

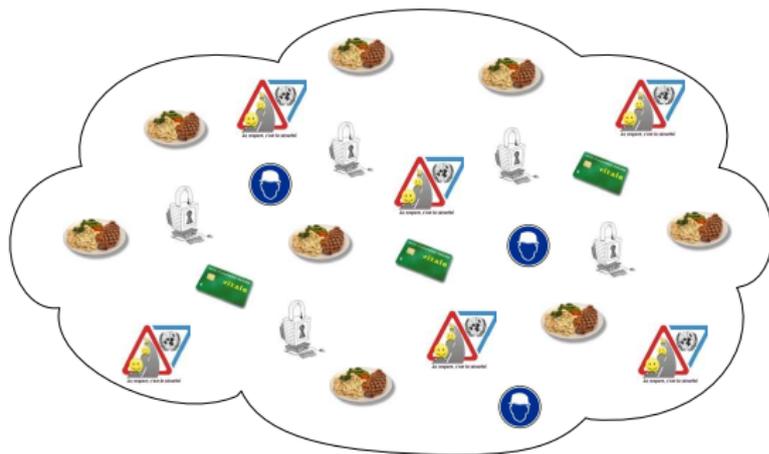
# Apprentissage de connaissances morphologiques pour l'acquisition automatique de ressources lexicales

Delphine Bernhard

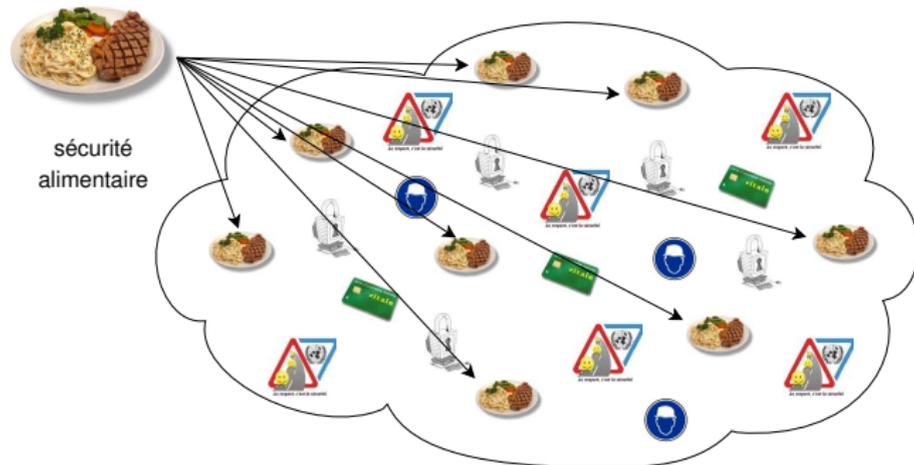
Laboratoire TIMC-IMAG, Grenoble

30 novembre 2006

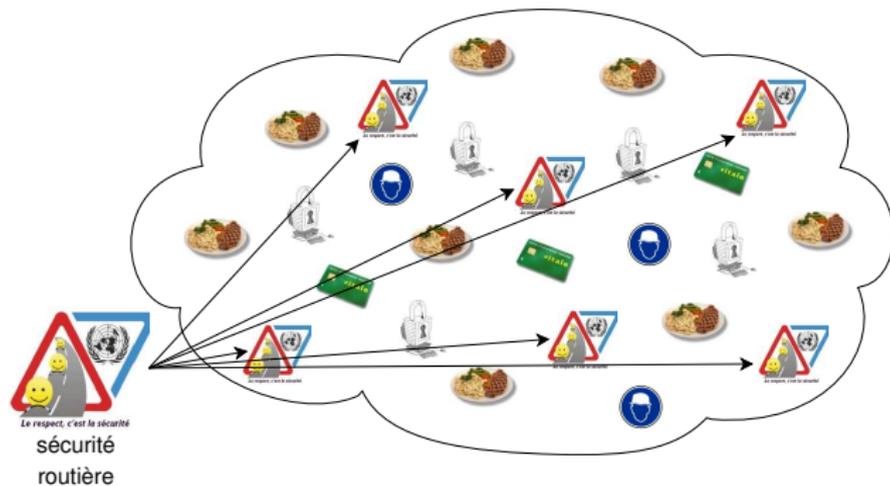
## Sécurité ?



## Sécurité ?

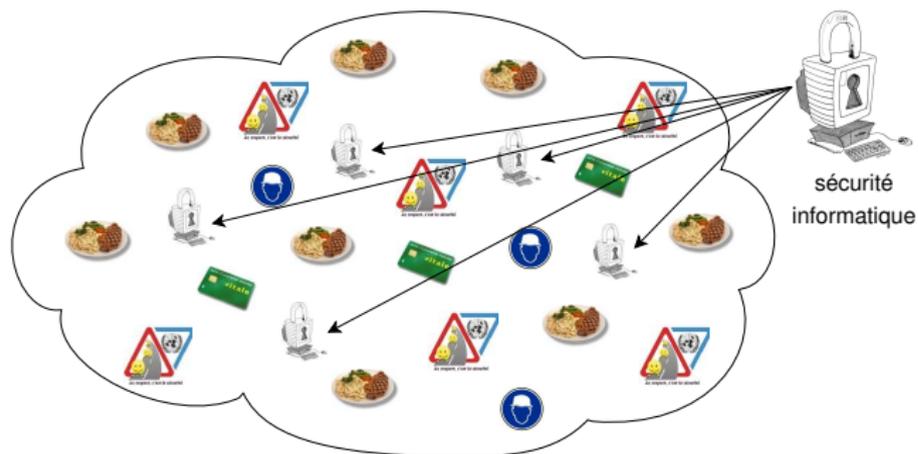


## Sécurité ?

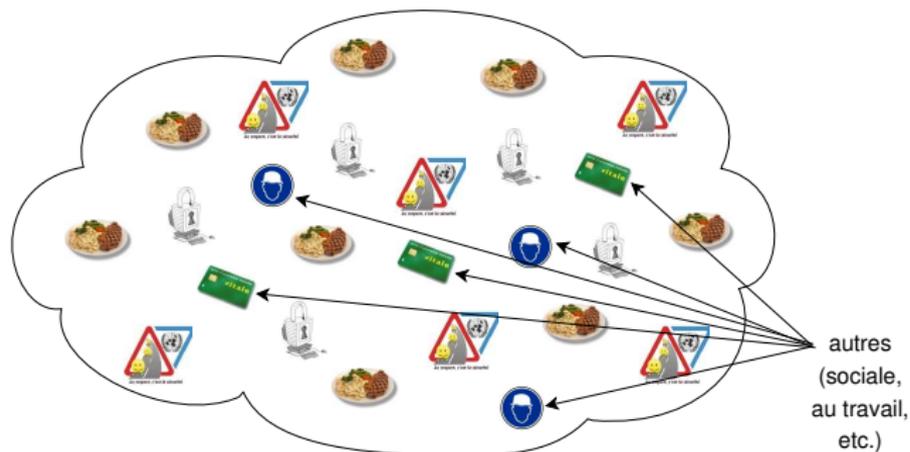


# Contexte : organisation des données

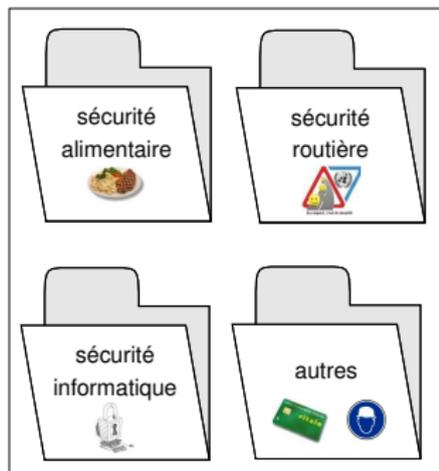
Sécurité ?



## Sécurité ?



## Sécurité ?



Organisation et structuration des informations via des ressources décrivant et classant les connaissances

# Quelles ressources pour décrire les connaissances ?

Ressources construites manuellement par des experts

dictionnaires, terminologies, thésaurus, ontologies

+ contrôle, précision

– coût, couverture

Descriptions collaboratives

indexation sociale (folksonomie) : tags ~ mots-clés

+ faible coût, couverture

– imprécision, absence de contrôle

Construction automatique de ressources

- ▶ Identification des termes représentant les concepts
- ▶ Acquisition de relations sémantiques pour la structuration des connaissances

## Patrons

- ▶ Termes complexes  
N + Adj  
éruption volcanienne
- ▶ Relations sémantiques  
SN tels que SN+ ou SN  
des phénomènes climatiques  
**tels que** la température ou  
les précipitations

## Statistiques

- ▶ Termes simples  
Mesures de comparaison  
Log du rapport de  
vraisemblance
- ▶ Termes complexes  
Mesures d'association  
Information mutuelle,  $\chi^2$ ,  
coefficient de Jaccard, etc.
- ▶ Similarité contextuelle  
Comparaison de vecteurs de  
co-occurrence

## Morphologie

Étude des morphèmes = les plus petites unités linguistiques qui allient forme et sens

## Morphologie

Étude des morphèmes = les plus petites unités linguistiques  
qui allient forme et sens

pneumonoultramicroscopicsilicovolcanoconiosis

## Morphologie

Étude des morphèmes = les plus petites unités linguistiques  
qui allient forme et sens

pneumonoultramicroscopicsilicovolcanoconiosis



## Morphologie

Étude des morphèmes = les plus petites unités linguistiques  
qui allient forme et sens

pneumonoultramicroscopicsilicovolcanoconiosis

extrême



## Morphologie

Étude des morphèmes = les plus petites unités linguistiques  
qui allient forme et sens

pneumonoultramicroscopicsilicovolcanoconiosis



microscopique

## Morphologie

Étude des morphèmes = les plus petites unités linguistiques  
qui allient forme et sens

pneumonoultramicroscopicsilicovolcanoconiosis

silicium



## Morphologie

Étude des morphèmes = les plus petites unités linguistiques  
qui allient forme et sens

pneumonoultramicroscopicsilicovolcanoconiosis

volcan



## Morphologie

Étude des morphèmes = les plus petites unités linguistiques  
qui allient forme et sens

pneumonoultramicroscopicsilicovolcanoconiosis

poussière



## Morphologie

Étude des morphèmes = les plus petites unités linguistiques  
qui allient forme et sens

pneumonoultramicroscopicsilicovolcanoconiosis

atteinte



## Morphologie

Étude des morphèmes = les plus petites unités linguistiques qui allient forme et sens

pneumonoultramicroscopicsilicovolcanoconiosis

maladie des poumons résultant de l'inhalation de poussières de silicium très fines produites par des volcans

## Morphologie

Étude des morphèmes = les plus petites unités linguistiques qui allient forme et sens

pneumonoultramicroscopicsilicovolcanoconiosis

maladie des poumons résultant de l'inhalation de poussières de silicium très fines produites par des volcans

## Identification de termes

Vocabulaire technique : utilisation fréquente de morphèmes caractéristiques comme méga+, micro+, +gramme ou +graphe.

## Morphologie

Étude des morphèmes = les plus petites unités linguistiques qui allient forme et sens

pneumonoultramicroscopicsilicovolcanoconiosis

maladie des poumons résultant de l'inhalation de poussières de silicium très fines produites par des volcans

## Identification de termes

Vocabulaire technique : utilisation fréquente de morphèmes caractéristiques comme méga+, micro+, +gramme ou +graphe.

## Extraction de relations sémantiques

cobalto**thérapie** est un type de **thérapie**

## Objectifs

Intégration de la morphologie dans le processus d'acquisition automatique de ressources lexicales à partir de textes

**Il est nécessaire de disposer de ressources morphologiques**

## Méthodologie et matériel

- ▶ Travail sur corpus
- ▶ Langue de spécialité
- ▶ Apprentissage et approche statistique
- ▶ Indépendance aux langues

Données textuelles



Internet

# Schéma global

Données textuelles



Internet



Corpus de  
textes de  
spécialité

# Schéma global

Données textuelles



Internet



Corpus de  
textes de  
spécialité



Liste des mots  
du corpus

# Schéma global

Données textuelles

Analyse morphologique  
non supervisée



Internet



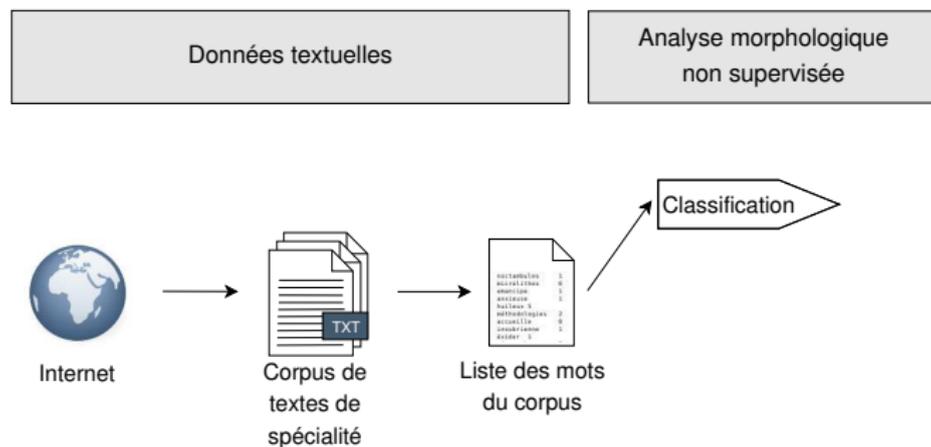
Corpus de  
textes de  
spécialité



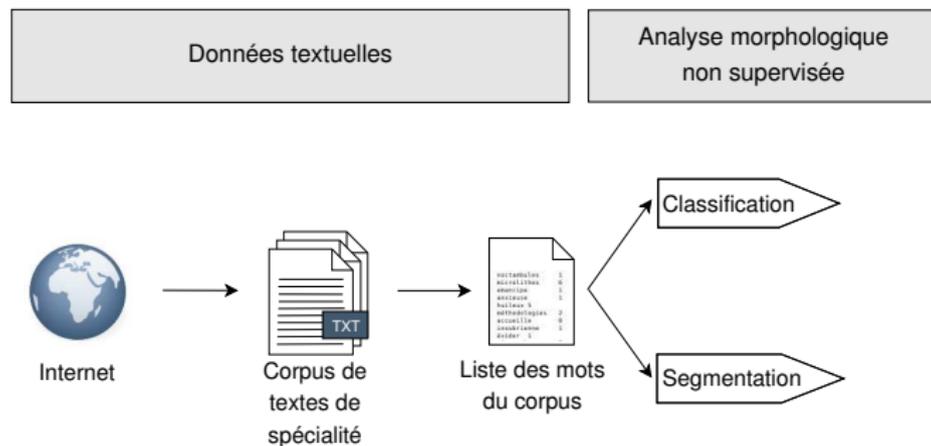
octobre	1
novembre	0
decembre	1
janvier	1
fevrier	1
mars	2
avril	0
mai	1
juin	1

Liste des mots  
du corpus

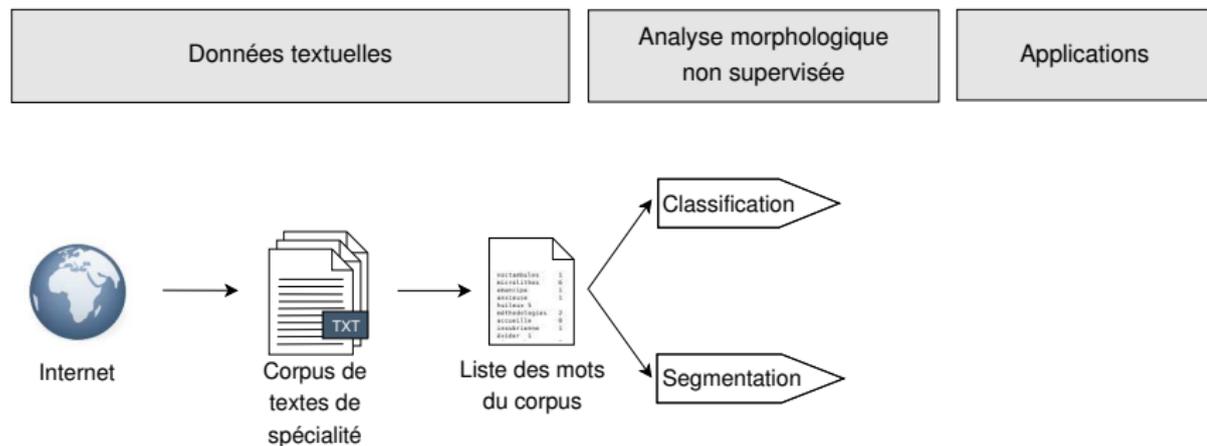
# Schéma global



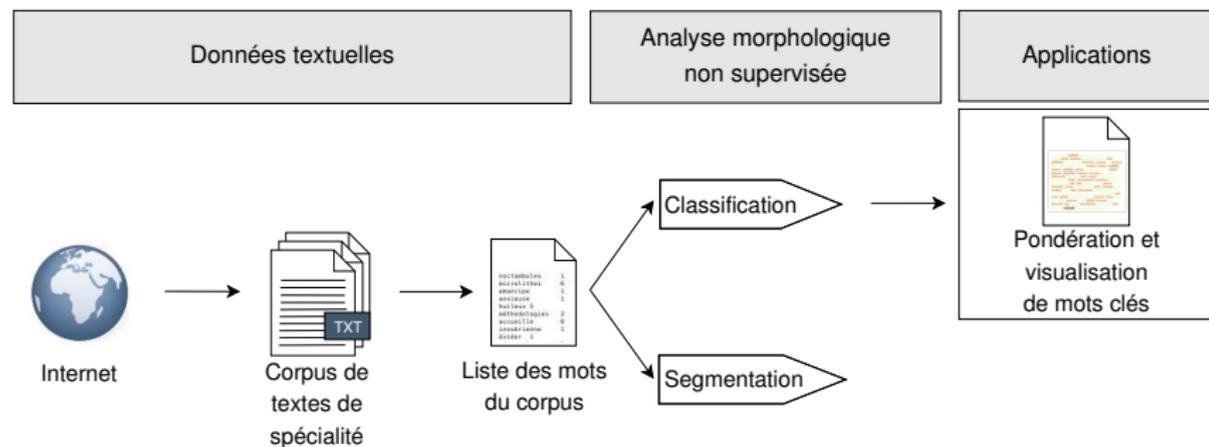
# Schéma global



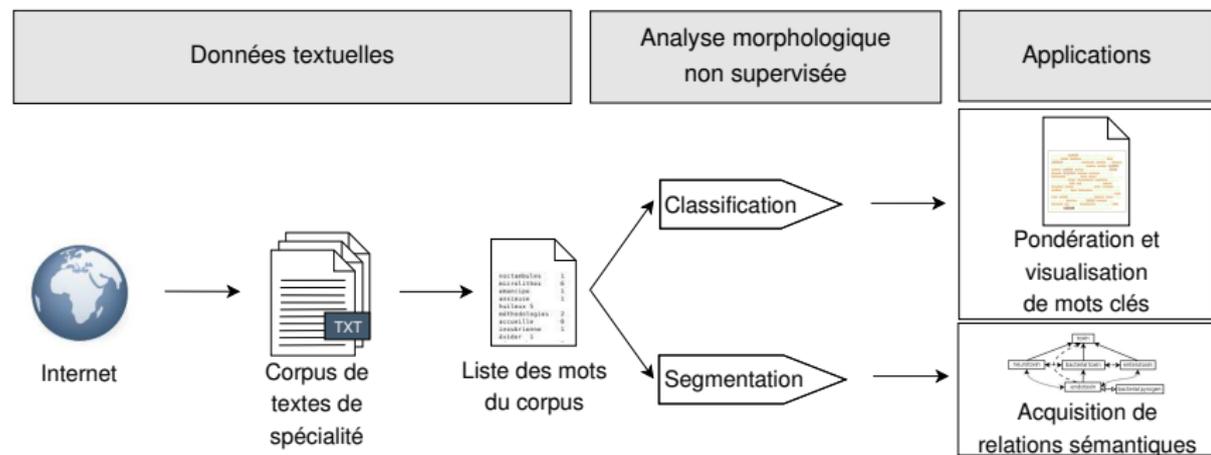
# Schéma global



# Schéma global



# Schéma global



Contexte et objectifs

## Apprentissage de connaissances morphologiques

- Construction de corpus

- Analyse morphologique par segmentation

- Analyse morphologique par classification

## Exploitation des résultats

- Pondération et visualisation de mots clés

- Acquisition de relations sémantiques

## Conclusion et perspectives

Contexte et objectifs

Apprentissage de connaissances morphologiques

- Construction de corpus

- Analyse morphologique par segmentation

- Analyse morphologique par classification

Exploitation des résultats

- Pondération et visualisation de mots clés

- Acquisition de relations sémantiques

Conclusion et perspectives

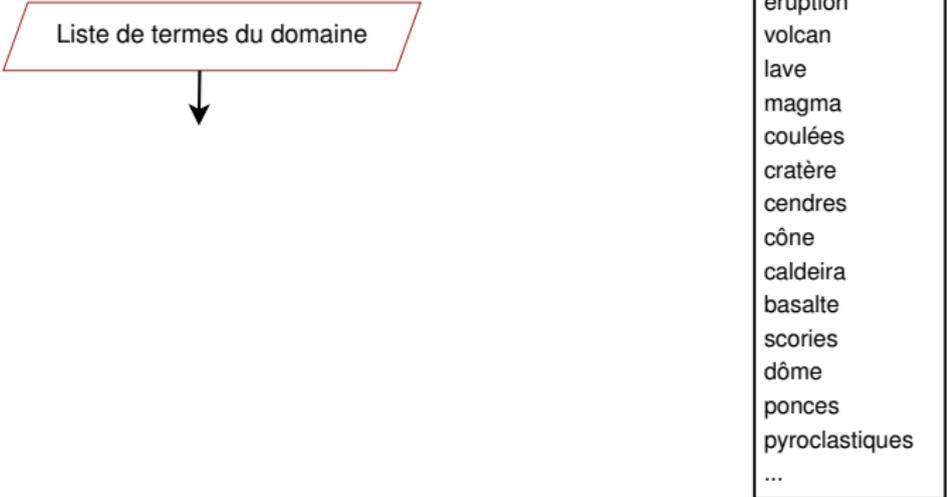
## Corpus existant

Surtout généralistes, inadaptés au vocabulaire technique

## Construction de corpus

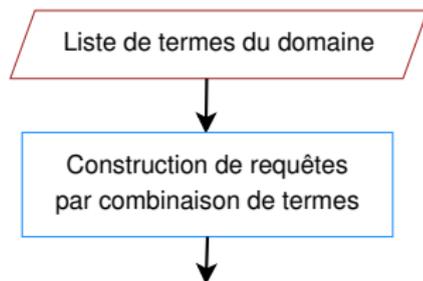
- ▶ Manuelle : processus long
- ▶ Automatique :
  - ▶ Source : le Web
  - ▶ Outils : inspirés de l'approche BootCat [Baroni et Bernardini, 2004]

Liste de termes du domaine



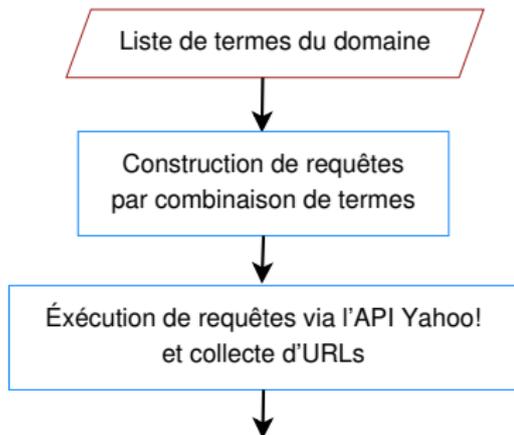
```
graph TD; A[/Liste de termes du domaine/] --> B[éruption  
volcan  
lave  
magma  
coulées  
cratère  
cendres  
cône  
caldeira  
basalte  
scories  
dôme  
ponces  
pyroclastiques  
...];
```

éruption  
volcan  
lave  
magma  
coulées  
cratère  
cendres  
cône  
caldeira  
basalte  
scories  
dôme  
ponces  
pyroclastiques  
...



éruption volcan lave  
éruption volcan magma  
magma coulées cratère  
magma coulées cendres  
magma coulées cône  
coulées cône caldeira  
cône caldeira basalte  
cône caldeira scories  
cône caldeira dôme  
scories dôme ponces  
scorie dôme pyroclastiques  
...

# Collecte d'URLs



éruption volcan lave

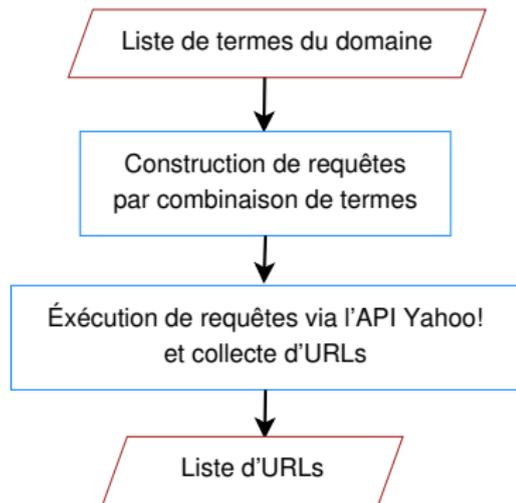


**YAHOO!**



[www.volcans.ch/pages/minute30\\_2002.html](http://www.volcans.ch/pages/minute30_2002.html)  
[www.runisland.com/volcan.html](http://www.runisland.com/volcan.html)  
[fr.wikipedia.org/wiki/Volcan](http://fr.wikipedia.org/wiki/Volcan)  
[www.volcanogeol.com/hawaii/lave.htm](http://www.volcanogeol.com/hawaii/lave.htm)  
[users.skynet.be/lave.belgique](http://users.skynet.be/lave.belgique)  
[www.fournaise.info](http://www.fournaise.info)

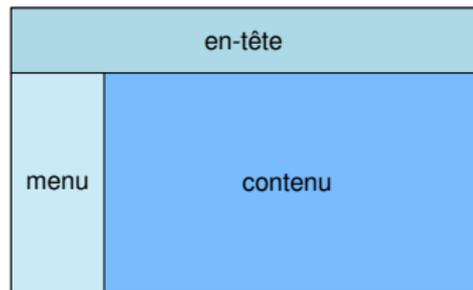
# Collecte d'URLs



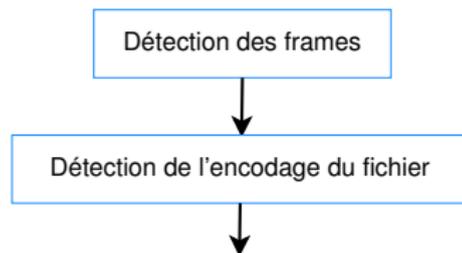
[www.volcans.ch/pages/minute30\\_2002.html](http://www.volcans.ch/pages/minute30_2002.html)  
[www.runisland.com/volcan.html](http://www.runisland.com/volcan.html)  
[fr.wikipedia.org/wiki/Volcan](http://fr.wikipedia.org/wiki/Volcan)  
[www.volcanogeol.com/hawaii/lave.htm](http://www.volcanogeol.com/hawaii/lave.htm)  
[users.skynet.be/lave.belgique](http://users.skynet.be/lave.belgique)  
[www.fournaise.info](http://www.fournaise.info)  
[www.ipgp.jussieu.fr/~aestp7/2002sicile.html](http://www.ipgp.jussieu.fr/~aestp7/2002sicile.html)  
[site.voila.fr/volcan](http://site.voila.fr/volcan)  
[www.volcans.info/juillet\\_2001.htm](http://www.volcans.info/juillet_2001.htm)  
[www.volcans2003.com/html/fiche/fiche1.htm](http://www.volcans2003.com/html/fiche/fiche1.htm)  
[www.vulcania.com/fr/reperes-volcaniques-66.html](http://www.vulcania.com/fr/reperes-volcaniques-66.html)  
[www.univ-ubs.fr/ecologie/volcanisme.html](http://www.univ-ubs.fr/ecologie/volcanisme.html)  
[www.volcan-actif.com/eruptions.htm](http://www.volcan-actif.com/eruptions.htm)  
[fr.rian.ru/russia/20050715/40914661.html](http://fr.rian.ru/russia/20050715/40914661.html)  
...

# Collecte et pré-traitement des fichiers

Détection des frames

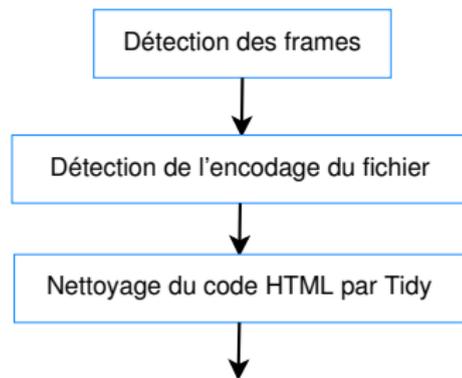


# Collecte et pré-traitement des fichiers



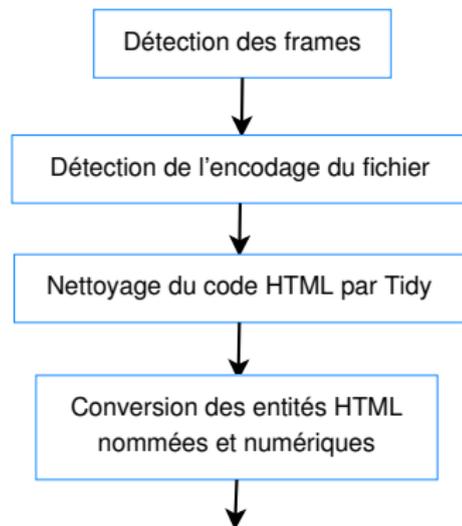
```
<html lang="fr">
<head>
<meta http-equiv="Content-Type"
content="text/html;
charset=iso-8859-1">
<title>...</title>
```

# Collecte et pré-traitement des fichiers

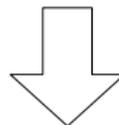


```
Warning: missing </h3> before <form>  
Warning: inserting implicit <font>  
Warning: discarding unexpected </font>  
Warning: discarding unexpected </h3>  
Warning: <spacer> is not approved by W3C  
Error: <csobj> is not recognized!  
221 warnings, 6 errors were found!
```

# Collecte et pré-traitement des fichiers

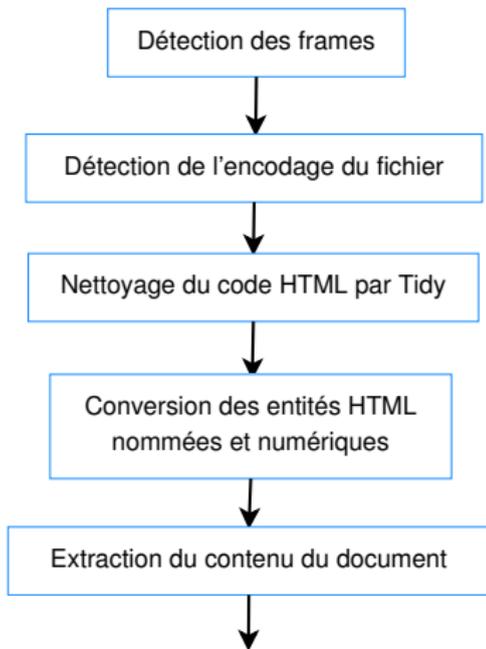


C'**était**, à l'encoignure de la rue de la Michodière et de la rue Neuve-Saint-Augustin, un magasin de nouveautés dont les **étalages** **éclataient** en notes vives, dans la douce et **pâle** journée d'octobre.



C'était, à l'encoignure de la rue de la Michodière et de la rue Neuve-Saint-Augustin, un magasin de nouveautés dont les **étalages** **éclataient** en notes vives, dans la douce et **pâle** journée d'octobre.

# Collecte et pré-traitement des fichiers



Une encyclopédie est un ouvrage où l'on trouve de toutes les pages exhaustives de l'ensemble du savoir humain. Par extension, le mot désigne également un ouvrage qui traite systématiquement d'un domaine de connaissance en particulier. Le mot encyclopédie vient du grec ancien ἐγκυκλιῶν περὶ τῆς ἀριθμητικῆς ἐπιστήμης (encyclopaedia) « résumé de toutes les sciences ». Il dérive d'encyclopaïde par une erreur de transcription, puis au XIX<sup>e</sup> siècle les savants le latinisèrent en encyclopædia.

Annexes Occidentales

[Liste alphabétique](#)

Cet article chez Wikipédia

#### Droit algérien

[Droit algérien des affaires](#) [Droit fiscal algérien](#) [Internationalisation](#)

#### La Suite OpenPortal - OIP

[Liste d'offres pour la gestion du Droit Individuel à la Formation](#) [www.OpenPortal](#)

voir aussi

Sinonyme:

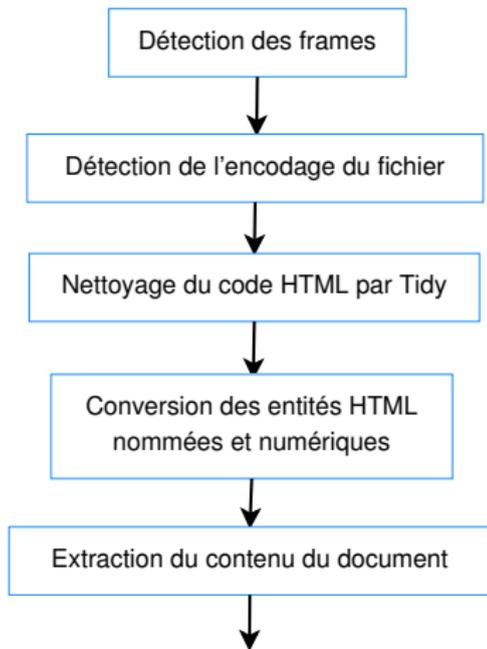
<b>Annexes Occidentales</b>
<a href="#">Droit des affaires</a>
<a href="#">Mémoires sciences</a>
<a href="#">Colloque d'art</a>
<a href="#">Droit algérien</a>
<a href="#">Droit</a>

Source de l'article : [Wikipédia](#). Le contenu est disponible selon les termes de la [Licence de documentation libre \(DCL\)](#).

Sciences naturelles et mathématiques • Chimie • Écologie • Mathématiques • Physique • Sciences de la Terre • Sciences de l'Univers • Sciences de la Vie • Statistiques • Sciences humaines • Anthropologie • Archéologie • Éducation • Géographie • Histoire • Management • Langues et linguistiques • Neurologie • Philosophie • Psychologie • Sciences exactes • Sciences de l'information et des bibliothèques • Sociologie • Politique, droit et société • Associations et organismes • Commerce • Culture et sécurité • Droit • Économie • Entreprise • Famille • Équation • Section de l'environnement • Médias • Théologie • Urbanisme • Paléontologie • Préhistoire • Agriculture • Christianisme • Économie • Médecine • Mythologie • Religion • Soudan • Spiritualité • Théologie • Arts et culture • Art • Art visuel • Arts du spectacle • Cinéma • Culture populaire • Danse • Littérature • Média • Musique • Techniques et sciences appliquées • Admiration • Agriculture • Architecture • Communication • Économie • Éducation • Énergie • Informatique • Internet • Ingénierie • Médecine • Technologie • Télécommunications • Transport • Vie quotidienne et loisirs • Biologie • Cuisine • Développement • Jardinage • Jeu • Substrat • Santé • Sport • Tourisme • Divers • Liste des lutes • Liste des pays du monde • Biographie • Arts appliqués

JAPANESE ENCYCLOPEDIA | WORLD ENCYCLOPEDIA | ENCYCLOPEDIA FRANÇAISE | ENCYCLOPEDIA PORTUGUESA | ENCYCLOPEDIA SPANOLA | ENCYCLOPEDIA THAIANA | ENCYCLOPEDIA TAMILIANA | ENCYCLOPEDIA URDU | ENCYCLOPEDIA VIETNAMESE

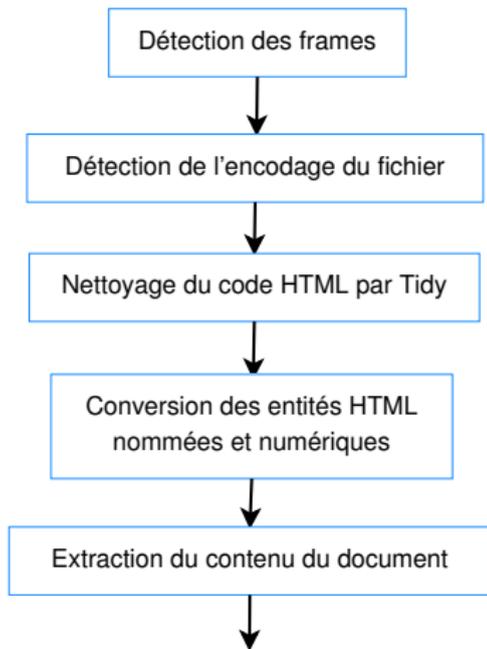
# Collecte et pré-traitement des fichiers



The screenshot shows the French Wikipedia homepage with the following elements:

- Header: **encyclopaedic.net** and **L'ENCYCLOPÉDIE FRANÇAISE**
- Navigation: [encyclopaedic française](#), [page d'accueil](#), [index alphabétique](#)
- Search: **Nom du site** (input field)
- Text: "Une encyclopédie est un ouvrage où figurent des connaissances humaines. Par extension, le mot désigne également un ouvrage qui traite systématiquement d'un domaine de connaissance en particulier. Le mot encyclopédie vient du grec ancien ἐγκυκλιῶν περὶ τῆς ἀγωγῆς τῆς παιδείας ἐν ἐγκυκλοπαιδίᾳ." (An encyclopaedia is a work where human knowledge is listed. By extension, the word also designates a work that systematically treats a particular field of knowledge. The word encyclopaedia comes from the ancient Greek ἐγκυκλιῶν περὶ τῆς ἀγωγῆς τῆς παιδείας ἐν ἐγκυκλοπαιδίᾳ.)
- Section: **Articles connexes** (with link [voir aussi](#))
- Section: **Droit algérien** (with link [lire également des affaires Droit fiscal algérien](#))
- Section: **La Suite OpenPortal - OIP** (with link [Lire l'étage pour la gestion du Droit Individuel à la Formation](#))
- Section: **voir aussi** (with link [Sémantique](#))
- Section: **Articles liés** (with links: [Droit des affaires](#), [Matières sciences](#), [Colloque d'ad](#), [Droit français](#), [Droit](#))
- Text: "Source de l'article : [Wikipédia](#). Le contenu est disponible selon les termes de la [Licence de documentation libre \(DCL\)](#)"
- Footer: [JAPANESE ENCYCLOPEDIA](#) | [WORLD ENCYCLOPEDIA](#) | [ENCYCLOPEDIA FRANÇAISE](#) | [ENCICLOPEDIA ESPAÑOLA](#) | [ENCICLOPEDIA PORTUGUESA](#) | [SHILSHI ENCYCLOPEDIA](#)

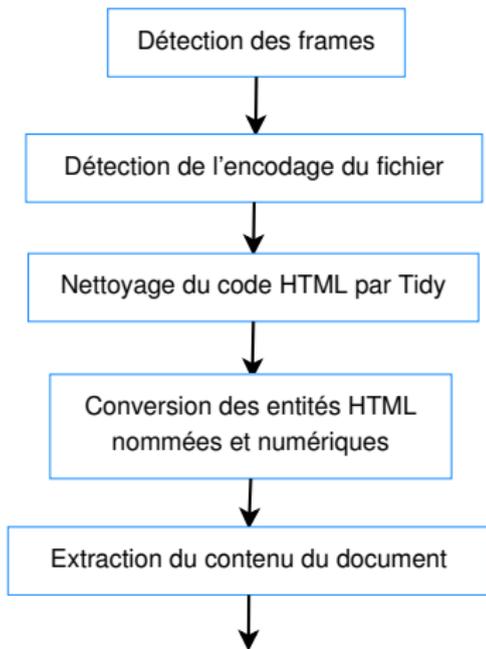
# Collecte et pré-traitement des fichiers



The screenshot shows the website **encyclopaedic.net**, which is a French encyclopedia. The page features a navigation bar with links for "Nom du site", "encyclopaedic française", "page d'accueil", and "index alphabétique". Below the navigation bar, there is a search bar and a section titled "Ejecta" with a description of volcanic ejecta. A "Publicités" (Advertisements) box is visible on the page. The footer contains a list of categories such as "Sciences naturelles et mathématiques", "Chimie", "Écologie", "Mathématiques", "Physique", "Sciences de la Terre", "Sciences de l'Univers", "Sciences de la Vie", "Statistiques", "Sciences humaines", "Andropogone", "Archéologie", "Éducation", "Géographie", "Histoire", "Management", "Langues et Linguistique", "Neurologie", "Pédagogie", "Philosophie", "Psychologie", "Sciences exactes", "Sciences de l'Information et des Bibliothèques", "Sociologie", "Politique", "Droit et Justice", "Associations et Organismes", "Commerce", "Diplomatie et Sécurité", "Droit", "Économie", "Entreprise", "Familia", "Gastronomie", "Santé et Bien-être", "Sciences de l'Environnement", "Médias", "Histoire et Littérature", "Philosophie et Théologie", "Agriculture", "Alchimie", "Christianisme", "Économie", "Mysticisme", "Mythologie", "Religion", "Santé", "Spiritualité", "Théologie", "Arts et Culture", "Art", "Art visuel", "Arts du spectacle", "Cinéma", "Culture populaire", "Danse", "Littérature", "Média", "Musique", "Techniques et Sciences appliquées", "Adaptabilité", "Agriculture", "Architecture", "Communication", "Économie", "Éducation", "Énergie", "Informatique", "Internet", "Ingénierie", "Médecine", "Technologie", "Télécommunications", "Transport", "Vie quotidienne et loisirs", "Biologie", "Cuisine", "Développement", "Jardinage", "Jeux", "Hobbies", "Santé", "Sport", "Tourisme", "Divers", "Liste des livres", "Liste des pays du monde", "Biographies", "Arts appliqués".



# Collecte et pré-traitement des fichiers

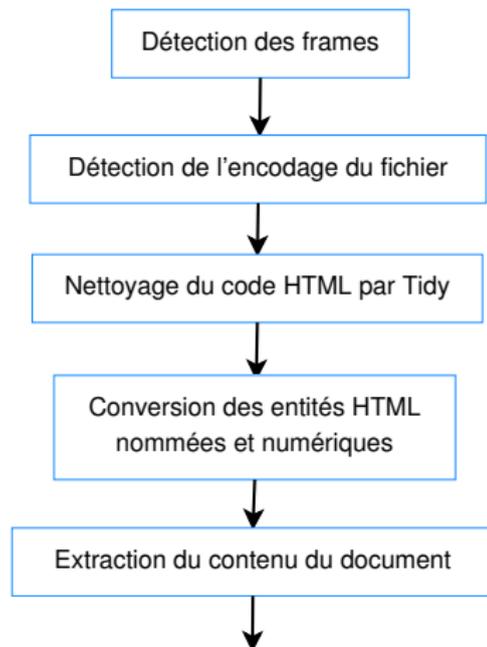


The screenshot shows the homepage of encyclopaedic.net, titled 'L'ENCYCLOPÉDIE FRANÇAISE'. The page features a search bar, navigation links, and various content sections. Annotations in colored boxes highlight specific elements:

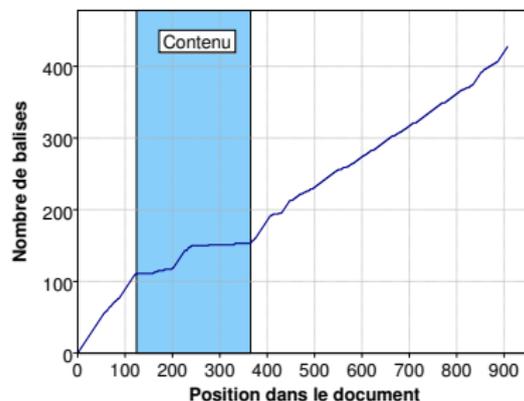
- Nom du site**: Points to the site title.
- Publicités**: Points to the 'Ejecta' section, which contains text about volcanic eruptions.
- Liens**: Points to the 'Liens' section, which lists various topics like 'Droit algérien' and 'La Suite OpenPortal - IIR'.
- Publicités**: Points to a small table of links in the 'Annonces gratuites' section.
- Informations légales**: Points to the footer area containing legal information and a list of related encyclopedias.



# Collecte et pré-traitement des fichiers

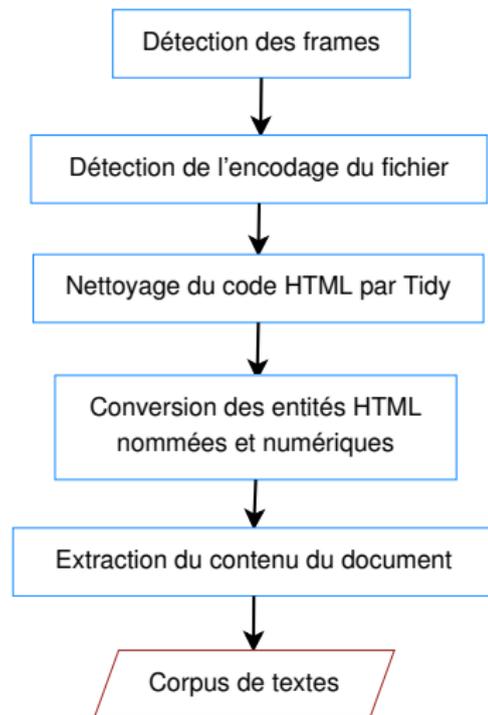


Méthode proposée par [Finn et al., 2001] :  
extraction de la sous-partie du document où la  
densité des mots est importante





# Collecte et pré-traitement des fichiers



Nombre de formes différentes		
	Anglais	Français
Cancer du sein	86 149	46 898
Volcanologie	47 789	59 768

## Deux approches

1. Segmentation : découpage des mots en segments morphémiques étiquetés
2. Classification : regroupement des mots dans des familles morphologiques

## Contraintes

- ▶ Prise en compte des procédés de formation suivants :
  - ▶ flexion : **carcinome carcinomes**
  - ▶ dérivation : **carcinome carcinomateux**
  - ▶ composition : **carcinome hépatocarcinome**
- ▶ Méthode utilisable pour d'autres langues que l'anglais et le français et pour divers domaines

## Données

Liste de mots

## Étapes

1. Apprentissage de préfixes et de suffixes
2. Acquisition de bases
3. Segmentation des mots par alignement et comparaison
4. Sélection de la meilleure segmentation

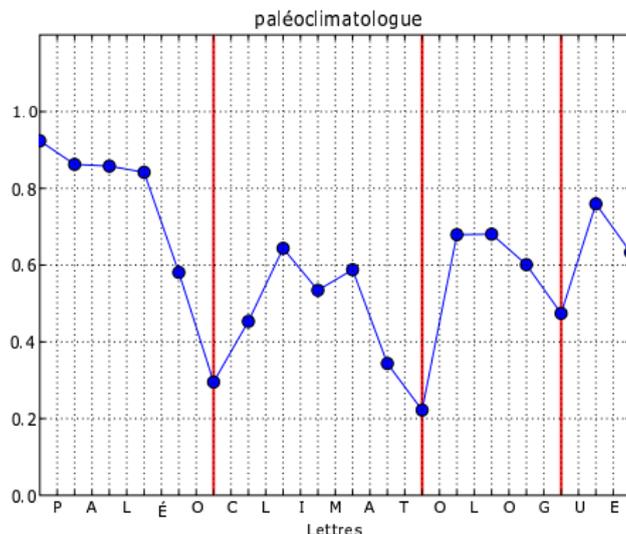
Entrée

Mots les  
plus longs

# Apprentissage de préfixes et de suffixes [1]

## Localisation de segments

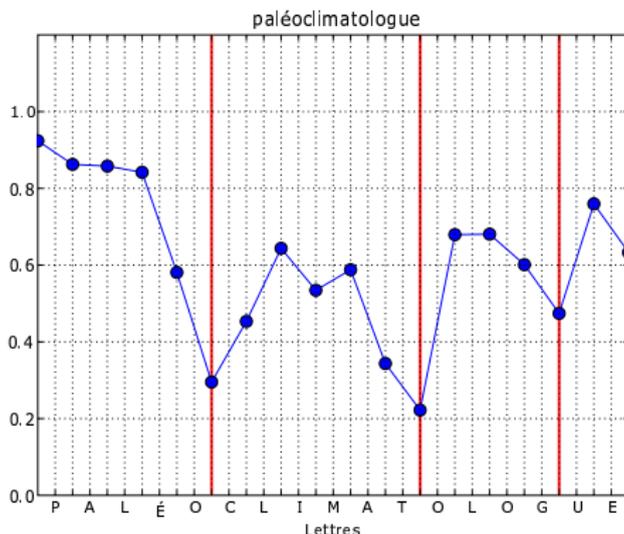
Entrée  
Mots les  
plus longs



# Apprentissage de préfixes et de suffixes [1]

## Localisation de segments

Entrée  
Mots les  
plus longs



Sortie  
Segments

# Apprentissage de préfixes et de suffixes [2]

Identification d'une base parmi les segments

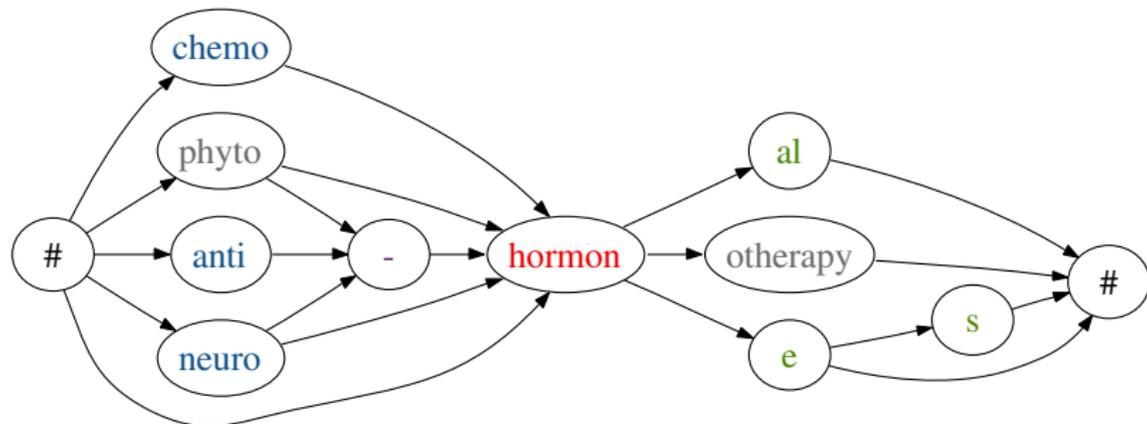
	paléo	climat	olog	ue
fréquence	68	> 17 <	288	1 348
longueur	5	< 6 >	4	2

Préfixes et suffixes

paléo		s
		isation
		isés
	climat	s
paléo		ologue
ac		ation
dendro		ologie

Retranchement des préfixes et des suffixes de  
tous les mots

# Alignement des segments de mots [1]



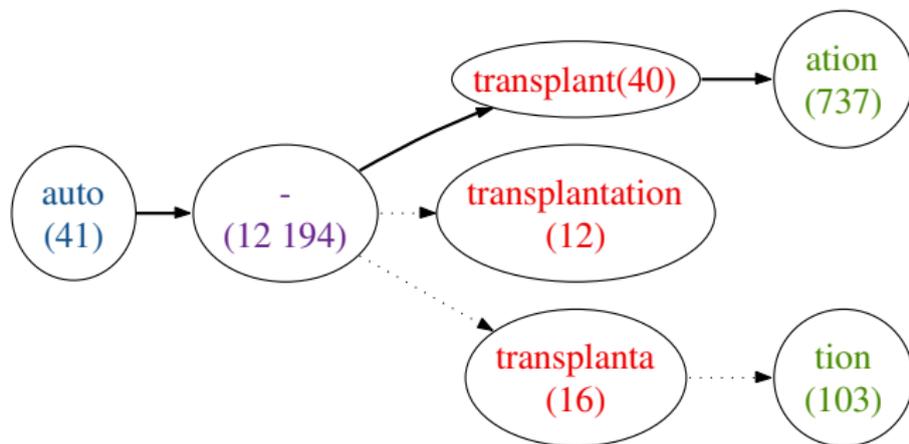
# Alignement des segments de mots [2]

## Validation des préfixes et suffixes inconnus

Mots	Suffixes connus $A_1$	Bases potentielles $A_2$	Nouveaux suffixes $A_3$
hormonal hormonotherapy hormone hormones	-al  -e	  -otherapy	   -es

$$\frac{|A_1| + |A_2|}{|A_1| + |A_2| + |A_3|} \geq a \text{ et } \frac{|A_1|}{|A_1| + |A_2|} \geq b$$

# Sélection de la meilleure segmentation



# Segmentation des mots absents du corpus d'apprentissage

- ▶ Sélection des segments qui minimisent le coût global
- ▶ Fonctions de coût utilisées :

$$\text{coût}_1(s_i) = -\log \frac{f(s_i)}{\sum_i f(s_i)}$$

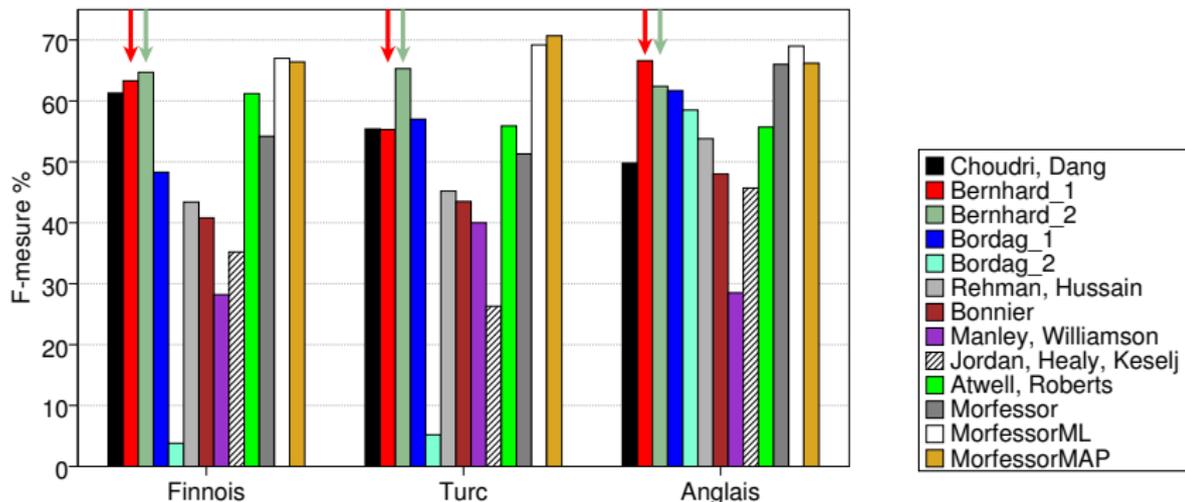
$$\text{coût}_2(s_i) = -\log \frac{f(s_i)}{\max_i [f(s_i)]}$$

# Exemples

Mots	Segmentations
allogreffe	allo + greffe
autogreffe	auto + greffe + ion
post-greffe	post + - + greffe
re-greffe	re + greffe + ion
greffe	trans + pla + n + t
greffe	trans + plant + ion
greffes	trans + plant + ion + s
greffé	trans + plant + é
greffées	trans + plant + é + e + s
greffés	trans + plant + é + s
greffes	trans + pla + n + t + s

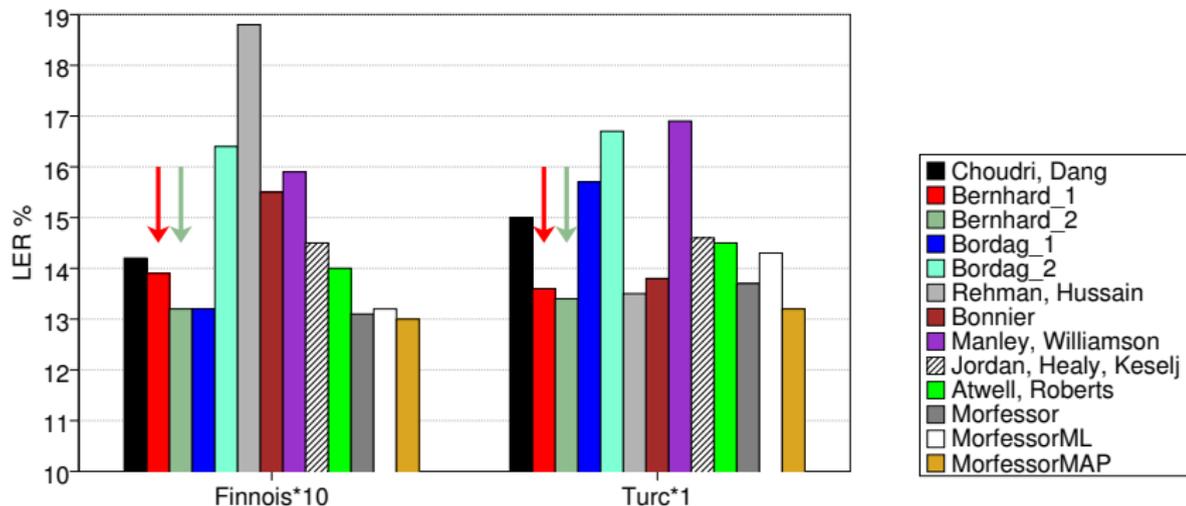
# Évaluation 1 : Morpho Challenge

## Compétition 1 : évaluation des segmentations



# Évaluation 1 : Morpho Challenge

## Compétition 2 : reconnaissance de la parole



## Evaluation 2 : Synthèse de la parole

- ▶ Evaluation effectuée par V. Demberg
- ▶ Contexte : utilisation de la morphologie pour améliorer les résultats d'un système de conversion de graphèmes en phonèmes en allemand
- ▶ Résultats décevants : pas d'amélioration des résultats de la conversion
- ▶ Montre que les systèmes de segmentation morphologique non supervisés n'obtiennent pas encore une F-mesure suffisante

# Évaluation 3 : Familles morphologiques

## Familles de référence

- ▶ CELEX pour l'anglais
- ▶ Familles construites manuellement pour l'anglais et le français

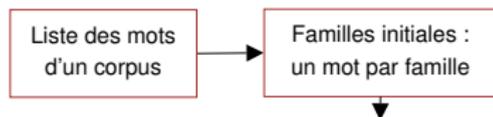
## Mesure d'évaluation

Prise en compte du nombre d'éléments corrects, insérés et supprimés dans une famille morphologique par rapport aux familles de références.

## Résultats obtenus

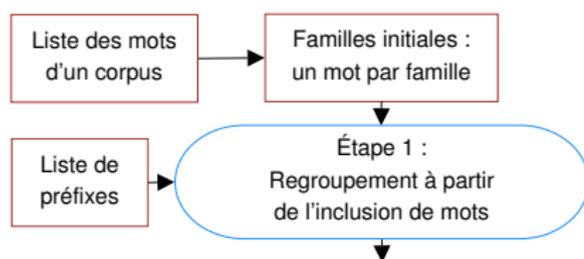
Proches des résultats de MorphoChallenge : F-mesure entre 60 et 70%.

# Analyse morphologique par classification



subvolcaniques  
sub-volcaniques  
post-volcaniques  
volcaniques  
paléo-volcan  
volcan  
paléovolcanique  
volcanique  
subocéanique  
océanique  
sub-océaniques  
océaniques  
océan  
océans

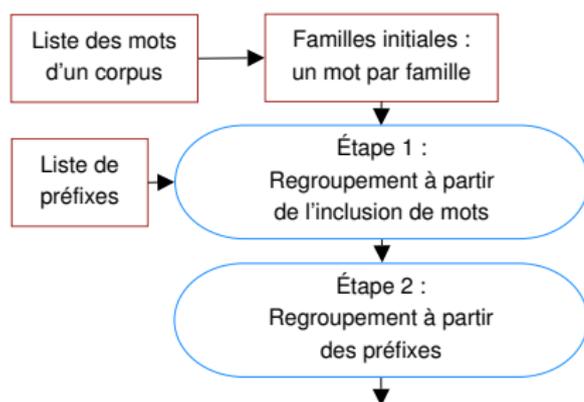
# Analyse morphologique par classification



subvolcaniques sub-volcaniques post-volcaniques volcaniques
paléo-volcan volcan
paléovolcanique volcanique
subocéanique océanique
sub-océaniques océaniques

océan  
océans

# Analyse morphologique par classification



subvolcaniques  
sub-volcaniques  
post-volcaniques  
volcaniques

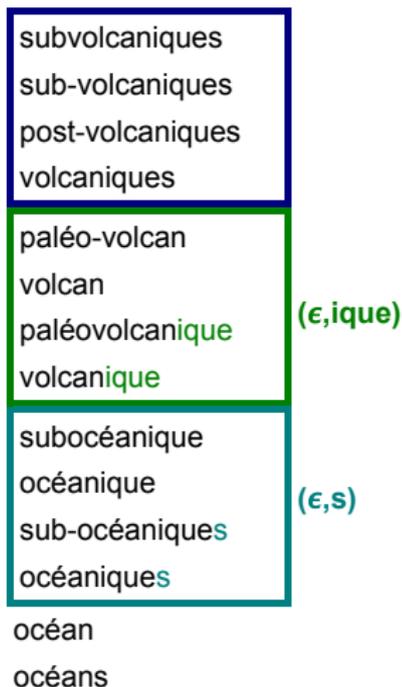
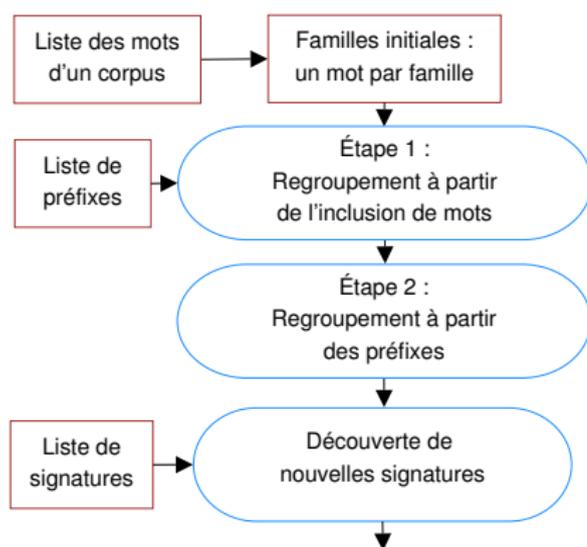
paléo-volcan  
volcan  
paléovolcanique  
volcanique

subocéanique  
océanique  
sub-océaniques  
océaniques

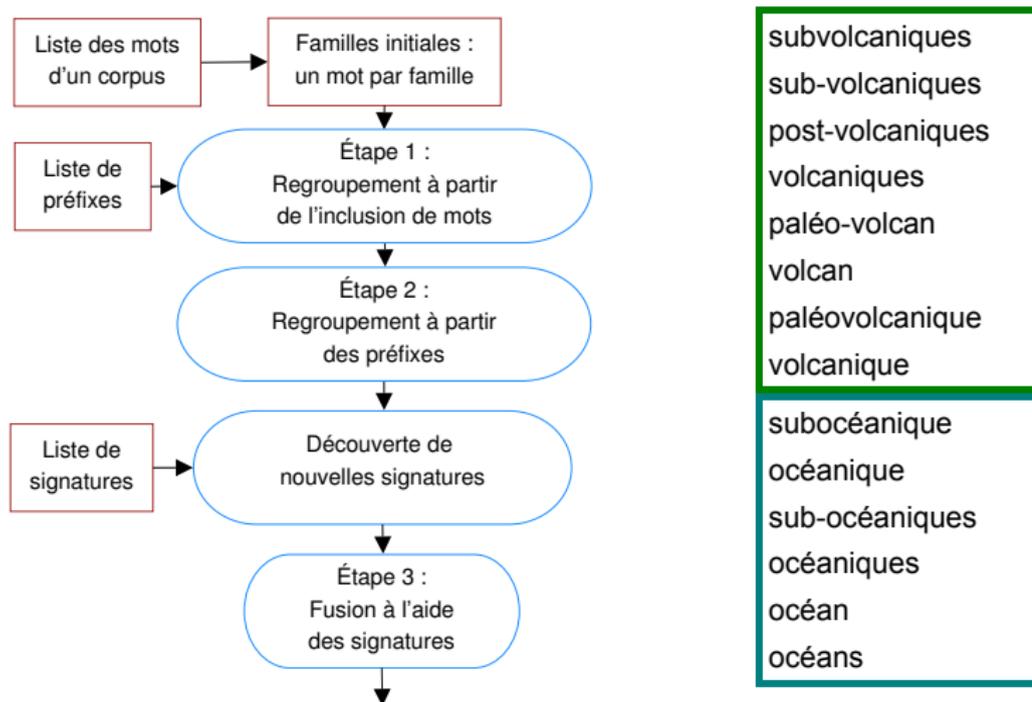
océan

océans

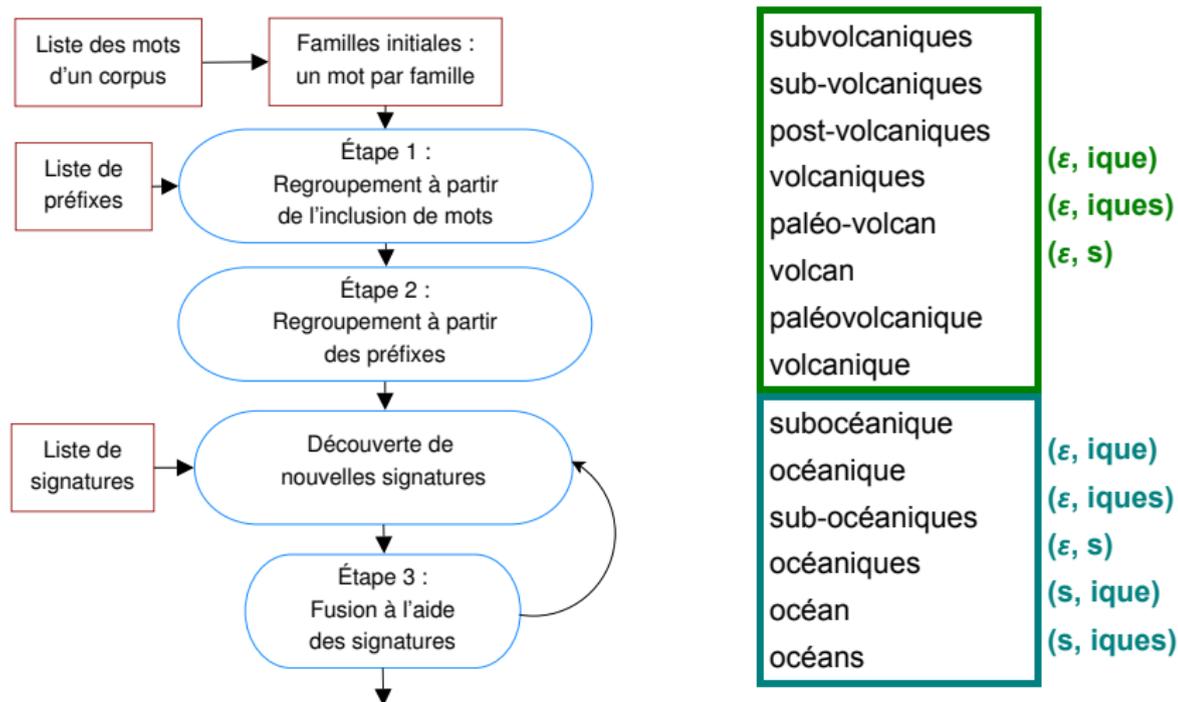
# Analyse morphologique par classification



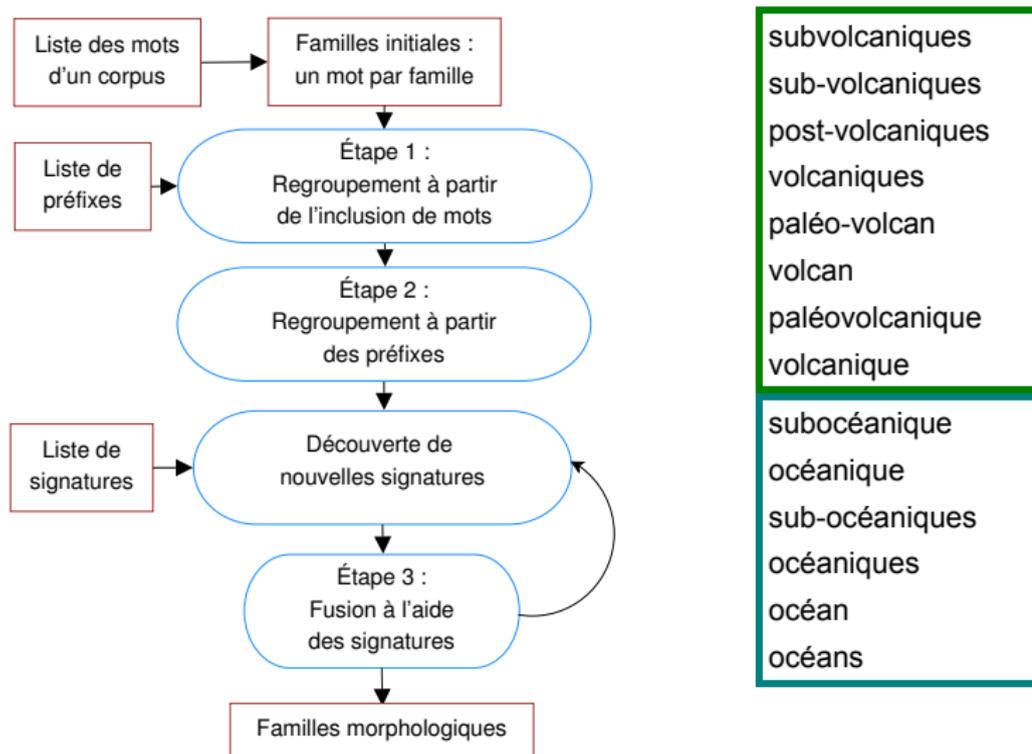
# Analyse morphologique par classification



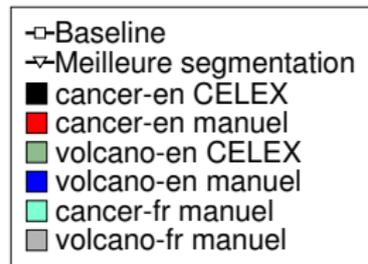
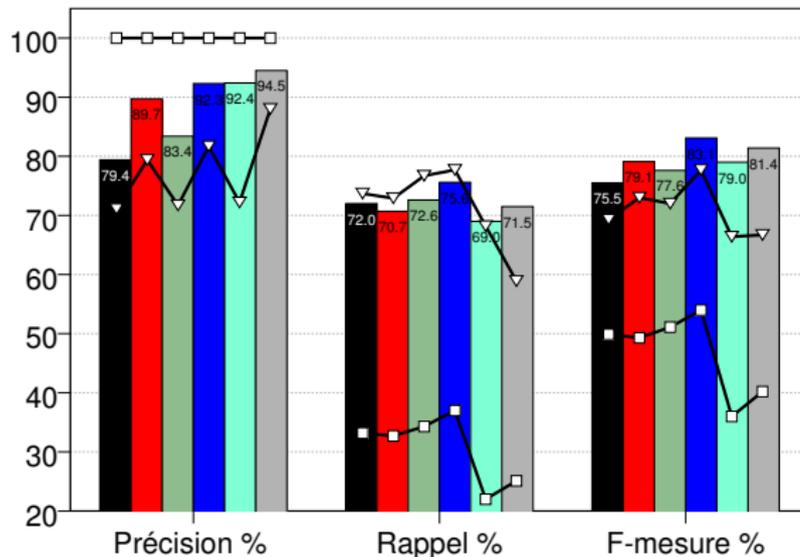
# Analyse morphologique par classification



# Analyse morphologique par classification



# Évaluation : Familles morphologiques



## Améliorations par rapport à l'analyse par segmentation

- ▶ Plus grande précision
- ▶ Doublement des consonnes en fin de radical  
[dimensionnement, dimension, dimensionnées, dimensions]
- ▶ Changements d'accentuation  
[crateres, pseudocratère, intra-cratère, pseudocratères,  
intra-cratérique, cratères, craters, crater, cratérique,  
pseudo-cratères, cratère, cratere, intracratère, intracratérique]

## Perspectives

- ▶ Évaluations complémentaires
- ▶ Classification multiple
- ▶ Déduction d'une segmentation des mots à partir de la classification

Contexte et objectifs

Apprentissage de connaissances morphologiques

- Construction de corpus

- Analyse morphologique par segmentation

- Analyse morphologique par classification

Exploitation des résultats

- Pondération et visualisation de mots clés

- Acquisition de relations sémantiques

Conclusion et perspectives

## Méthode

- ▶ **Mots clés** = mots qui décrivent le mieux le contenu d'un document ou d'un corpus
- ▶ Identification des familles de mots spécifiques au corpus étudié : combinaison d'indices structurels (familles morphologiques) et statistiques (fréquence)

## Mesures de pondération

- ▶ Fréquence de surface : nombre d'occurrences du mot dans le document ou le corpus considéré
- ▶ Fréquence cumulée : somme des fréquences de surface des mots appartenant à une même famille morphologique
- ▶ Comparaison des fréquences (de surface et cumulée) : log du rapport de vraisemblance

## Liste pondérée au format HTML

- ▶ **Liste pondérée** : la taille et la couleur d'un élément dépendent de son poids
- ▶ Carte des thématiques les plus importantes du corpus
- ▶ Représentation des familles par le mot le plus fréquent

above activity and andesite area ash avalanches basalt c  
caldera cinder cloud collapse cone crater crust debris  
deposits dome during earth earthquakes ejected emissions  
**eruption** explosions flank flows formed  
fragments from fumaroles gases geological hazards helens hot island  
kilauea kilometers km lahars lake large lava layer located m  
**magma** material meters mount mountain observed occurred of  
pinatubo plume pumice pyroclastic rim river rock seismic  
slopes small steam summit surface tephra these thick tremor type  
usgs valley vent **volcano** water zone

andesitic	1
andesine	5
andesit	1
andesite	1336
andesite-based	1
andesite-dacite	15
andesites	132
andesitic	581
andesiting	1
basalt-andesite	4
basalt-andesite-dacite	3
dacitic-andesitic	4
	<b>2084</b>

## Mots pondérés par la fréquence de surface

a about above activity also an **and** are area as ash **at** be  
been but **by** caldera can cone crater deposits dome during earthquakes  
**eruption** **eruptions** flow flows for **from** has have high **in**  
into **is** it its km lake large lava m magma may more most mount  
new not **of** **on** one OR pyroclastic rock small some summit surface  
than **that** **the** these they this **to** two up vent volcanic  
volcano volcanoes WAS water were when which with years

## Familles pondérées par la fréquence cumulée

a about activity an **and** are area as ash at basalt be been  
but **by** caldera can cone continued crater deposits dome during  
earthquakes east **eruption** events explosions flows for formed  
**from** has have high **in** into **is** it km lahars lake large lava m  
magma may more most mount not occurred **of on** one OR other  
produced pyroclastic report rock seismic small summit surface than **that**  
**the** these this time **to** vent **volcano** was were west  
which with years

erupt	926
erupted	2658
erupting	659
eruption	14923
eruption-induced	12
eruptions	9237
eruptive	2656
eruptives	6
eruptive-type	6
erupts	390
noneruption-induced	5
noneruptive	21
non-eruptive	23
noneruptively	3
posteruption	62
post-eruption	35
posteruptive	4
post-eruptive	5
preeruption	160
pre-eruption	89
preeruptive	18
pre-eruptive	20
	<b>31947</b>

## Mots pondérés par le log du rapport de vraisemblance

above active **activity** and area **ash** basalt basaltic  
**caldera** cinder cone cones **crater** debris deposits  
dome during earth earthquake **earthquakes** **erupted**  
**eruption** **eruptions** **eruptive** events explosions  
explosive flank **flow flows** formed fragments from gases  
geological helens hot kilauea kilometers **km** lahar lahars lake large  
**lava** layer **m magma** material meters **mount** mountain  
observed occurred **of** pinatubo plume pumice **pyroclastic** river  
rock rocks seismic seismicity small steam summit surface  
tephra these tremor usgs valley vent vents **volcanic**  
**volcano volcanoes** water zone

## Familles pondérées par le log du rapport de vraisemblance

above activity and andesite area ash avalanches basalt c  
caldera cinder cloud collapse cone crater crust debris  
deposits dome during earth earthquakes ejected emissions  
**eruption** explosions flank flows formed  
fragments from fumaroles gases geological hazards helens hot island  
kilauea kilometers km lahars lake large lava layer located m  
magma material meters mount mountain observed occurred of  
pinatubo plume pumice pyroclastic rim river rock seismic  
slopes small steam summit surface tephra these thick tremor type  
usgs valley vent **volcano** water zone

erupt	926
erupted	2658
erupting	659
eruption	14923
eruption-induced	12
eruptions	9237
eruptive	2656
eruptives	6
eruptive-type	6
erupts	390
noneruption-induced	5
noneruptive	21
non-eruptive	23
noneruptively	3
posteruption	62
post-eruption	35
posteruptive	4
post-eruptive	5
preeruption	160
pre-eruption	89
preeruptive	18
pre-eruptive	20
	<b>31947</b>

## Quelles relations sémantiques ?

- ▶ Relations d'inclusion et d'identité :
  - ▶ Synonymie : livre - bouquin
  - ▶ Hyper-/Hyponymie (EST-UN) : chien - animal
  - ▶ Méronymie (PARTIE-DE) : bras - corps
- ▶ Antonymie : chaud - froid
- ▶ Co-hyponymie : chien - chat

## Relations structurelles basées sur les segments morphémiques

### 1. Inclusion

a. Expansion gauche : *lymphedema* – *edema*

[lymph + edema] – [edema]

b. Insertion : *hepatosplenomegaly* – *hepatomegaly*

[hepat + o + splen + o + mega + ly] – [hepat + o + mega + ly]

### 2. Substitution : *osteosarcoma* – *chondrosarcoma*

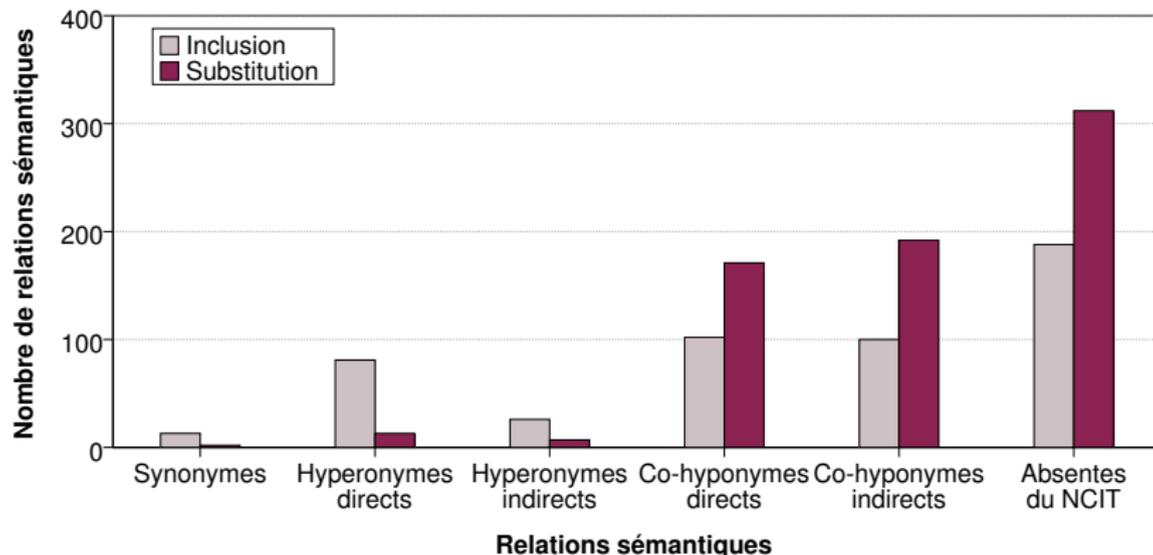
[oste + sarcoma] – [chondro + sarcoma]

## Déduction de liens sémantiques

- ▶ Recherche de paires de mots liés par les relations structurelles précédentes
- ▶ Hypothèses :
  - ▶ Inclusion : hyper-/hyponymie
  - ▶ Substitution : co-hyponymie

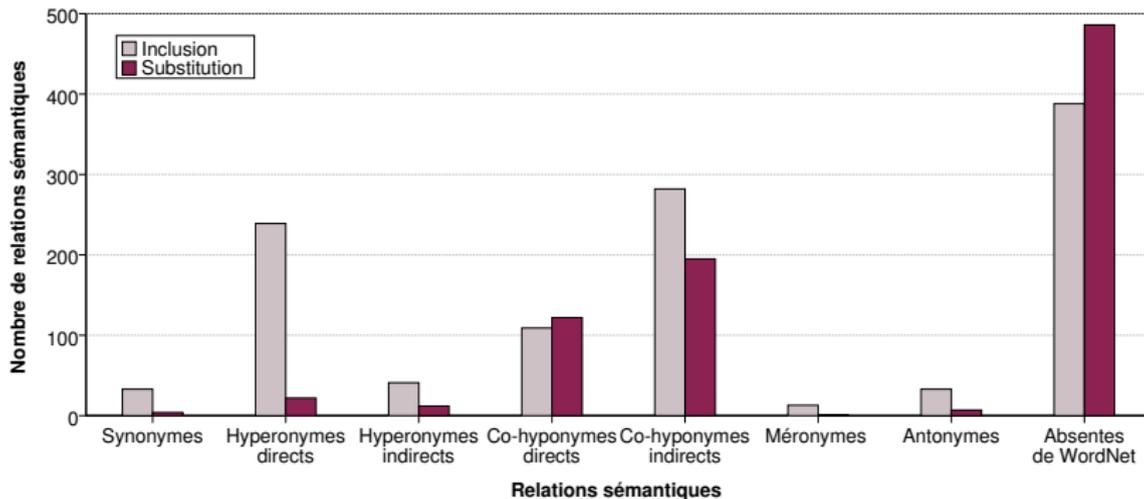
# Relations sémantiques identifiées

Comparaison avec le thésaurus du National Cancer Institute (NCIT)



# Relations sémantiques identifiées

## Comparaison avec WordNet



# Analyse des résultats

## Synonymie : inclusion

*paper, newspaper*

*mistrust, distrust*

## Hyper-/Hyponymie : inclusion

*conductor > semiconductor*

## Co-hyponymie : inclusion **et** substitution

*hypothalamus* et *thalamus* sont co-hyponymes de *neural structure* dans WordNet et co-hyponymes de *Brain\_Part* dans le NCIT

## Méronymie : préfixes

*half-hour, hour*

*midnight, night*

## Antonymie : préfixes

*disagreement, agreement*

*hypertension, hypotension*

Contexte et objectifs

Apprentissage de connaissances morphologiques

- Construction de corpus

- Analyse morphologique par segmentation

- Analyse morphologique par classification

Exploitation des résultats

- Pondération et visualisation de mots clés

- Acquisition de relations sémantiques

Conclusion et perspectives

## Apprentissage non supervisé de connaissances morphologiques

Deux approches différentes :

1. Découpage des mots en segments morphémiques
2. Regroupement des mots dans des familles morphologiques

## Applications

1. Identification et visualisation des mots clés d'un corpus
2. Acquisition de relations sémantiques

- ▶ Travail sur corpus
  - ▶ Corpus construits automatiquement
  - ▶ Données réalistes
- ▶ Langue de spécialité
  - ▶ Deux thématiques : médecine et sciences de la terre
- ▶ Apprentissage et approche statistique
  - ▶ Pas de données externes au corpus
- ▶ Indépendance aux langues
  - ▶ français et l'anglais
  - ▶ + finnois, turc et allemand pour le système d'analyse par segmentation

- ▶ Amélioration des systèmes d'apprentissage de connaissances morphologiques
- ▶ Utilisation des informations contextuelles
- ▶ Évaluation pour d'autres applications et d'autres langues
- ▶ Morpho Challenge 2007

Merci pour votre attention

# Rôle du corpus de référence [1]

Corpus de référence : collection de corpus de l'université de Leipzig (presse)

above activity and andesite area ash avalanches basalt c  
caldera cinder cloud collapse cone crater crust debris  
deposits dome during earth earthquakes ejected emissions  
**eruption** explosions flank flows formed  
fragments from fumaroles gases geological hazards helens hot island  
kilauea kilometers km lahars lake large lava layer located  
m magma material meters mount mountain observed  
occurred of pinatubo plume pumice pyroclastic rim  
river rock seismic slopes small steam summit surface tephra  
these thick tremor type usgs valley vent **volcano**  
water zone

# Rôle du corpus de référence [2]

Corpus de référence : liste de mots anglais de Morpho Challenge (projet Gutenberg, corpus Gigaword et Brown)

above activity andesite area ash avalanches basalt caldera  
cinder collapse composition cone crater crust data debris  
deposits dome during earthquakes east ejected emissions  
**eruption** explosions flank flows formed  
fragments from fumaroles gas gases geological hazards helens june  
kilauea kilometers km lahars lake large lava layer level  
located m magma material meters mount mountain occurred  
peak photo pinatubo plume pumice pyroclastic report  
rhyolite rim river rock seismic slopes steam summit surface  
survey tephra tremor type usgs valley vent **volcano**  
west zone

# Rôle du corpus de référence [3]

## Corpus de référence : corpus médical

above activity andesite ash at avalanches basalt caldera  
cinder cloud collapse composition cone continued crater  
debris deposits dome during earth earthquakes east  
ejected **eruption** explosions fall feet flank  
flows formed fragments from fumaroles gas gases geological  
hazards helens island june kilauea kilometers km lahars lake  
large lava layer m magma meters miles mount mountain  
near north on pinatubo plume pumice pyroclastic river  
rock seismic slopes south steam summit surface survey  
tephra the tremor usgs valley vent **volcano** water  
west zone

## Même niveau dans la hiérarchie

- ▶ Antonymie  
**in**activity, activity  
**non**smoker, smoker
- ▶ Unités de mesures :  
**kilo**volt, volt  
**table**spoon, spoon
- ▶ Position :  
**hypo**thalamus, thalamus  
**para**thyroid, thyroid

L'absence de segment morphémique est porteuse de sens