a schema.*

We conclude this section by showing the equivalence of the chain analysis in the infinite and finite case.

**Theorem 9.2.10.** *Let  $d$  be a DTD,  $q$  a query and  $u$  an update, with at most only one free variable  $x$  and  $\Gamma = \{x \mapsto s\}$  a static environment. Then,  $q \perp_{C_d} u$  iff  $q \perp_{C_d^k} u$ .*

*Proof.* Direction  $\Rightarrow$ ) is straightforward since  $C_d^k \subseteq C_d$ . We show now direction  $\Leftarrow$ ) by showing that  $q \not\perp_{C_d} u$  implies  $q \not\perp_{C_d^k} u$ .

By structural induction on  $u$ .

*Base.*

$[u = ()]$  No chain is inferred for the empty update, hence it is independent wrt any query.

[ $u = \text{delete } q_0$ ] In this case  $q \not\vdash_{C_d} u$  implies

$$\text{confl}(r_q, U_u) \cup \text{confl}(U_u, r_q) \cup \text{confl}(U_u, v_q) \neq \emptyset$$

Since  $U_u = r_{q_0}$  (up to ":" separator), let  $k = k_q + k_{q_0}$ .

By lemma 4.4.5

$$\text{confl}(r_q, r_{q_0}) \neq \emptyset \quad \text{implies} \quad \text{confl}(r_q^k, r_{q_0}^k) \neq \emptyset$$

$$\text{confl}(r_{q_0}, r_q) \neq \emptyset \quad \text{implies} \quad \text{confl}(r_{q_0}^k, r_q^k) \neq \emptyset$$

$$\text{confl}(r_{q_0}, v_q) \neq \emptyset \quad \text{implies} \quad \text{confl}(r_{q_0}^k, v_q^k) \neq \emptyset$$

and therefore we conclude  $q \not\vdash_{C_d^k} u$ .

[ $u = \text{rename } q_0 \text{ as } a$ ] In this case  $q \not\vdash_{C_d} u$  implies

$$\text{confl}(r_q, U_u) \cup \text{confl}(U_u, r_q) \cup \text{confl}(U_u, v_q) \neq \emptyset$$

Recall that  $U_u = U_1 \cup U_2$  where

$$U_1 = \{ c:b \mid c.b \in r_{q_0} \} \quad \text{and} \quad U_2 = \{ c:a \mid c.b \in r_{q_0} \}$$

For all conflicting update chain belonging to  $U_1$  we conclude as for the **delete** case.

For all conflicting update chain belonging to  $U_2$  we make a case analysis on the kind of conflict.

**confl( $r_q, U_2$ )** For all  $(c_q, c_u) \in r_q \times U_2$  such that  $c_q \leq c_u$  we have two cases to consider:  $c_q < c_u$  and  $c_q = c_u$ . If  $c_q < c_u$  then  $\text{confl}(r_q, U_2) \neq \emptyset$  if and only if  $\text{confl}(r_q, U_1) \neq \emptyset$  because the conflict arises on a label preceding the renamed one. Then reasoning as for  $U_1$  we conclude  $q \not\vdash_{C_d^k} u$ . If  $c_q = c_u$ , we distinguish again if  $\text{confl}(r_q, U_1) \neq \emptyset$  or not. If  $\text{confl}(r_q, U_1) \neq \emptyset$  then we conclude  $q \not\vdash_{C_d^k} u$  as before. If  $\text{confl}(r_q, U_1) = \emptyset$  then in order for the update to be valid wrt the schema **a** is required to be a sibling of a type returned by  $q_0$ . Since  $\text{confl}(r_q, U_1) = \emptyset$  then  $c.b \notin r_0$  and this means that  $q_0$  specifies a label  $b \neq a$  in the step generating  $c.b$ . We modify  $q_0$  substituting  $b$  with  $a$  and we conclude by Lemma 4.4.5 on  $q, q_0$  noticing that adding a **a** tag in  $q_0$  raises a label frequency of 1 at most.

**confl( $U_2, r_q$ )** For all  $(c_u, c_q) \in U_2 \times r_q$  such that  $c_u \leq c_q$  we only consider:  $c_u < c_q$  since  $c_q = c_u$  has been addressed in the case of **confl( $r_q, U_2$ )**. With a case analysis on the emptiness of **confl( $U_1, r_q$ )** and by substituting **a** with **b** in  $q_0$ , we conclude  $q \not\vdash_{C_d^k} u$  as before.

**confl( $U_2, v_q$ )** Identical to **confl( $U_2, r_q$ )**.

$[u = \text{insert } q_1 \text{ pos } q_0]$  In this case  $q \not\vdash_{C_d} u$  implies

$$\text{confl}(r_q, U_u) \cup \text{confl}(U_u, r_q) \cup \text{confl}(U_u, v_q) \neq \emptyset$$

We consider the case  $\text{pos} = \text{into}$  since other cases are similar. Recall that  $U = U_1 \cup U_2$  where

$$\begin{aligned} U_1 &= \{ c:e \mid c \in r_{q_0} \quad e \in e_{q_1} \} \\ U_2 &= \{ c:\alpha.c' \mid c \in r_{q_0} \quad c_1.\alpha \in r_{q_1} \quad c_1.\alpha.c' \in C_d \} \end{aligned}$$

We consider first conflicts generated by  $U_1$

$\text{confl}(r_q, U_1)$  For all  $(c_q, c_u) \in r_q \times U_1$  such that  $c_q \leq c_u$ , let  $c_u = c:e$  we have two cases to consider depending on whether  $c_q \leq c$  or  $c < c_q$ .

If  $c_q \leq c$ , by Lemma 4.4.5 on  $q, q_0$  we can fold  $c_q$  to  $c'_q$  and  $c$  to  $c'$  with  $c'_q, c'$  both  $(k_q + k_{q_0})$ -chains. Since  $e$  is created by element construction it is completely specified and it is easy to see that  $c'.e$  is a  $(k_q + k_u)$ -chain.

If  $c < c_q$  by Lemma 4.4.5 on  $q, q_0$  we can fold  $c_q$  to  $c'_q$  and  $c$  to  $c'$  with  $c'_q, c'$  both  $(k_q + k_{q_0})$ -chains. Let  $\bar{e} \leq e$  such that  $c.\bar{e} = c_q$ , since  $e$  is completely specified by query syntax by in the element construction then we have  $c''_q = c''.\bar{e}$  as a folding of  $c_q$  that preserves  $\bar{e}$  as a suffix, with  $c''$  a folding of  $c$  as prefix, and  $c''_q$  is a  $(k_q + k_{q_0})$ -chain. Also, of course  $c''.e$  is a  $(k_q + k_u)$ -chain and we conclude  $q \not\vdash_{C_d} u$ .

$\text{confl}(U_1, r_q)$  Analogous to  $\text{confl}(r_q, U_1)$ .

$\text{confl}(U_1, v_q)$  Analogous to  $\text{confl}(r_q, U_1)$ .

Concerning  $U_2$ , we want to show that we can restrict the conflict analysis on a subset of  $U_2$  namely  $U'_2$  defined as

$$U'_2 = \{ c:\alpha \mid c \in r_{q_0} \quad c_1.\alpha \in r_{q_1} \}$$

- If  $\text{confl}(r_q, U_2) = \emptyset$  then either  $\text{confl}(r_q, U'_2)$  or  $\text{confl}(U'_2, r_q)$ .
- If  $\text{confl}(U_2, r_q) = \emptyset$  then  $\text{confl}(U'_2, r_q)$ .
- If  $\text{confl}(U_2, v_q) = \emptyset$  then  $\text{confl}(U'_2, v_q)$ .

The proof for  $\text{confl}(r, U'_2)$  follows as for the **rename** case.

$[u = \text{replace } q_0 \text{ with } q]$  The statement follows as a combination of **delete** and **insert** cases.

*Induction.*

$[u = u_1, u_2]$  By inductive hypothesis on  $u_1, u_2$ .



$[u = \text{for } x \text{ in } q \text{ return } u_1]$  By inductive hypothesis on the length of the derivation for  $q$  and then for  $u$ . Analogous to the case of Lemma 9.2.5.

$[u = \text{let } x := q \text{ return } u_1]$  Analogous to the for case.

$[u = \text{if } (q) \text{ then } u_1 \text{ else } u_2]$  By inductive hypothesis on the length of the derivation for  $q, u_1$  and  $u_2$ .

□

**Theorem 9.2.11** (Soundness of  $C_d^k$  Independence). *Let  $d$  be a DTD,  $q$  a query and  $u$  an update, with at most only one free variable  $x$ , which is always bound, during static analysis, to the root type  $s_d$  and, during query evaluation, to the root location  $l_t$  of a valid tree  $t \in d$ . Let  $k = k_q + k_u$  as defined above, then  $q \perp_{C_d^k} u$  implies  $q \Vdash_d u$*

*Proof.* Straightforward since by Theorem 9.2.10  $q \perp_{C_d} u$  iff  $q \perp_{C_d^k} u$ , and by Theorem 4.4.2  $q \perp_{C_d} u$  implies  $q \Vdash_d u$ . □

# Bibliography

- [AAD<sup>+</sup>12] Serge Abiteboul, Yael Amsterdamer, Daniel Deutch, Tova Milo, and Pierre Senellart. Finding optimal probabilistic generators for XML collections. In *ICDT*, 2012.
- [AFL02] Marcelo Arenas, Wenfei Fan, and Leonid Libkin. On verifying consistency of XML specifications. In *PODS*, 2002.
- [AYCLS01] Sihem Amer-Yahia, SungRan Cho, Laks V. S. Lakshmanan, and Divesh Srivastava. Minimization of tree pattern queries. In *SIGMOD*, 2001.
- [BBC<sup>+</sup>11] Mohamed Amine Baazizi, Nicole Bidoit, Dario Colazzo, Noor Malla, and Marina Sahakyan. Projection for XML update optimization. In *EDBT*, 2011.
- [BC09a] Michael Benedikt and James Cheney. Schema-based independence analysis for XML updates. In *VLDB*, 2009.
- [BC09b] Michael Benedikt and James Cheney. Semantics, types and effects for XML updates. In *DBPL*, 2009.
- [BC10] Michael Benedikt and James Cheney. Destabilizers and independence of XML updates. *PVLDB*, 3(1), 2010.
- [BCCN06] Véronique Benzaken, Giuseppe Castagna, Dario Colazzo, and Kim Nguyen. Type-based XML projection. In *VLDB*, 2006.
- [BCF03] V. Benzaken, G. Castagna, and A. Frisch. CDuce: an XML-centric general-purpose language. In *ICFP*, 2003.
- [Ben10] Michael Benedikt. Analysis of declarative updates: invited talk. In *EDBT/ICDT Workshops*, 2010.
- [BGMM09] Henrik Björklund, Wouter Gelade, Marcel Marquardt, and Wim Martens. Incremental XPath evaluation. In *ICDT*, 2009.
- [BK06] Michael Benedikt and Christoph Koch. Interpreting tree-to-tree queries. In *ICALP*, 2006.

- [BK09] Michael Benedikt and Christoph Koch. From XQuery to relational logics. *ACM TODS*, 34(4), 2009.
- [BNSV10] Geert Jan Bex, Frank Neven, Thomas Schwentick, and Stijn Vansummeren. Inference of concise regular expressions and DTDs. *ACM TODS*, 2010.
- [BPSM<sup>+</sup>06] Tim Bray, Jean Paoli, C. M. Sperberg-McQueen, Eve Maler, François Yergeau, and John Cowan. Extensible Markup Language (XML) 1.1 (Second Edition). Technical report, W3C Consortium, 2006.
- [BTCU10a] Nicole Bidoit-Tollu, Dario Colazzo, and Federico Ulliana. Detecting XML query-update independence. *26ème journée des Bases des données Avancées*, 2010.
- [BTCU10b] Nicole Bidoit-Tollu, Dario Colazzo, and Federico Ulliana. Detecting XML query-update independence. *International Formal Methods Workshop*, 2010.
- [BTCU12] Nicole Bidoit-Tollu, Dario Colazzo, and Federico Ulliana. Type-based detection of XML query-update independence. *PVLDB*, 5(9), 2012.
- [CCF<sup>+</sup>06] Don Chamberlin, Michael Carey, Daniela Florescu, Donald Kossmann, and Jonathan Robie. XQueryP: Programming with XQuery. In *XIME-P*, 2006.
- [CGM11] Federico Cavalieri, Giovanna Guerrini, and Marco Mesiti. Updating XML schemas and associated documents through exup. In *ICDE*, 2011.
- [CGMS06] Dario Colazzo, Giorgio Ghelli, Paolo Manghi, and Carlo Sartiani. Static analysis for path correctness of XML queries. *Journal of Functional Programming*, 16, 2006.
- [CGS11] Dario Colazzo, Giorgio Ghelli, and Carlo Sartiani. Schemas for Safe and Efficient XML Processing, 2011. Tutorial - ICDE.
- [Che08] James Cheney. FLUX: FunctionaL Updates for XML. In *ICFP*, 2008.
- [Che09] James Cheney. XQuery Update analysis tools version 0.1, 2009. <http://homepages.inf.ed.ac.uk/jcheney/programs/>.
- [CM01] James Clark and Makoto Murata. RELAX NG specification, 2001.
- [DFF<sup>+</sup>10] D Draper, P Fankhauser, M Fernandez, A Malhotra, K Rose, M Rys, J Siméon, and P Wadler. XQuery 1.0 and XPath 2.0 Formal Semantics. Technical report, World Wide Web Consortium, 2010.
- [DM06] Bruno Dutertre and Leonardo De Moura. The Yices SMT solver. Technical report, 2006.
- [FCG04] Wenfei Fan, Chee Yong Chan, and Minos N. Garofalakis. Secure XML querying with security views. In *SIGMOD*, 2004.

- [Fra05] M. Franceschet. XPathMark - An XPath benchmark for XMark generated data. In *XSym*, 2005.
- [GL08] Pierre Genevès and Nabil Layaïda. XML reasoning solver user manual. Technical report, R. Report 6726, INRIA, 2008.
- [GLS07] Pierre Genevès, Nabil Layaïda, and Alan Schmitt. Efficient static analysis of xml paths and types. *PLDI*, 2007.
- [GORS08] Giorgio Ghelli, Nicola Onose, Kristoffer Rose, and Jerome Simeon. XML query optimization in the presence of side effects. In *SIGMOD*, 2008.
- [Gra03] Hervé Grall. *Deux critères de sécurité pour l'exécution de code mobile*. Thèse, École des Ponts ParisTech, 2003.
- [GRS06] Giorgio Ghelli, Christopher Ré, and Jérôme Siméon. XQuery!: An XML query language with side effects. In *EDBT Workshops*, 2006.
- [GRS08] Giorgio Ghelli, Kristoffer Høgsbro Rose, and Jérôme Siméon. Commutativity analysis for XML updates. *ACM TODS*, 2008.
- [GW97] Roy Goldman and Jennifer Widom. Dataguides: Enabling query formulation and optimization in semistructured databases. In *VLDB*, 1997.
- [Hid03] Jan Hidders. Satisfiability of XPath expressions. In *DBPL*, 2003.
- [HKL05] Beda Christoph Hammerschmidt, Martin Kempa, and Volker Linnemann. On the intersection of XPath expressions. In *IDEAS*, 2005.
- [HU00] John E. Hopcroft and Jeffrey D. Ullman. *Introduction to Automata Theory, Languages and Computation, Second Edition*. Addison-Wesley, 2000.
- [HVP05] Haruo Hosoya, Jérôme Vouillon, and Benjamin C. Pierce. Regular expression types for XML. *ACM Trans. Program. Lang. Syst.*, 27(1), 2005.
- [JGL12] Muhammad Junedi, Pierre Genevès, and Nabil Layaïda. XML query-update independence analysis revisited. *DocEng*, 2012.
- [KM11] Nils Klarlund and Anders Moller. MONA version 1.4 user manual. Technical report, BRICS, 2011.
- [MLMK05] Makoto Murata, Dongwon Lee, Murali Mani, and Kohsuke Kawaguchi. Taxonomy of XML schema languages using formal language theory. *ACM Trans. Internet Techn.*, 5(4):660–704, 2005.
- [MS03] A. Marian and J. Siméon. Projecting XML documents. In *VLDB*, 2003.
- [MS04] Gerome Miklau and Dan Suciu. Containment and equivalence for a fragment of XPath. *Journal of ACM*, 2004.

- [PAG10] Jorge Pérez, Marcelo Arenas, and Claudio Gutierrez. nSPARQL: A navigational language for RDF. *Journal of Web Semantics*, 2010.
- [PV00] Yannis Papakonstantinou and Victor Vianu. DTD inference for views of XML data. In *PODS*, 2000.
- [RCD<sup>+</sup>11] J Robie, D Chamberlin, M Dyck, D Florescu, J Milton, and J Simeon. XQuery update facility 1.0. Technical report, W3C Consortium, 2011.
- [RS06] Mukund Raghavachari and Oded Shmueli. Conflicting XML updates. In *EDBT*, 2006.
- [SCF<sup>+</sup>07] Jérôme Siméon, Don Chamberlin, Daniela Florescu, Scott Boag, Mary F. Fernández, and Jonathan Robie. XQuery 1.0: An XML query language. W3C recommendation, W3C, January 2007. <http://www.w3.org/TR/2007/REC-xquery-20070123/>.
- [SM73] L. J. Stockmeyer and A. R. Meyer. Word problems requiring exponential time (preliminary report). In *STOC*, 1973.
- [SWK<sup>+</sup>02] A. Schmidt, F. Waas, M. L. Kersten, M. J. Carey, I. Manolescu, and R. Busse. XMark: A benchmark for XML data management. In *VLDB*, 2002.
- [TBMM04] Henry S. Thompson, David Beech, Murray Maloney, and Noah Mendelsohn. XML Schema Part 1: Structures Second Edition. Technical report, World Wide Web Consortium, Oct 2004. W3C Recommendation.
- [Vor96] Sergei G. Vorobyov. An improved lower bound for the elementary theories of trees. In *CADE*, 1996.