

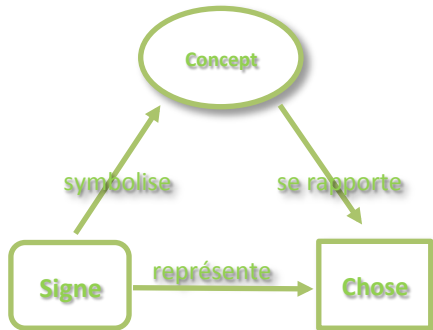
La 4^{ème} édition des journées
« Big Data Mining and Visualization »

Approche **sémantique** pour la **recommandation de documents** textuels dans un contexte **Big Data** appliquée à la veille économique.



Christophe CRUZ | 19 juin 2015 | Lyon

Modélisation et représentation des connaissances (MRC)



Web sémantique

Système centré utilisateur (SCU)



Système de recommandation

Big Data Value (BDV)



Véracité et valeur

Approche Numérique

- Taille, Dimension, poids
- Couleur, Type de matériaux
- Valeur financière

Approche Symbolique

- Espaces et usages
- Organisation architecturale
- Confort, dessert et facilités d'accès



Evaluation quantitative
entre éléments spécifiques

Evaluation qualitative
Sémantique sur un environnement



SEMIOTIQUE

Point de vue sémiotique :

Sémiotique: étudie le processus de signification c'est-à-dire la production, la codification et la communication de signes.

Syntaxe: un ensemble de signes

Sémantique: relation entre les signes et des entités du monde réel

Pragmatique: quels signes sont utilisés dans quel but ?

Social : qui utilise quels signes ?

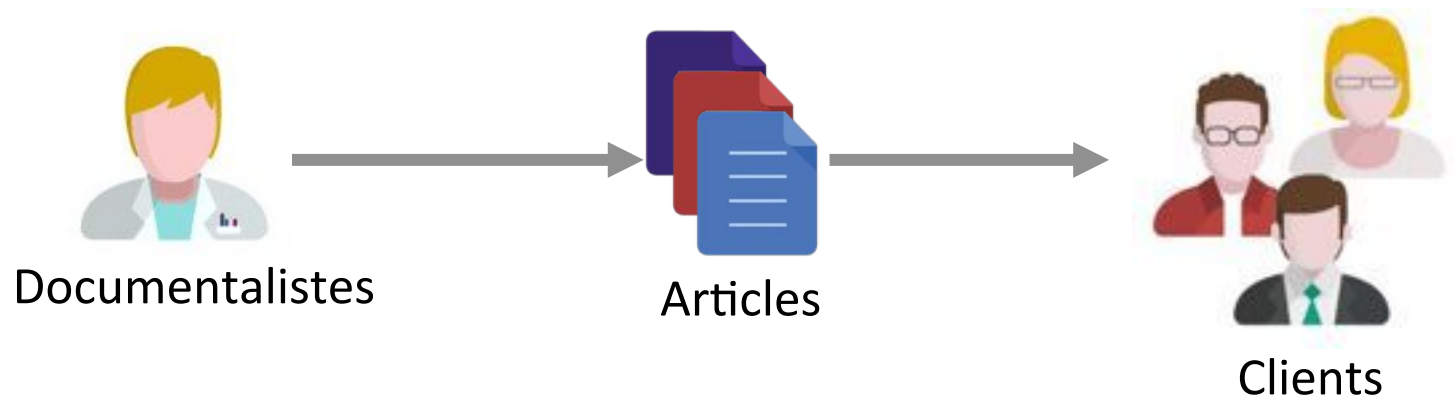


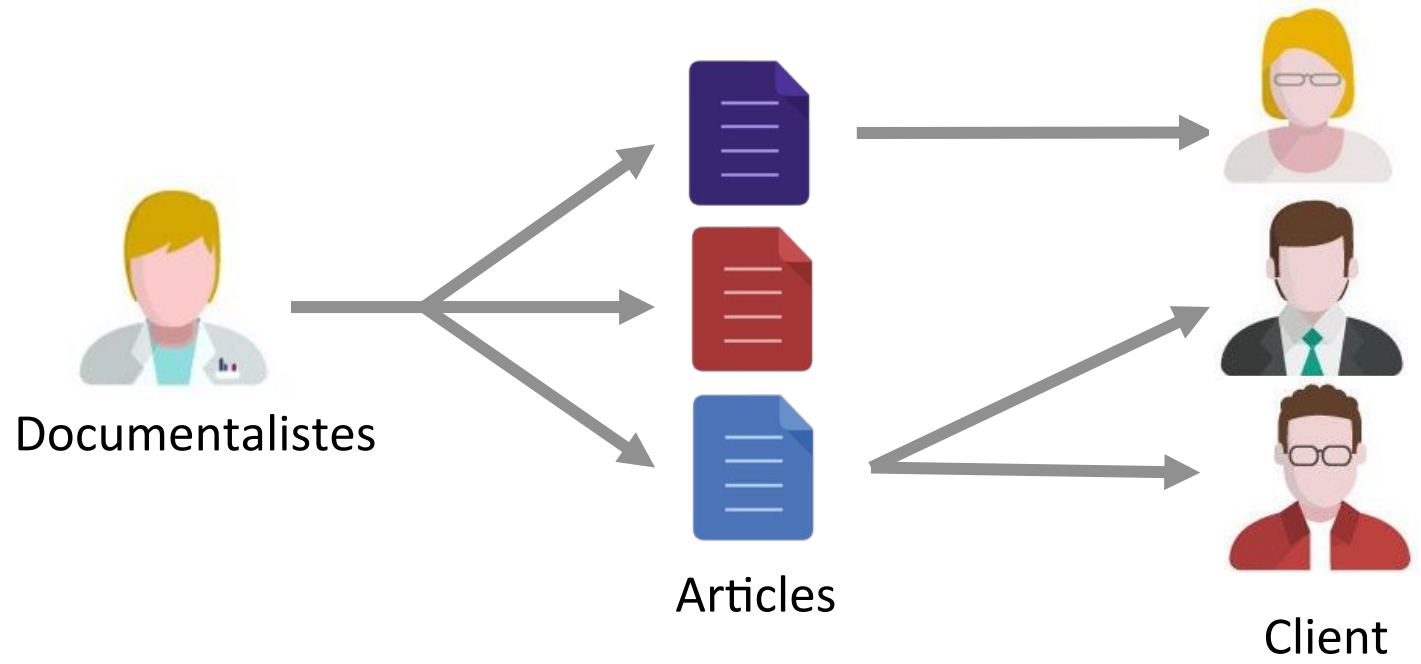
1. Le contexte industriel
2. Le transfert de technologie
3. La quête de la sémantique
4. La problématique Big Data
5. La solution First Eco Pro'fil
6. Conclusion

Le contexte industriel ...

La solution First Pro'Fil



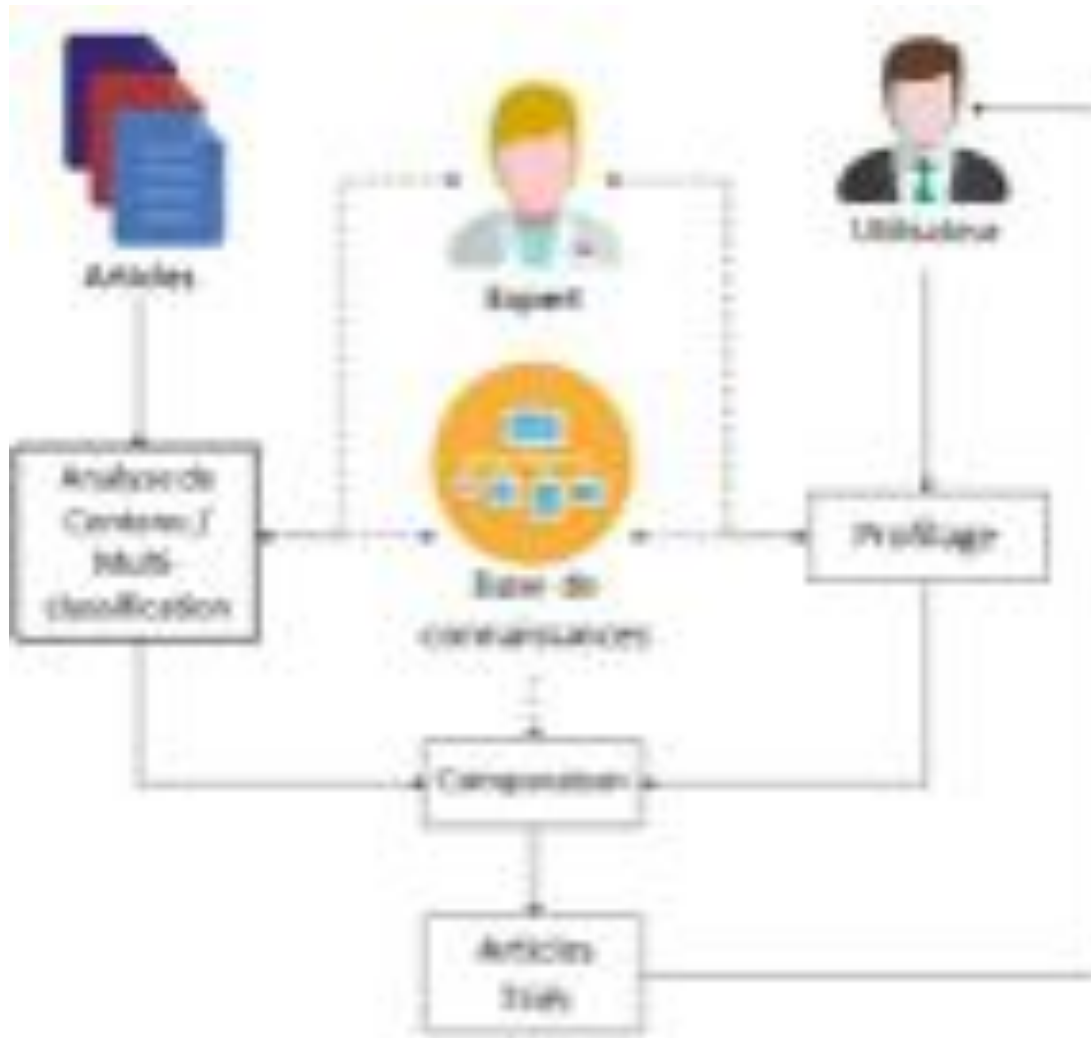






L'UNIQUE PORTAIL DE VEILLE ÉCONOMIQUE ET COMMERCIALE À VOTRE IMAGE

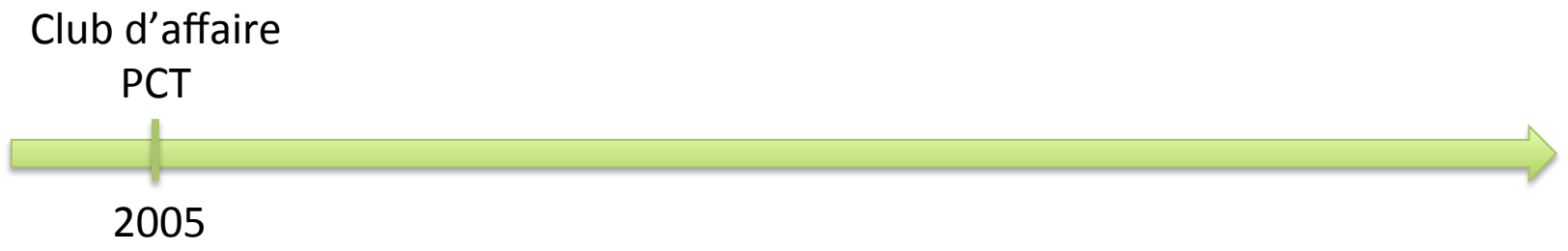


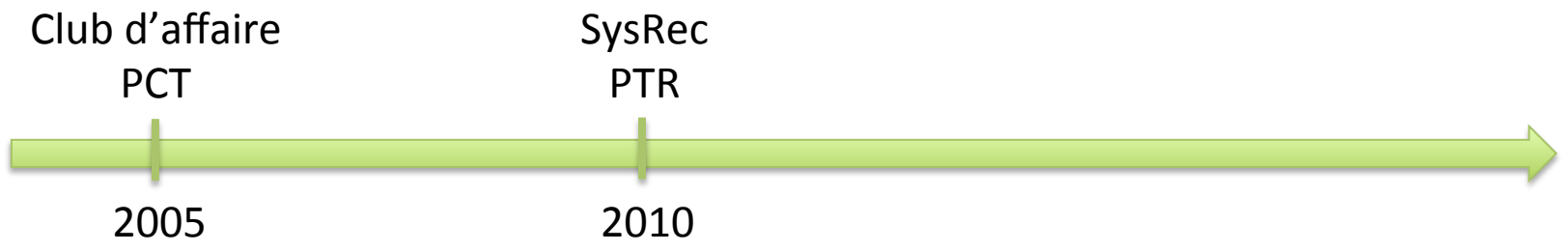


Transfert de technologie ...

Actualis SARL



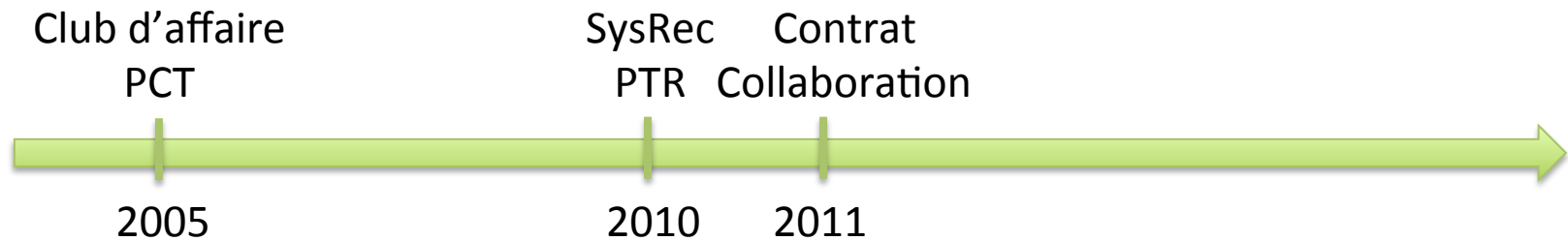




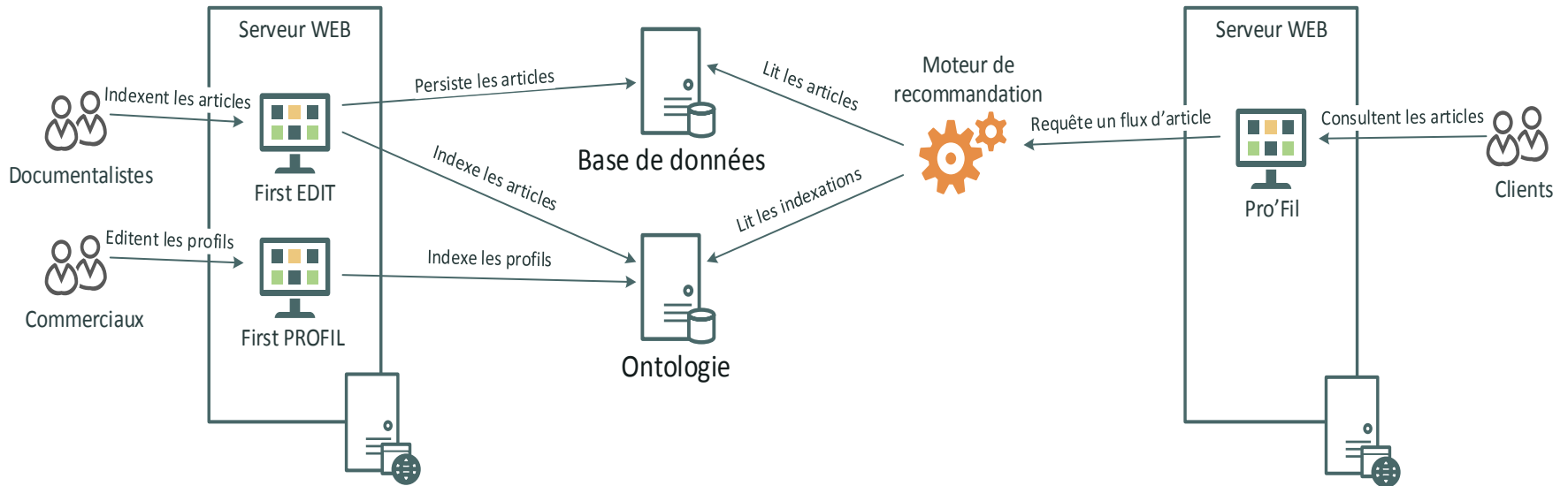


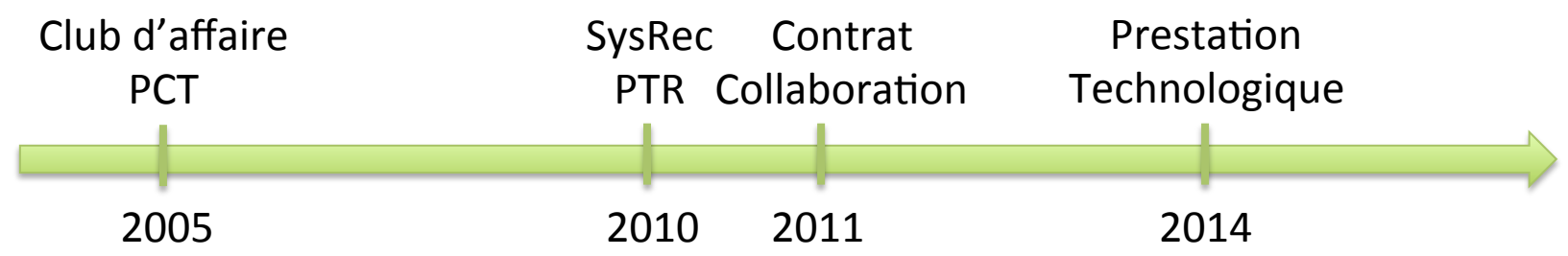
Rupture de politique - R&D

- Équipe de développement
- Changement du « process » métier
- Montée en compétence
- Homogénéisation des compétences et savoirs



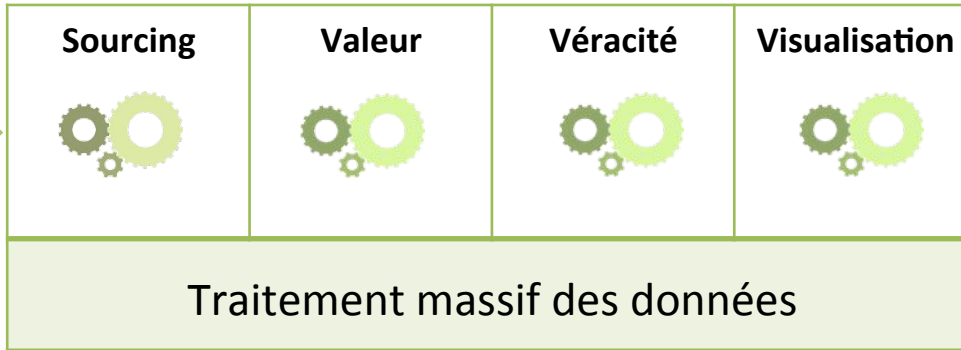
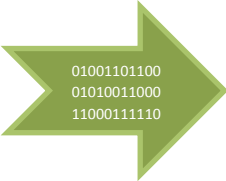
Architecture complète





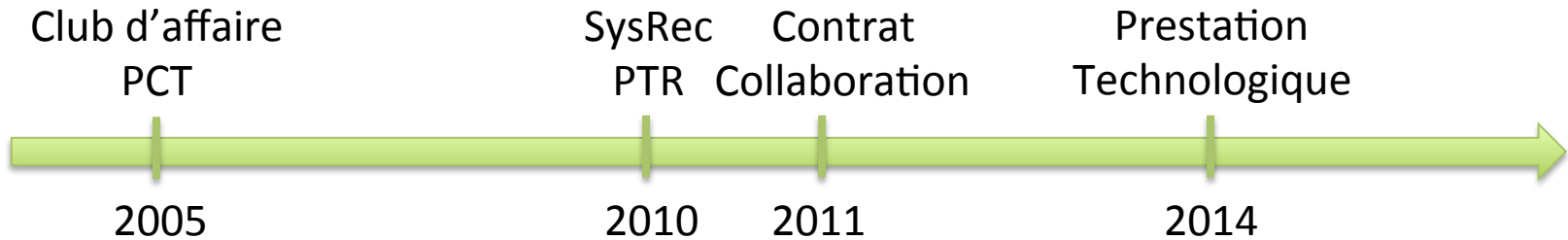


Web et Flux



First ECO
| Profil

Plateforme
Actualis



Transfert de technologie ...

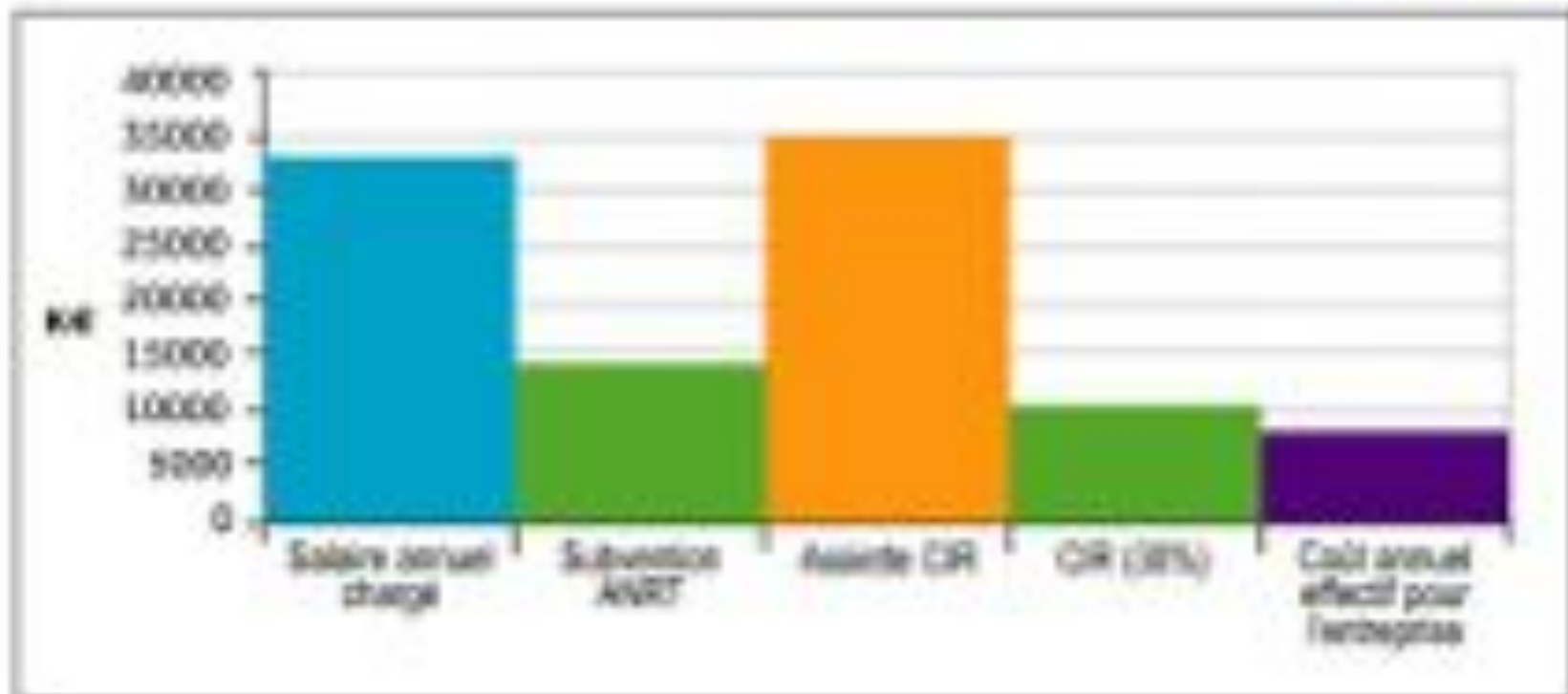
Parlons argent !



Illustration à l'aide d'un exemple simplifié

Un exemple permet de concrétiser l'impact pour une entreprise de l'embauche d'un doctorant dans le cadre d'une CIFRE.

Pour un salaire chargé du doctorant de 32 878 €, le coût effectif pour l'entreprise est de **8 103 €**.



Coûts	Année 1	Année 2	Année 3	Informations
Doctorant	33 000,00	33 000,00	33 000,00	
CRI	-24 700,00	-24 700,00	-24 700,00	dont 14k sub ANRT + CRI

Coûts	Année 1	Année 2	Année 3	Informations
Doctorant	33 000,00	33 000,00	33 000,00	
CRI	-24 700,00	-24 700,00	-24 700,00	dont 14k sub ANRT + CRI
Contrat de collab.	42 000,00	42 000,00	42 000,00	126 000 pour labo
Taux 2,0	84 000,00	84 000,00	84 000,00	doublement somme éligible
CRI	-25 200,00	-25 200,00	-25 200,00	30%/an

Coûts	Année 1	Année 2	Année 3	Informations
Doctorant	33 000,00	33 000,00	33 000,00	
CRI	-24 700,00	-24 700,00	-24 700,00	dont 14k sub ANRT + CRI
Contrat de collab.	42 000,00	42 000,00	42 000,00	126 000 pour labo
Taux 2,0	84 000,00	84 000,00	84 000,00	doublement somme éligible
CRI	-25 200,00	-25 200,00	-25 200,00	30%/an
Dépenses (2 ing.)	200 000,00	200 000,00	200 000,00	incompressible
Taux 1,5	300 000,00	300 000,00	300 000,00	éligible
CRI	-90 000,00	-90 000,00	-90 000,00	30%/an

Coûts	Année 1	Année 2	Année 3	Informations
Doctorant	33 000,00	33 000,00	33 000,00	
CRI	-24 700,00	-24 700,00	-24 700,00	dont 14k sub ANRT + CRI
Contrat de collab.	42 000,00	42 000,00	42 000,00	126 000 pour labo
Taux 2,0	84 000,00	84 000,00	84 000,00	doublement somme éligible
CRI	-25 200,00	-25 200,00	-25 200,00	30%/an
Dépenses (2 ing.)	200 000,00	200 000,00	200 000,00	incompressible
Taux 1,5	300 000,00	300 000,00	300 000,00	éligible
CRI	-90 000,00	-90 000,00	-90 000,00	30%/an
Coût total	275 000,00	275 000,00	275 000,00	825 000,00
CRI total	-139 900,00	-139 900,00	-139 900,00	-419 700,00
Coût total - CRI	135 100,00	135 100,00	135 100,00	0,50
Coût total Rech.	50 300,00	50 300,00	50 300,00	

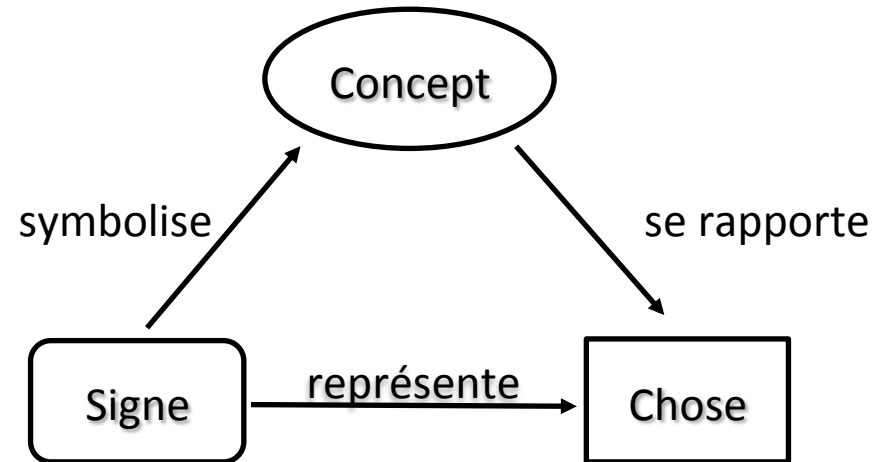
La quête de la sémantique ...



Complexe



Simple



Ogden et Richards (1923)

Modèles abstraits



Raisonnement

Du Web-documentaire ...



au Web des données

Echange et partage



(1906, 1978)

Kurt Gödel

“Systèmes logiques”

“Systèmes de décisions”

Alan Turing



(1912, 1954)

Décidabilité

$\exists \text{Algo} \rightarrow V \cup F$

Calculabilité

$\exists \text{Algo}$

pour $f(x)$

en temps fini



Phase 0 (1965-1980)

Pre-DL - introduction réseaux sémantiques et frames.

Critiques liées au manque de **sémantique formelle**.

KL-One: premier système de logique de description

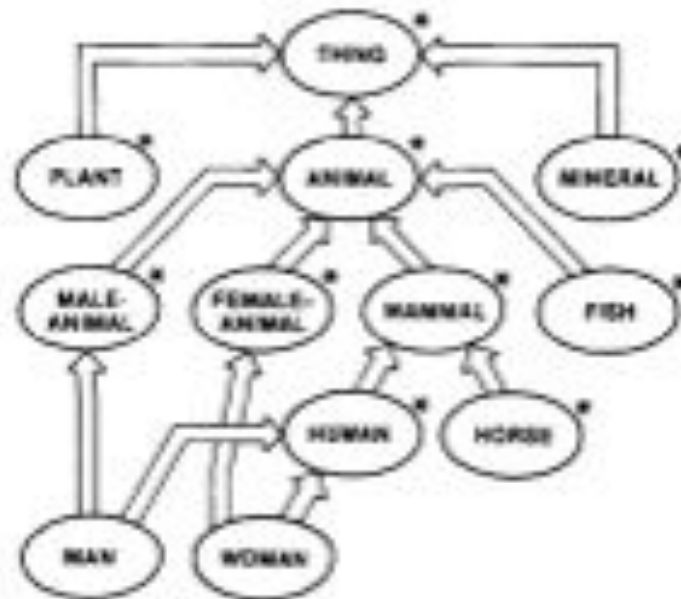


Figure 1. A simple KL-ONE network of Generic Concepts.



Phase 1 (1980-1990)

Algorithmes de subsomption structurale

KL-One, K-Rep, Krypton, Back, Loom

Très efficaces, mais **non complets**,
voire même **indécidables**

sauf pour des fragments très pauvres

English	FOL	DL
Dog with a Spot (DWS)	$DWS(x) \Leftrightarrow$ $Dog(x) \wedge (\exists y. has(x,y)$ $\quad \wedge Spot(y))$	$DWS \Leftrightarrow$ $Dog \sqcap \exists has. Spot$
Large Dog with a Dark Spot (LDWDS)	$LDWDS(x) \Leftrightarrow$ $(Dog(x) \wedge Large(x)) \wedge$ $(\exists y. has(x,y)$ $\quad \wedge (Spot(y) \wedge Dark(y)))$	$LDWDS \Leftrightarrow$ $Dog \sqcap Large \sqcap$ $\exists has. (Spot \sqcap Dark)$



Phase 2 (1990-1995)

Introduction d'algorithmes fondés sur les **tableaux**.

Permet de décider de la **cohérence**

Premiers systèmes utilisant ces méthodes: Kris et Krack.
 Implémentations efficaces, même si la complexité dans le pire des cas n'est plus polynomiale (**TBox/ABox**)

<i>TBox</i>	<i>ABox</i>
$Femelle \sqsubseteq \top \sqcap \neg M\grave{a}le$	$Humain(Anne)$
$M\grave{a}le \sqsubseteq \top \sqcap \neg Femelle$	$Femelle(Anne)$
$Animal \equiv M\grave{a}le \sqcup Femelle$	$Femme(Sophie)$
$Humain \sqsubseteq Animal$	$Humain(Robert)$
$Femme \equiv Humain \sqcap Femelle$	$\neg Femelle(Robert)$
$Homme \equiv Humain \sqcap \neg Femelle$	$Homme(David)$
$M\grave{e}re \equiv Femme \sqcap \exists relationParentEnfant$	$relationParentEnfant(Sophie, Anne)$
$P\grave{e}re \equiv Homme \sqcap \exists relationParentEnfant$	$relationParentEnfant(Robert, David)$
$M\grave{e}reSansFille \equiv M\grave{e}re \sqcap$ $\forall relationParentEnfant. \neg Femelle$	
$relationParentEnfant \sqsubseteq \top_R$	



Thomas Gruber

“A Translation Approach to Portable Ontology”
1993

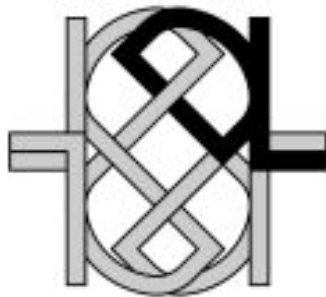
(1959,)

Concepts, propriétés
relations, fonctions,
contraintes, axiomes,
sont définis

Sans ambiguïté

Une ontologie est une spécification explicite d'une
conceptualisation

Modèle abstrait et vue simplifiée
d'un phénomène réel que nous
souhaitons représenter



Phase 3 (1995-2000)

LD très expressives basées sur les tableaux

Exploration des liens avec les logiques modales

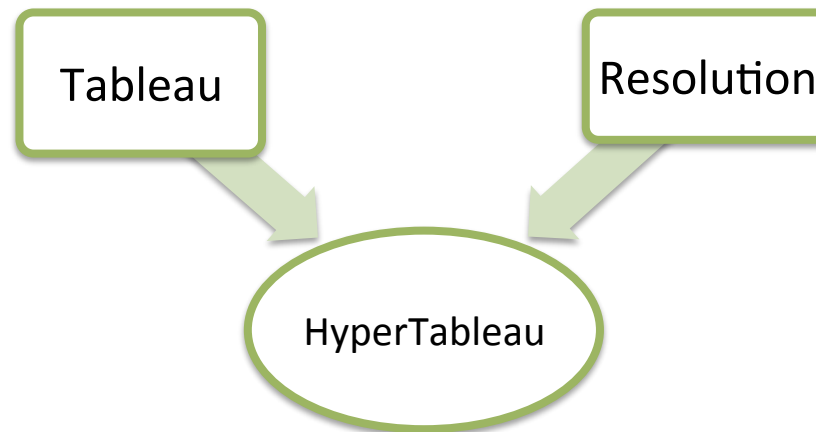
$\mathcal{ALC} ::= \perp \mid \top \mid A \mid \neg C \mid C \sqcap D \mid C \sqcup D \mid \exists R.C \mid \forall R.C$

<p>Concept constructors:</p> <ul style="list-style-type: none"> ⊖ F - functionality²: $\{x \mid xR\}$ ⊖ N - (unqualified) number restrictions: $\{2 \geq R\}$, $\{5 \leq R\}$ ⊖ Q - qualified number restrictions: $\{2 \geq R.C\}$, $\{5 \leq R.C\}$ ⊖ D - nominals: $\{A\}$ or $\{A_1, \dots, A_n\}$ ("one-of") ⊖ μ - least fixpoint operator: $\mu X.C$ <p><i>Footnote: complex roles² in number restrictions³</i></p>	<p>Role constructors:</p> <ul style="list-style-type: none"> ⊖ I - role inverse: R^{-1} ⊖ ∩ - role intersection²: $R \cap S$ ⊖ ∪ - role union: $R \cup S$ ⊖ ¬ - role complement: $\neg R$ (w/ \neg) ⊖ ∘ - role chain (composition): $R \circ S$ ⊖ * - reflexive-transitive closure⁴: R^* ⊖ id - concept identity: $\text{id}(C)$
<p>TBox (concept axioms):</p> <ul style="list-style-type: none"> • empty TBox • acyclic TBox ($A \sqsubseteq C$, A is a concept name; no cycles) • general TBox ($C \sqsubseteq D$, for arbitrary concepts C and D) 	<p>RBox (role axioms):</p> <ul style="list-style-type: none"> ⊖ S - role transitivity: $\text{Tr}(R)$ ⊖ H - role hierarchy: $R \sqsubseteq S$ ⊖ R - complex role inclusions: $R \circ S \sqsubseteq R$, $R \circ S \sqsubseteq S$ ⊖ x - some additional features (click to see them)

OWL-Lite, OWL-DL, OWL 1.1

Reasoner You have selected a Description Logic: \mathcal{ALC}

Systems	Algorithm	Sound	Complete	Rule Support	Expressivity
Pellet	Tableau	Yes	Yes	Yes (SWRL)	SROIQ (D)
FaCT++	Tableau	Yes	Yes	No	SROIQ (D)
Racer	Tableau	Yes	Yes	Yes (SWRL)	SHIQ
Kaon2	Resolution	Yes	Yes	Yes (SWRL)	SHIQ
Hermit	HyperTableau	Yes	Yes	Yes (SWRL)	SROIQ (D)



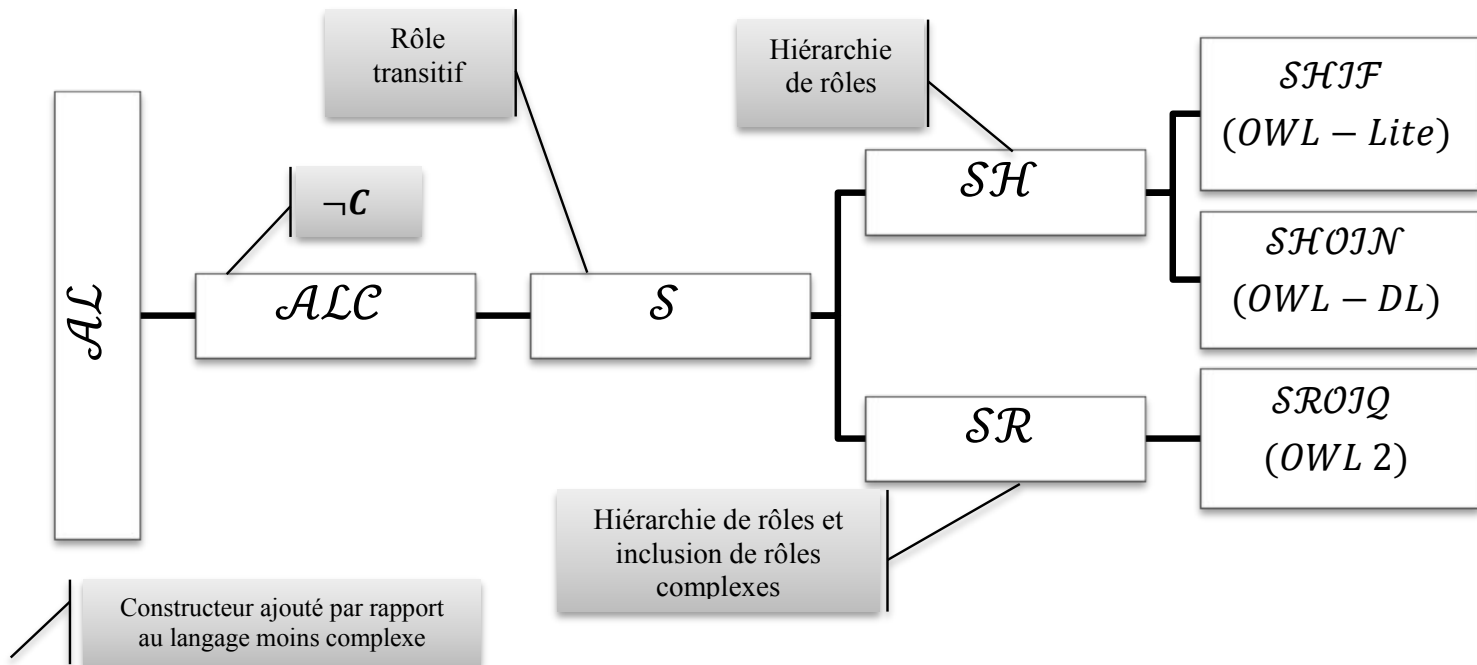


Phase 4 (2000-2012)

Web sémantique

Systèmes d'information, ...

Utilisation DL moins expressives...





« Le web sémantique ne fonctionnera jamais ! »

[James Hendler « Why the semantic Web will never work », ESWC'11](#)

Numérique

Théologie

Botanique

Physique

Politique



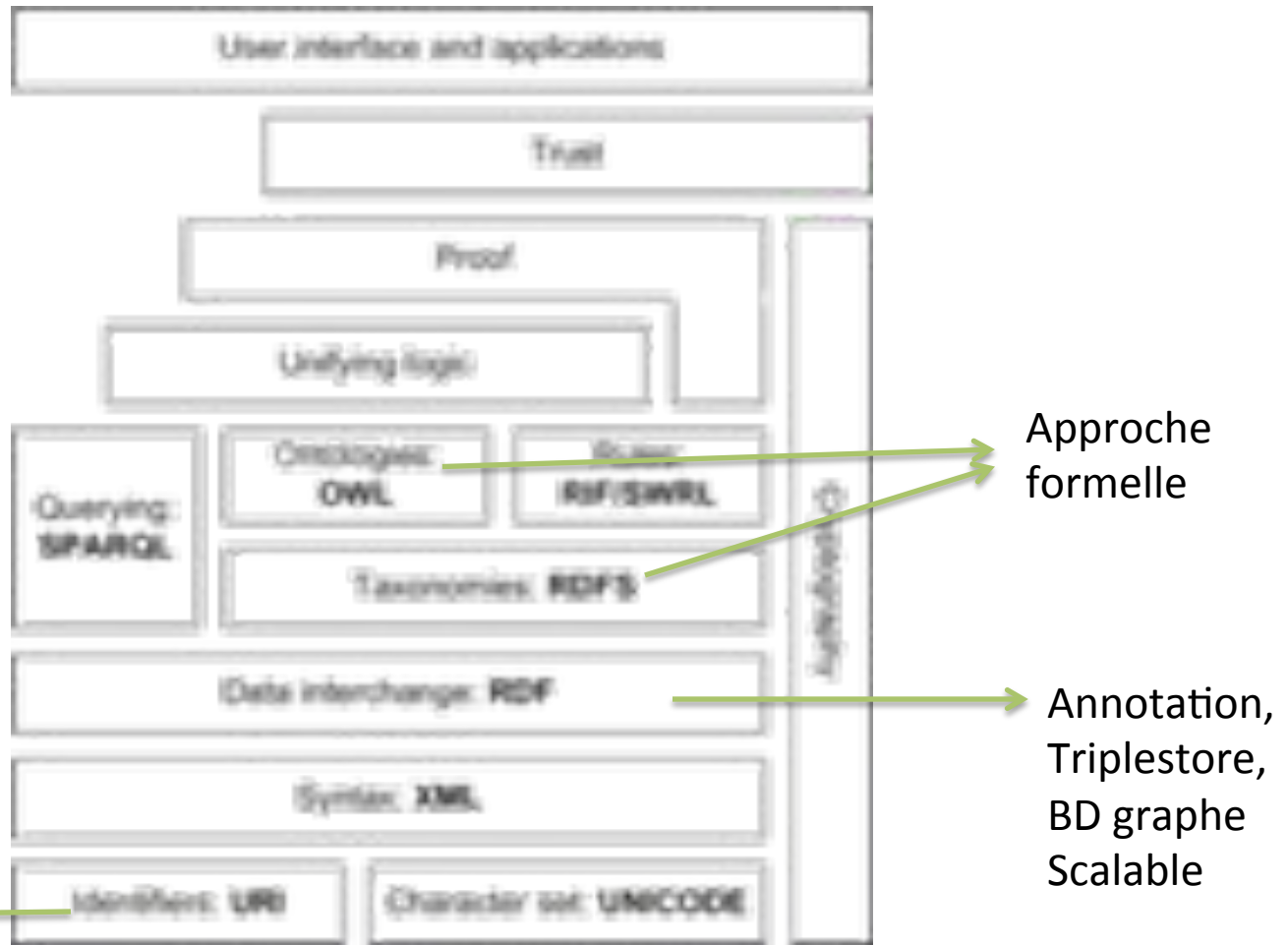
Sciences de la
nutrition

Agriculture

Culture

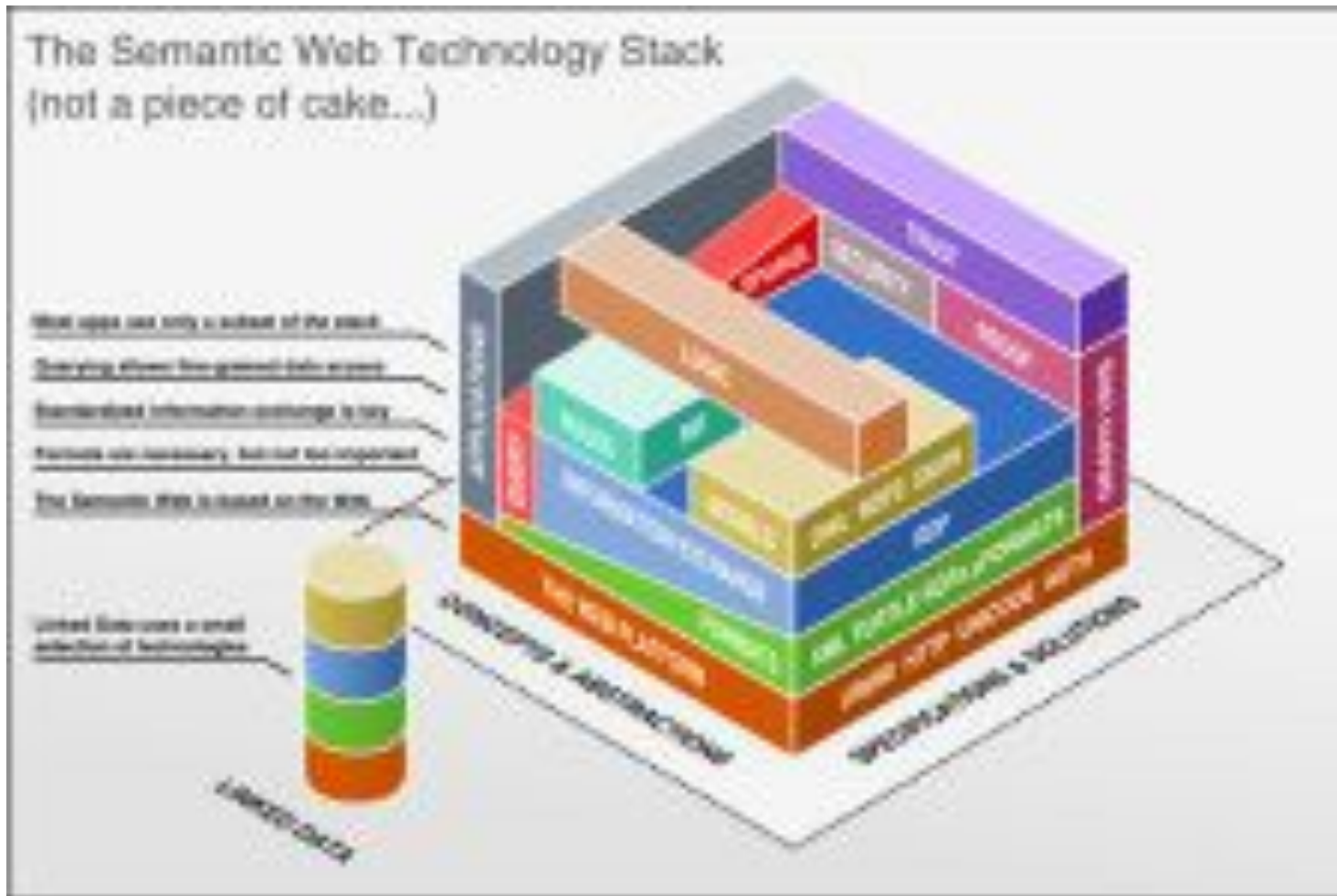
Semantic Error 409 – Ontology Not Found

You've encountered a " Ontology Not Found" error while trying to access a semantic term grounded on the [University of Burgundy Laboratory LE2I](#) Cheksem web server.



14 ans après ...

« La pile technologique du Web sémantique »



14 ans après ...

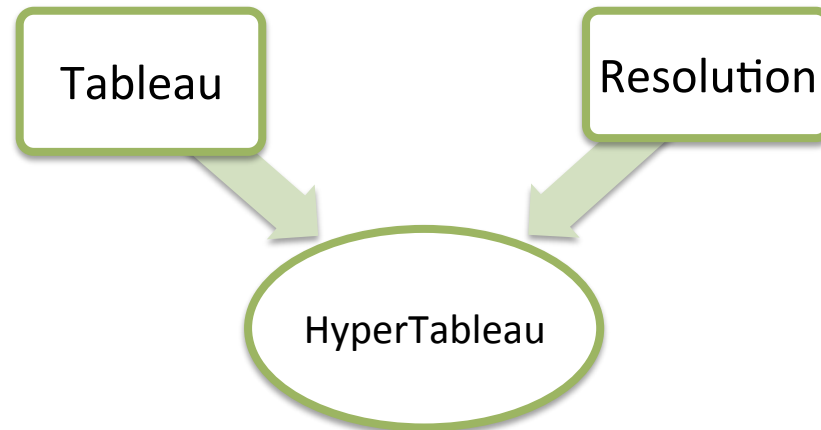
« La pile technologique du Web sémantique »



“I’m sorry Dave,
I’m afraid I can’t do that.”

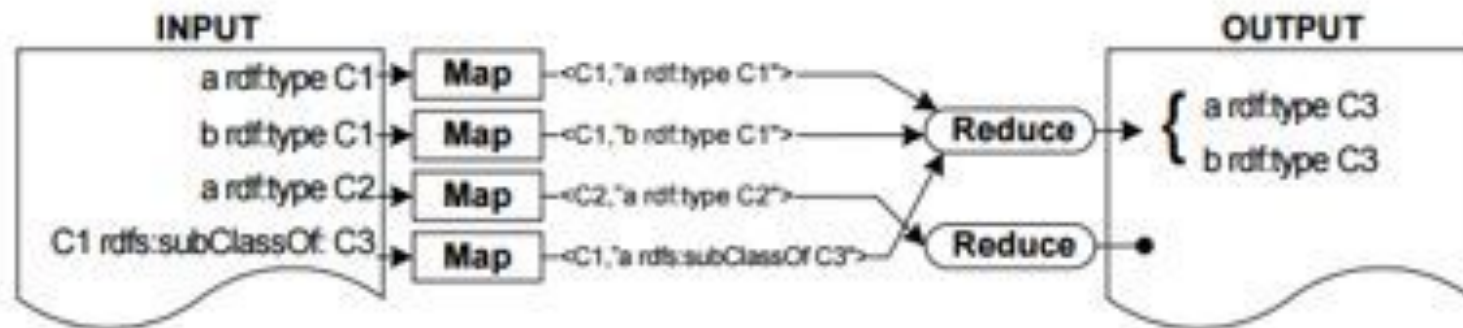
Phase 5 (2012-20..)

Application à l’échelle du web (sémantique)
« Web Reasoning »



Passage à l'échelle difficile ...

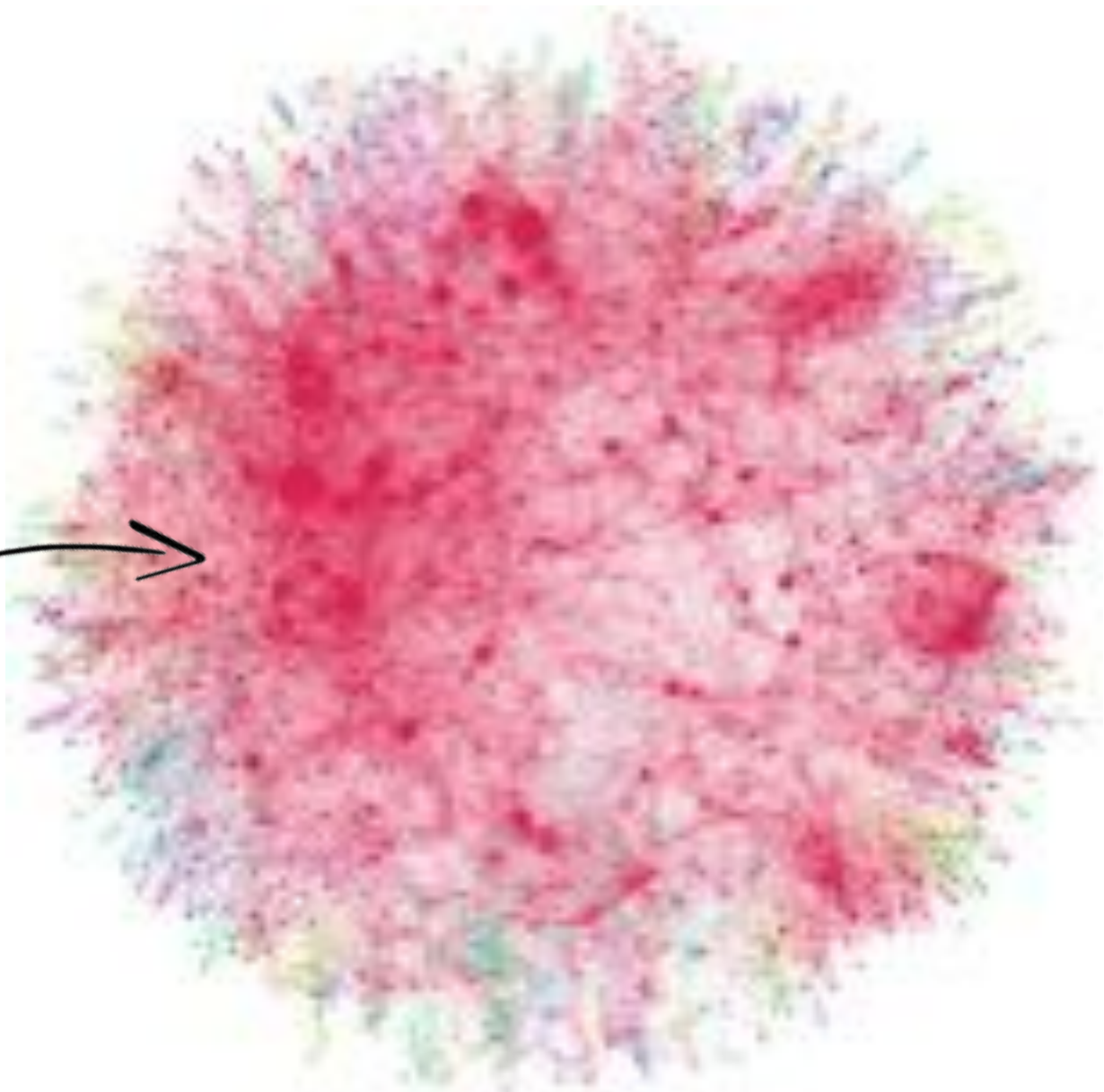
WebPIE - Raisonnement basé sur les règles ...



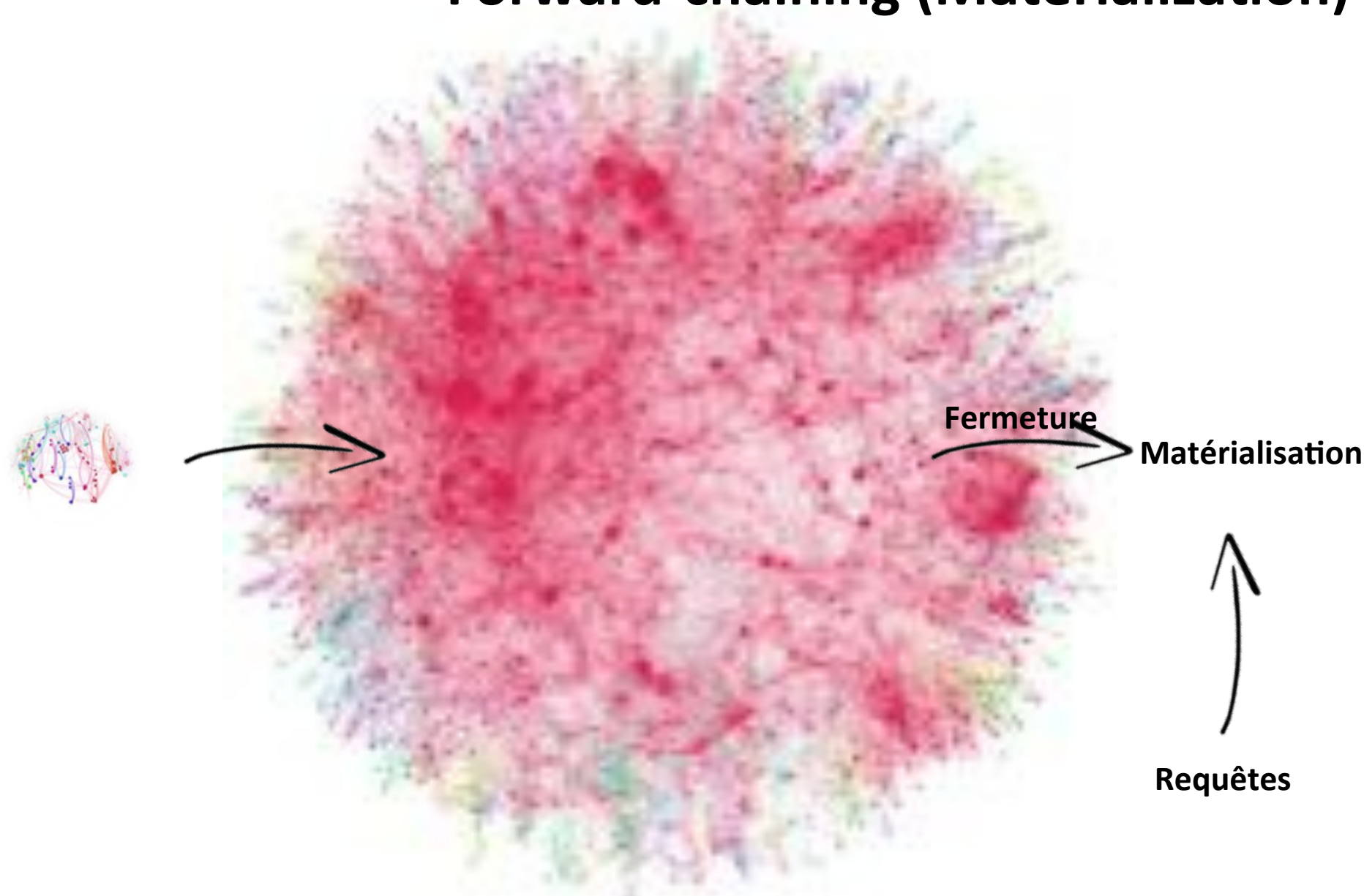
et Map Reduce

Urbani, J., Kotoulas, S., Maassen, J., Van Harmelen, F., & Bal, H. (2012). WebPIE: A Web-scale parallel inference engine using MapReduce. *Web Semantics: Science, Services and Agents on the World Wide Web*.

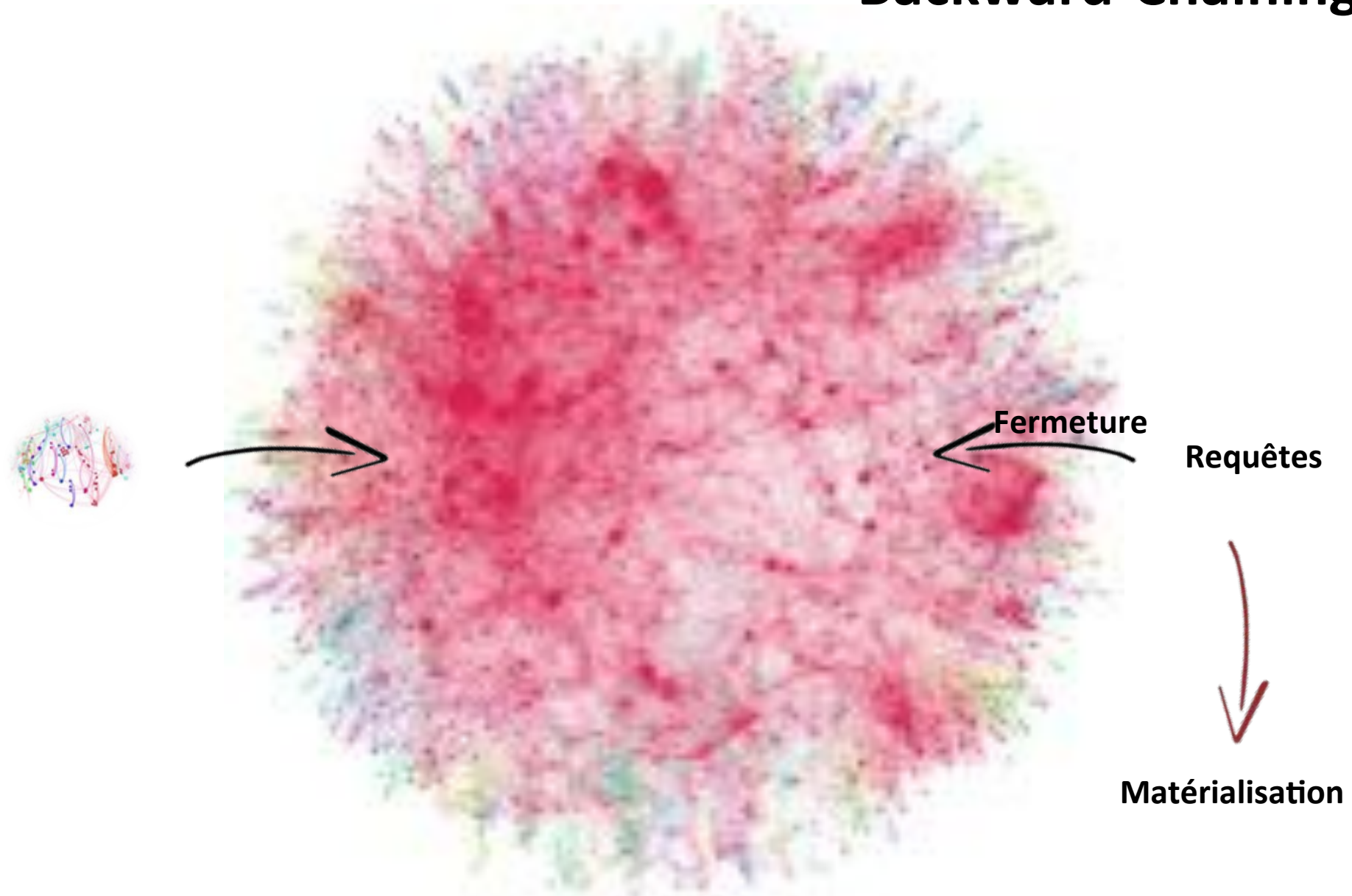
Urbani, J., Kotoulas, S., Oren, E., & Harmelen, F. (2009). **Scalable Distributed Reasoning** Using MapReduce. In - *ISWC 2009 SE - 40* (Vol. 5823, pp. 634–649).



Forward-chaining (Materialization)



Backward-Chaining



Règles pré-matérialisées

QueryPIE

Matérialisation
Backward-chaining
OWL-RL

	rules	<pre>118 owl:reducedProperty rdfs:Property 119 owl:reducedClass rdfs:Class 120 owl:domain owl:Thing 121 owl:rdfs:Property rdfs:Property</pre>
	rule	<pre>122 owl:type owl:Property owl:Property owl:Property owl:Property owl:Property 123 owl:subClassOf owl:Class owl:Class 124 owl:type owl:Class owl:Property owl:Property owl:Property owl:Property 125 owl:subPropertyOf owl:Property 126 owl:type owl:Class owl:Class 127 owl:type owl:SymmetricProperty owl:Property 128 owl:type owl:Property owl:Property 129 owl:equivalentClass owl:Class 130 owl:subPropertyOf owl:Property 131 owl:subPropertyOf owl:Property 132 owl:subPropertyOf owl:Property 133 owl:subPropertyOf owl:Property 134 owl:subPropertyOf owl:Property 135 owl:subPropertyOf owl:Property 136 owl:subPropertyOf owl:Property 137 owl:subPropertyOf owl:Property 138 owl:subPropertyOf owl:Property 139 owl:subPropertyOf owl:Property 140 owl:subPropertyOf owl:Property 141 owl:subPropertyOf owl:Property 142 owl:subPropertyOf owl:Property 143 owl:subPropertyOf owl:Property 144 owl:subPropertyOf owl:Property 145 owl:subPropertyOf owl:Property 146 owl:subPropertyOf owl:Property 147 owl:subPropertyOf owl:Property 148 owl:subPropertyOf owl:Property</pre>

Conférence ESWC'15

Prototype sur OWL 2 EL

Distributed and Scalable OWL EL Reasoning

Nguyen Huuhoang¹, Pascal Struder², ProfRadu Mamon², and Freddy Léves²

¹ Wright State University, OH, USA.

{huhoang.n, pascal.struder, prof@radu.mamon}@wright.edu

² SmartCityX Technology Centre, IBM Research, Dublin, Ireland.
freddy.levess@ie.ibm.com

Abstract. OWL 2 EL is one of the tractable profiles of the Web Ontology Language (OWL) which is a W3C recommended standard. OWL 2 EL provides sufficient expressivity to model large hierarchical ontologies as well as streaming data such as traffic, while at the same time allows for efficient reasoning services. Existing reasoners for OWL 2 EL, however, use only a single machine and are thus constrained by memory and computational power. At the same time, the automated generation of ontological information from streaming data and text can lead to very large ontologies which can exceed the capacities of these reasoners. We thus describe a distributed reasoning system that scales well using a cluster of commodity machines. We also apply our system to a real case on city traffic data and show that it can handle volumes which cannot be handled by current single machine reasoners.

Slider: an Efficient Incremental Reasoner

J. Chevalier (Fragments RDFS et pD*) - streaming



Raisonner en OWL DL à
l'échelle du Web

Problème non résolu

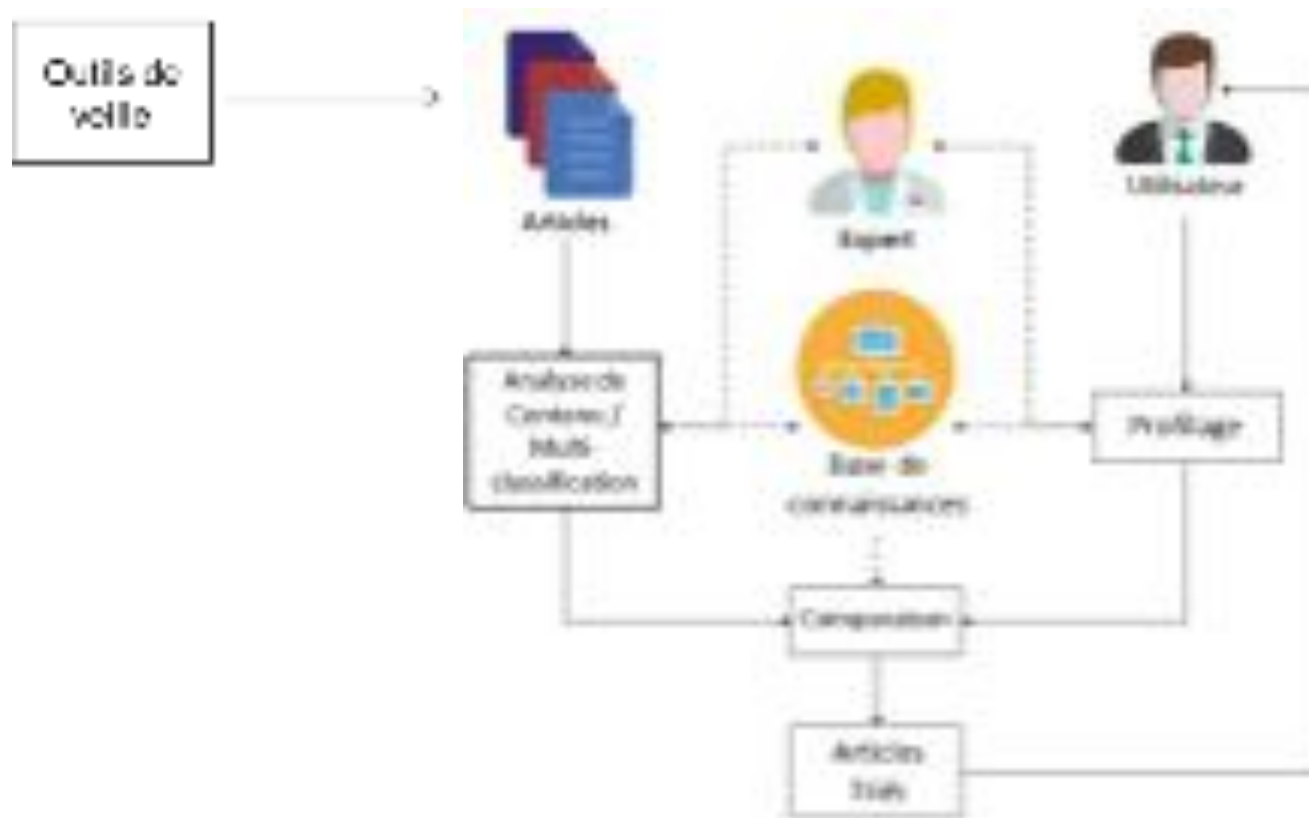


Le projet de l'entreprise ...

La sémantique, la veille et le web



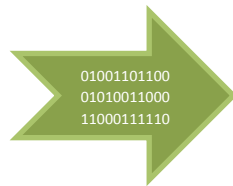
Outils actuels de veille



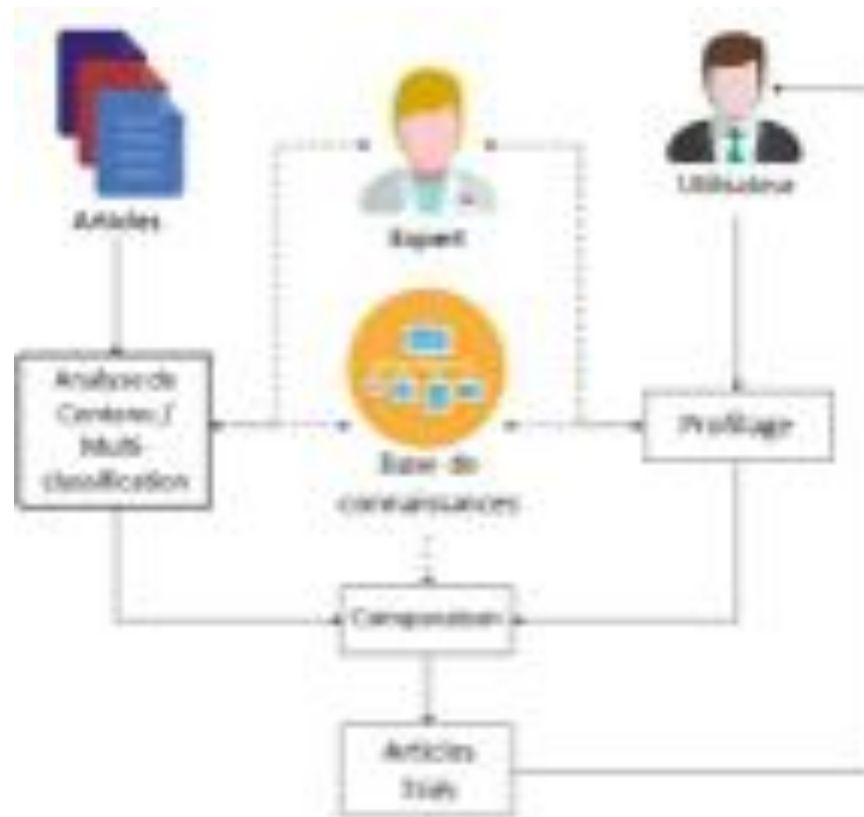
Quels sont les leviers pour faciliter cette étape ?

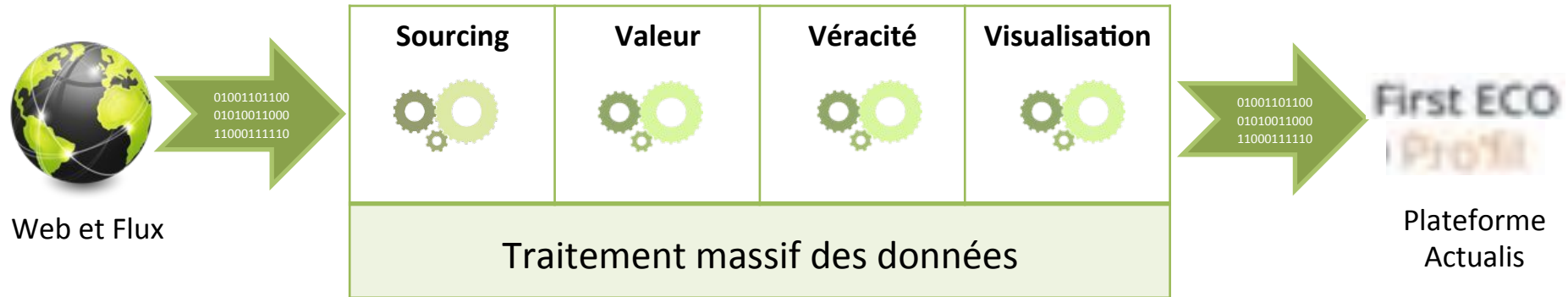


Web et Flux



Volume !
!!





Sourcing sémantique Croisement de l'information



Processus de **classification multi-label hiérarchique sémantique** pour le **Big Data**



Processus de **classification multi-label hiérarchique sémantique** pour le **Big Data**



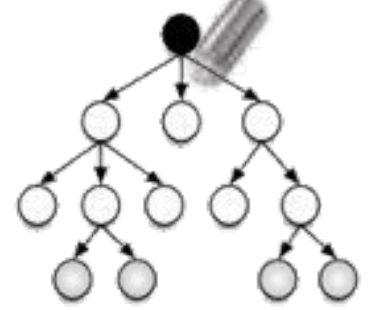
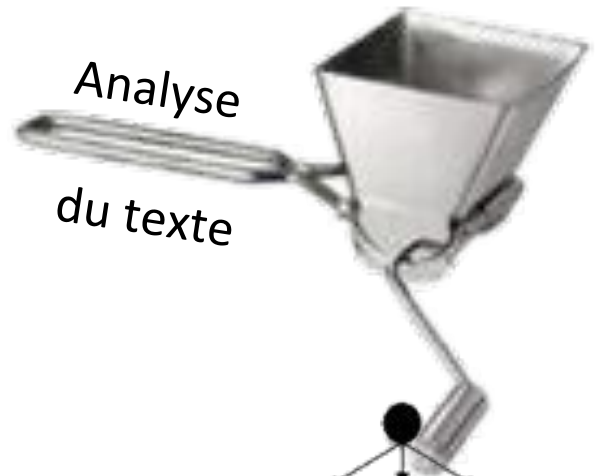
Processus de **classification multi-label hiérarchique sémantique** pour le **Big Data**



=



Comment générer
cette taxonomie à
l'échelle du Web ?



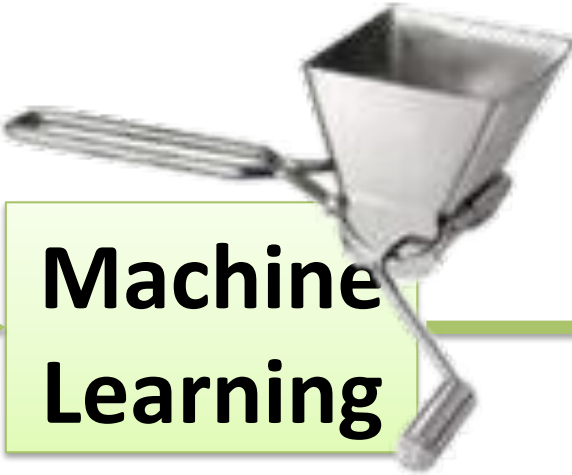
Modèle prédictif

Le problème ...

Comment coder cette moulinette ?



Exemples



Machine Learning



Programme

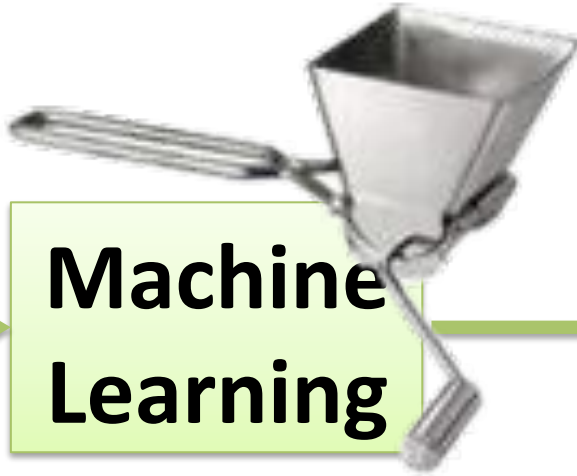
Par exemple ...

Donner un sens au sens des mots !

- La **veille** sanitaire est l'action de surveiller l'état de santé d'une population afin de prévenir des menaces d'épidémies.
- La **veille** des jours fériés, la journée de travail finit une heure plus tôt.
- J'ai entendu plusieurs versions sur les conditions pour avoir le jour férié rémunéré : il faut travailler la **veille** du jour et le lendemain et avoir un ...
- Ainsi la garde de nuit doit être apte au cours de sa **veille** nocturne, ...
- La garde de nuit **veille** sur la ville. Archivé. Quatre agents, deux chiens : la police municipale veillera sur la tranquillité publique le ...



**Phrases
+
Dictionnaires
& Définitions**



**Machine
Learning**



**Décision du
sens du mot**

« La **veille**
concurrentielle es
t l'activité
continue et en
grande partie
itérative qui vise à
une surveillance
active de ... »



**Décision du
sens du mot**



Veille³











jour, précédent, un, autre



fait, de, ne, pas, dormir, éveil



garde, fait, pendant, la, nuit



jour, précédent, un, autre



fait, de, ne, pas, dormir, éveil



garde, fait, pendant, la, nuit

Le café noir très fort
qu'il a pris au début
de la **veille**,
empêche Arsène
André de dormir.



jour, précédent, un, autre



fait, de, ne, pas, **dormir**, éveil



garde, fait, pendant, la, nuit

Le café noir très fort
qu'il a pris au début
de la **veille**,
empêche Arsène
André de **dormir**.



jour, précédent, un, autre



fait, de, ne, pas, dormir, éveil, *café, noir, très, fort,*
début, la, empêche, Arsène, André



garde, fait, pendant, la, nuit

Le café noir très fort
qu'il a pris au début
de la **veille**,
empêche Arsène
André de **dormir**.



George E. P. Box

Essentially, all models are wrong, but some are useful.

Empirical Model-Building and Response Surfaces (1987) p. 424



jour (35), précédent (62), un (36), autre (21), finis (3),
feux(1), longues(33), parlottes (2), précédent(40),
immédiatement(5), avant(37), etc.



fait (3), de (21) , ne (10), pas(30), dormir(64), éveil(45),
café(2), noir(10), très (3), fort(2), début(1), la(21),
empêche(40), Arsène(1), André(2), etc.



garde(50), fait(12), pendant(23), la(37), nuit(15),
préparez(22), concurrentielle(6), à(38), rester(9),
éveillé(11), difficile(17), alerte(55), surveillance(19), etc.

« La **veille**
concurrentielle

est l'activité
continue et en
grande partie
itérative qui
vise à une

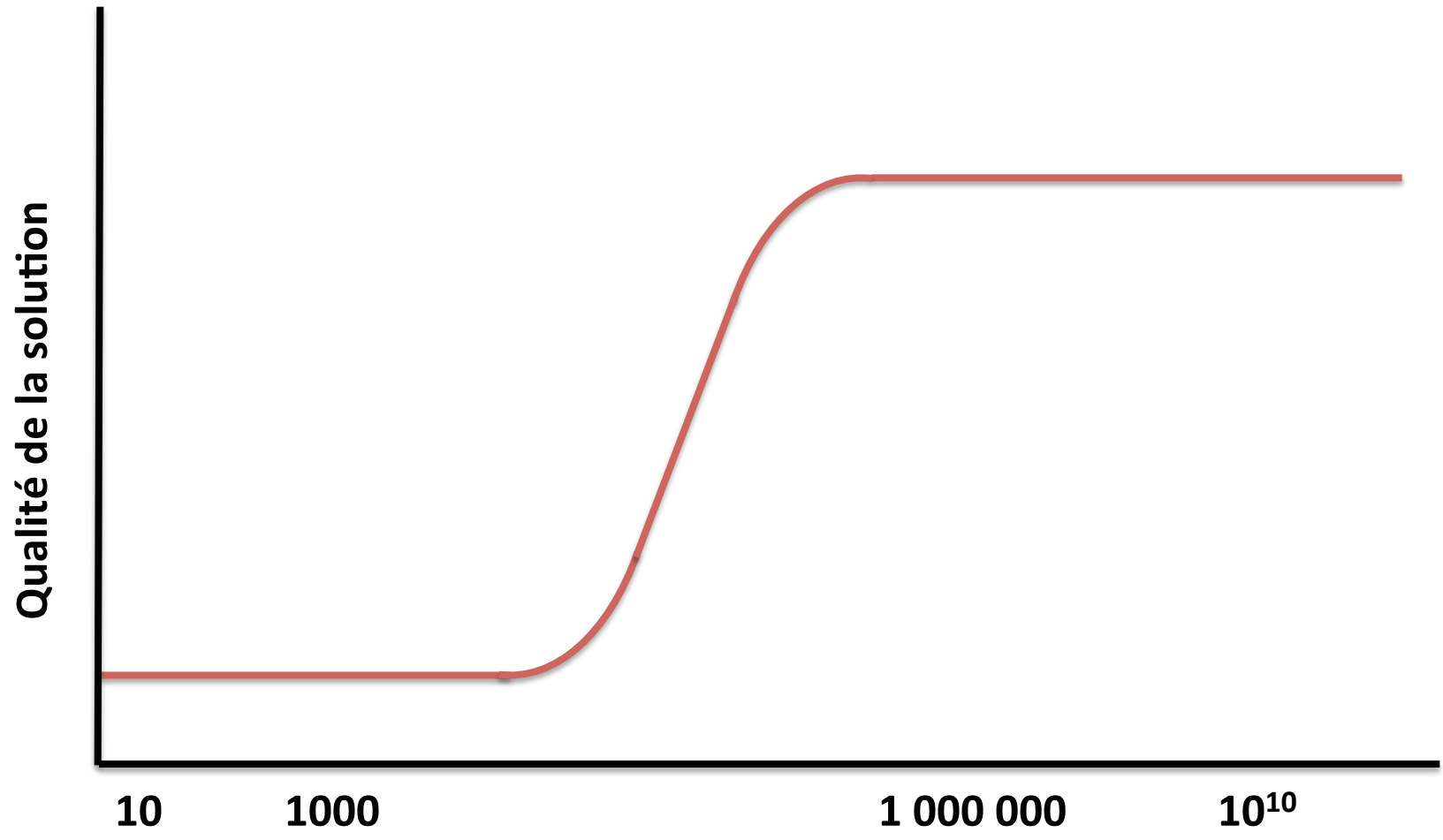
surveillance

active de ... »



garde(50), fait(12), pendant(23), la(37), nuit(15),
préparez(22), **concurrentielle**(6), à(38), rester(9),
éveillé(11), difficile(17), alerte(55), **surveillance**(19), etc.

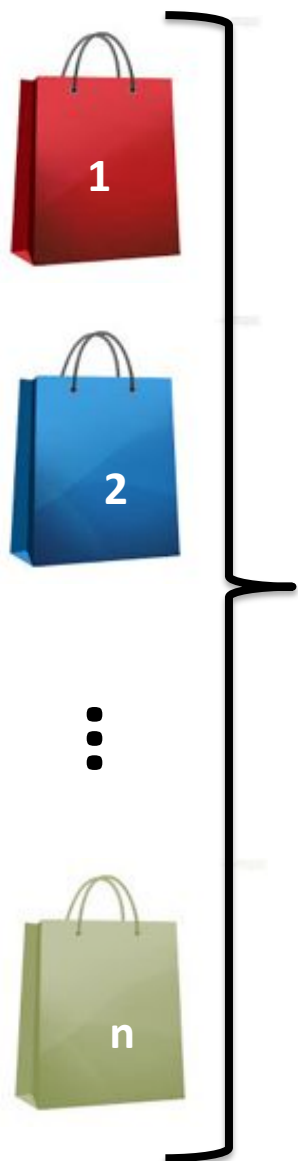
Seuils des données



Et notre problème ...

Développement d'une solution Big Data

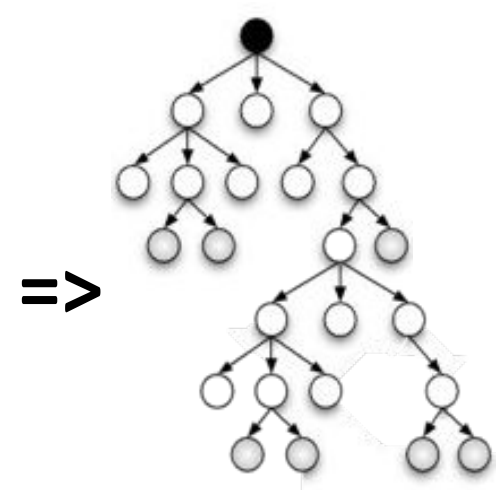




50 000 sacs de 50 000 mots
 $>10^{10}$ cases mémoires

$10^{10} = 10 \times 10 \times 10 \times 10 \times 10 \times 10 \times 10 \times 10 \times 10 \times 10$

%	W_1	W_2	W_3	W_4	W_5	W_6	W_7
Tax ₁	0	0	5	0	5	25	25
Tax ₂	0	75	0	0	0	75	5
Tax ₃	0	0	75	0	25	0	0
Tax ₄	5	25	25	0	5	93	25
Tax ₅	95	0	0	0	60	0	5
Tax ₆	0	60	0	95	0	0	90
Tax ₇	5	98	5	60	25	0	79



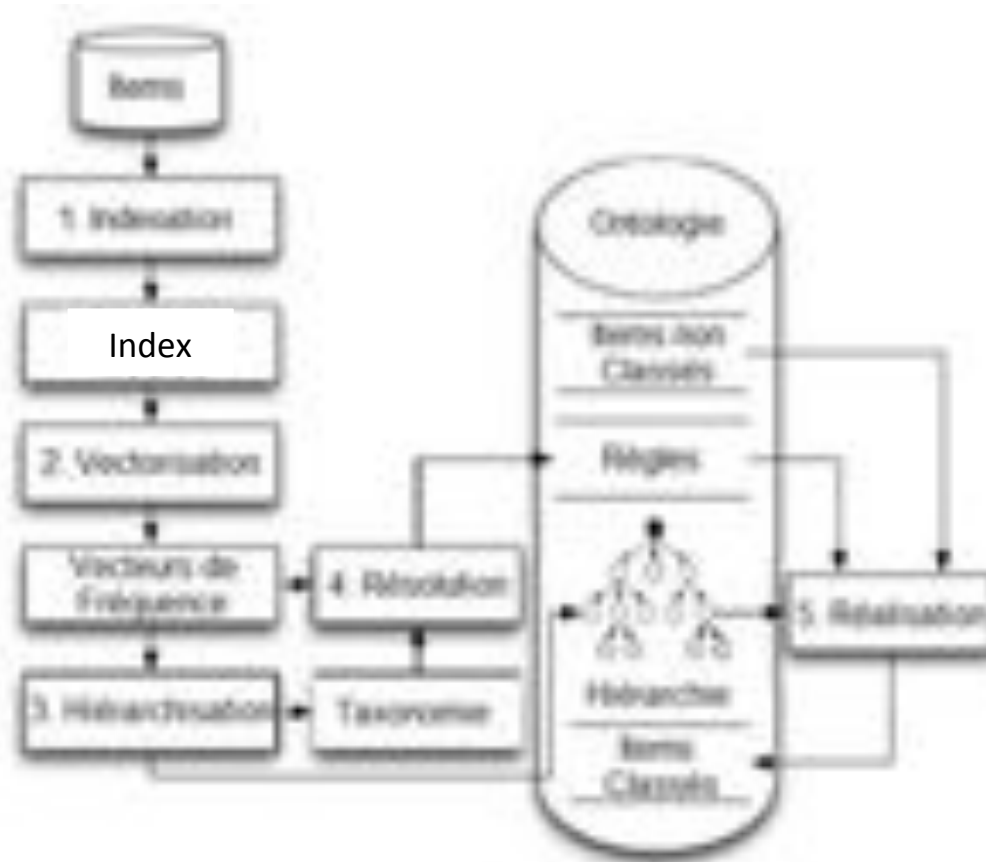
La méthode ...

En cinq phases

Usage des logiques descriptives

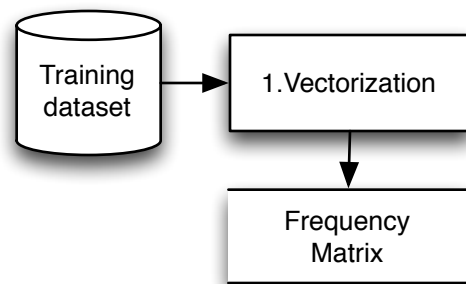
Prototype **vs 1.0**





Ontologies OWL et moteur d'inférence

1. Vectorisation

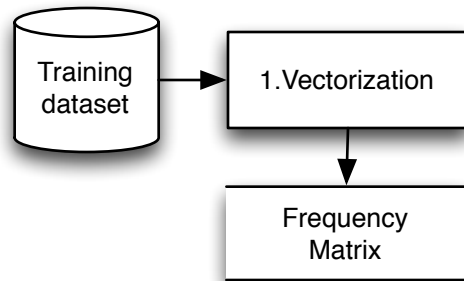


%	W_1	W_2	W_3	W_4	W_5	W_6	W_7
Tax₁	0	0	5	0	5	25	25
Tax₂	0	75	0	0	0	75	5
Tax₃	0	0	75	0	25	0	0
Tax₄	5	25	25	0	5	93	25
Tax₅	95	0	0	0	60	0	5
Tax₆	0	60	0	95	0	0	90
Tax₇	5	98	5	60	25	0	79

Tax_i sont les termes de la taxonomie

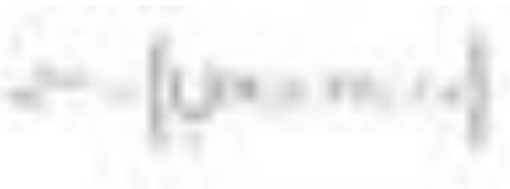
W_i sont les mots fréquents dans les textes

1. Vectorisation



%	W ₁	W ₂	W ₃	W ₄	W ₅	W ₆	W ₇
Tax ₁	0	0	5	0	5	25	25
Tax ₂	0	75	0	0	0	75	5
Tax ₃	0	0	75	0	25	0	0
Tax ₄	5	25	25	0	5	93	25
Tax ₅	95	0	0	0	60	0	5
Tax ₆	0	60	0	95	0	0	90
Tax ₇	5	98	5	60	25	0	79

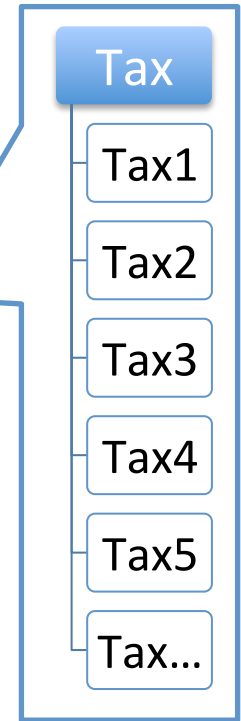
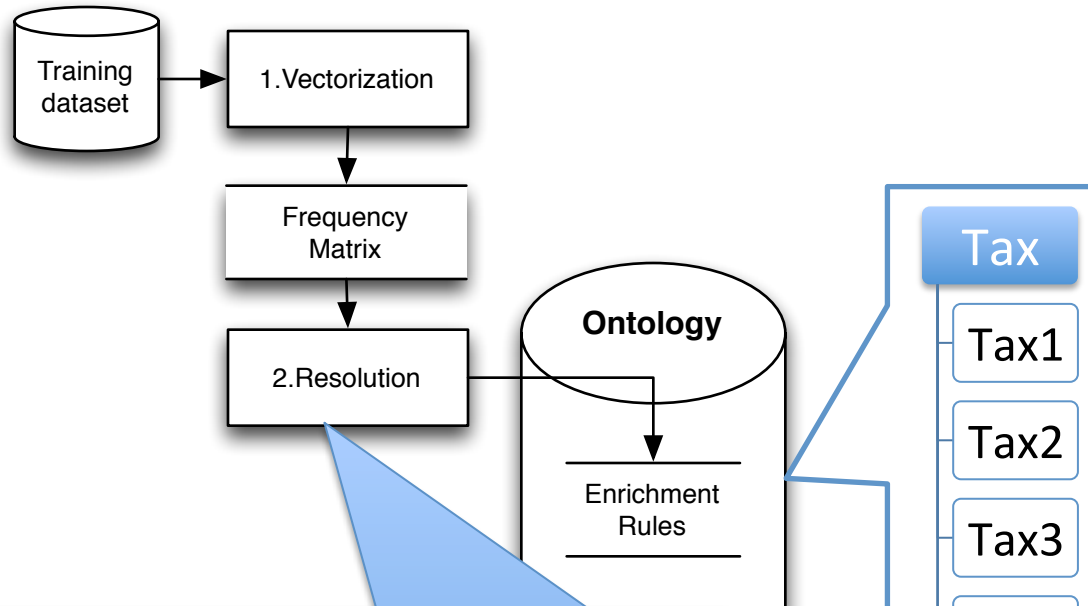
91



70

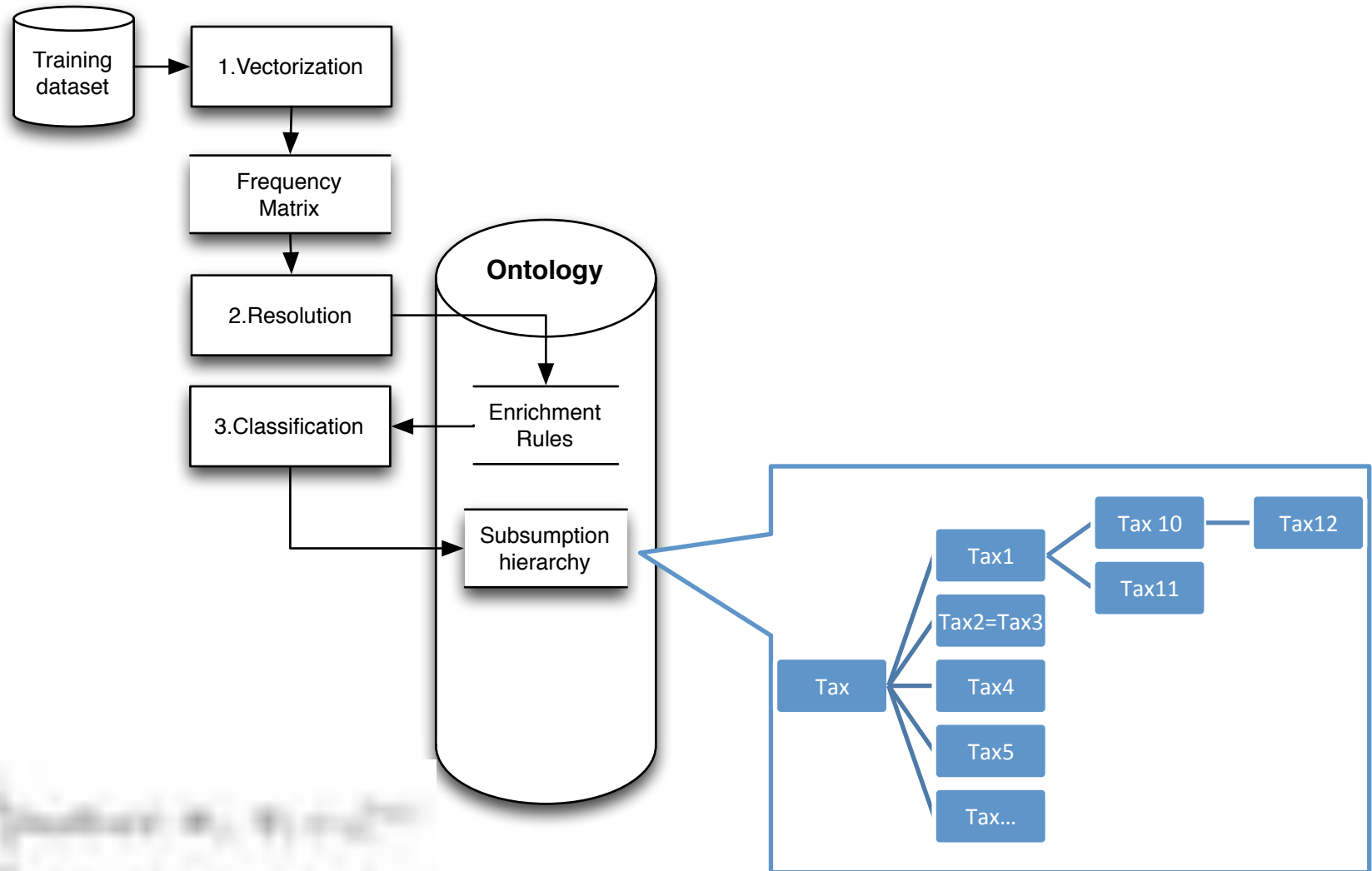


2. Résolution



%	W ₁	W ₂	W ₃	W ₄	W ₅	W ₆	W ₇
Tax ₁	0	0	5	0	5	25	25
Tax ₂	0	75	0	0	0	75	5
Tax ₃	0	0	75	0	25	0	0
Tax ₄	5	25	25	0	5	93	25
Tax ₅	95	0	0	0	60	0	5
Tax ₆	0	60	0	95	0	0	90
Tax ₇	5	98	5	60	25	0	79

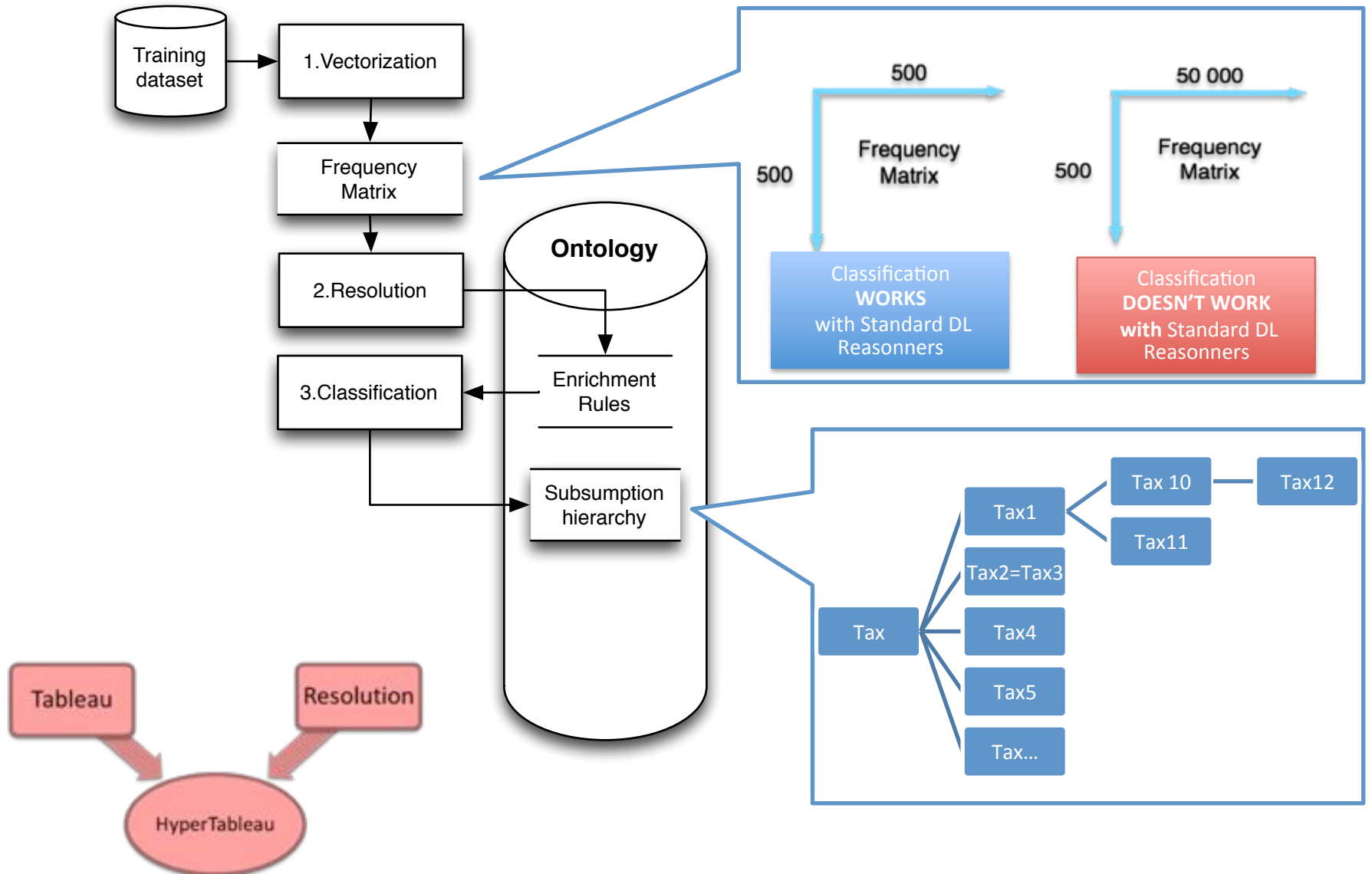
Faint, illegible text at the bottom of the slide, possibly bleed-through from the reverse side.



α-rules	FaCT++	HermiT	Pellet
i7 4Go DDR3 Xeon E3 24Go DDR3	50 s -	n.e.m. ¹ 8 h	n.e.m. 18 h
$\alpha\beta$-rules	FaCT++	HermiT	Pellet
i7 4Go DDR3 Xeon E3 24Go DDR3 Xeon E5 128Go DDR3	n.e.m. n.e.m. 2 h / out ³	n.e.m. out ² out	n.e.m. out out

Evaluation	Precision	Recall	F1-Mesure
Proposal with A-type rules	30%	6%	10%
HOMER [16]	-	-	25%

*Our way to create rules (i.e. with an **average of 10 terms** for all rules) has the consequence the creation of some ruleless classes. With this method, for our **983 classes, only 427 have rules**. There are 556 classes without labeling rule (obviously, these classes should have had β rules). So there are classes that the predictive model can not affect. This impacts very negatively the Recall.*



La méthode Big Data

En cinq phases

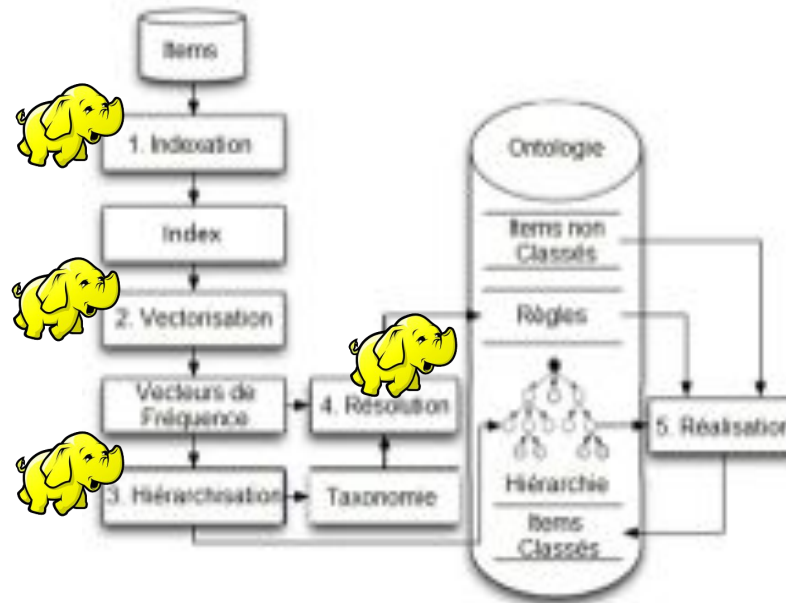
Usage des LD et des règles logiques

Prototype **vs 2.0**



Montée en charge

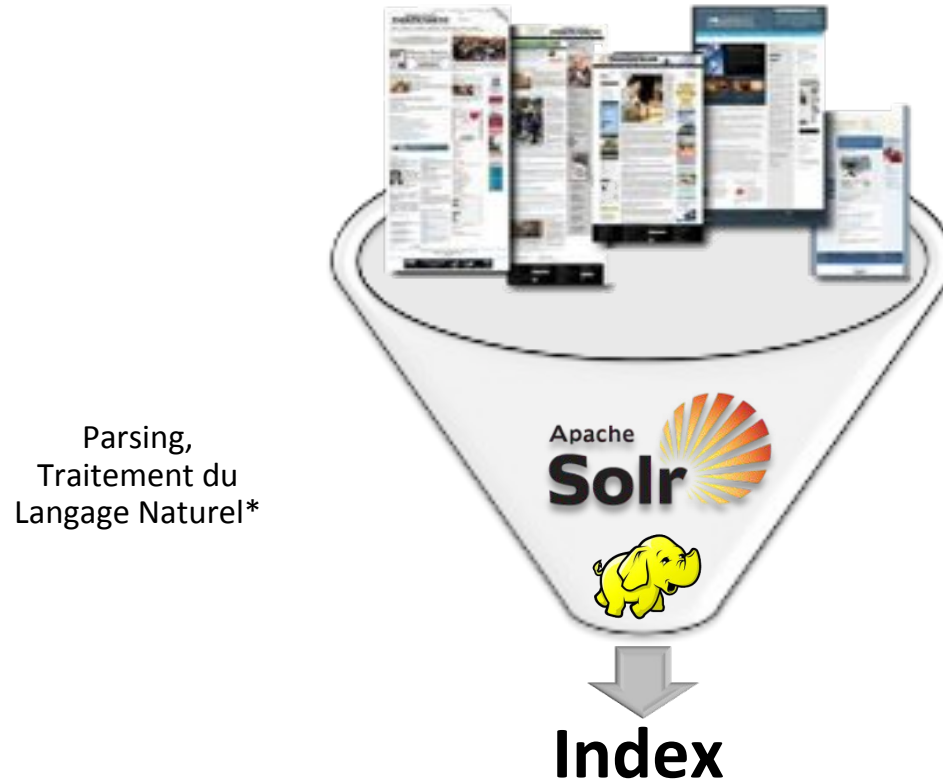
Utilisation du **framework Hadoop** et du modèle **MapReduce**



Comment distribuer chacune des phases du processus ?

Dean, J., & Ghemawat, S. (2008). MapReduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1), 107-113.

Génération d'un index des termes



* Tokenisation, Lemmatisation, Suppression des mots vides

Détection des termes pertinents basée sur l'index

%	w_1	w_2	w_3	w_4	w_5	w_6	w_7
Tax_1	0	0	5	0	5	25	25
Tax_2	0	75	0	0	0	75	5
Tax_3	0	0	75	0	25	0	0
Tax_4	5	25	25	0	5	93	25
Tax_5	95	0	0	0	60	0	5
Tax_6	0	60	0	95	0	0	90
Tax_7	5	98	5	60	25	0	79

Les termes dont la fréquence est inférieure à un certain seuil sont rejetés

- Uni-grams (termes)
- N-grams (collocations) - Amélioration

Génération des vecteurs de fréquences

Liste de collocation
{ uni-grams + n-grams }

%	w_1	w_2	w_3	w_4	w_5	w_6	w_7
Tax ₁	0	0	5	0	5	25	25
Tax ₂	0	75	0	0	0	75	5



%
Tax ₁
Tax ₂
Tax ₃
Tax ₄
Tax ₅
Tax ₆
Tax ₇

Liste des fréquences

%	w_1	w_2	w_3	w_4	w_5	w_6	w_7
Tax ₁	0	0	5	0	5	25	25
Tax ₂	0	75	0	0	0	75	5
Tax ₃	0	0	75	0	25	0	0
Tax ₄	5	25	25	0	5	93	25
Tax ₅	95	0	0	0	60	0	5
Tax ₆	0	60	0	95	0	0	90
Tax ₇	5	98	5	60	25	0	79

Seuil de pertinence

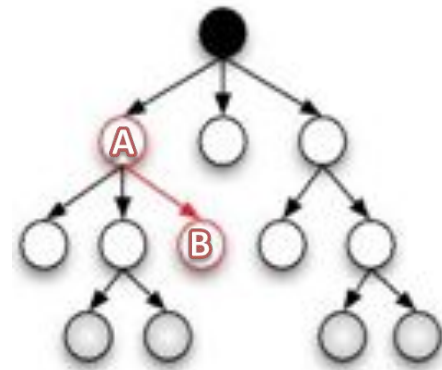


Algorithme des Stripes 



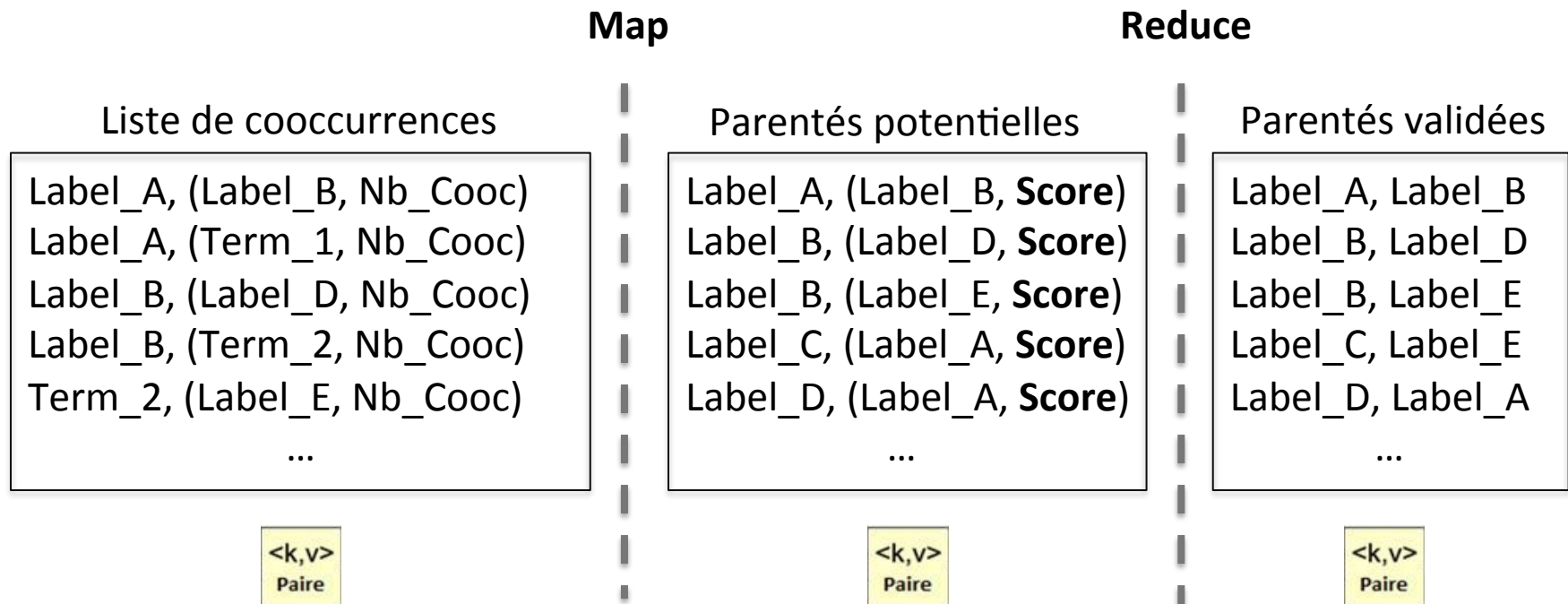
Algorithme de subsumption

%	W_1	W_2	W_3	W_4	W_5	W_6	W_7
Tax ₁	0	0	5	0	5	25	25
Tax ₂	0	75	0	0	0	75	5
Tax ₃	0	0	75	0	25	0	0
Tax ₄	5	25	25	0	5	93	25
Tax ₅	95	0	0	0	60	0	5
Tax ₆	0	60	0	95	0	0	90
Tax ₇	5	98	5	60	25	0	79



$$P(p|x) \geq t, P(x|p) < t$$

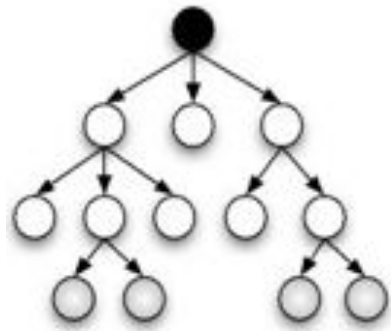
Algorithme de subsomption MapReduce



$$P(p|x) \geq t, P(x|p) < t$$

Score calculé à partir des cooccurrences

Hiérarchie chargée dans un triplestore



Hiérarchie



Stardog

Génération de règles de classification

Clauses de Horn (SWRL et Stardog)

%	W ₁	W ₂	W ₃	W ₄	W ₅	W ₆	W ₇
Tax ₁	0	0	5	0	5	25	25
Tax ₂	0	75	0	0	0	75	5
Tax ₃	0	0	75	0	25	0	0
Tax ₄	5	25	25	0	5	93	25
Tax ₅	95	0	0	0	60	0	5
Tax ₆	0	60	0	95	0	0	90
Tax ₇	5	98	5	60	25	0	79

Seuils α et β

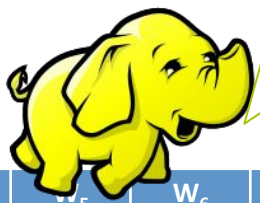
$\alpha : \text{Item}(?d) \wedge \text{Word}(?w1) \wedge \text{hasWord}(?d, ?w1) \rightarrow \text{Tax}(?d1)$

$\beta : \quad \text{Item}(?d) \wedge \text{Word}(?w1) \wedge \text{hasWord}(?d, ?w1) \wedge$
 $\text{Word}(?w2) \wedge \text{hasWord}(?d, ?w2) \rightarrow \text{Tax}(?d1)$

Item: nouveau document

Génération de règles de classification

Approche MapReduce
« diviser pour régner »



%	W ₁	W ₂	W ₃	W ₄	W ₅	W ₆	W ₇
Tax ₁	0	0	5	0	5	25	25
Tax ₂	0	75	0	0	0	75	5
Tax ₃	0	0	75	0	25	0	0
Tax ₄	5	25	25	0	5	93	25
Tax ₅	95	0	0	0	60	0	5
Tax ₆	0	60	0	95	0	0	90
Tax ₇	5	98	5	60	25	0	79

Tax₂, (Term₂, Nb_Cooc)
 Tax₂, (Term₆, Nb_Cooc)
 Tax₂, (Term₈, Nb_Cooc)
 Tax₂, (Term₉, Nb_Cooc)
 ...

<k,v>
Paire

Génération de règles de classification

$Item(?d) \wedge Word(?w1) \wedge hasWord(?d, ?w1)$
 $\rightarrow Tax(?d1)$



Règles

Stardog

Utilisation des règles pour classer

Règles

$Item(?d) \wedge Word(?w1) \wedge hasWord(?d, ?w1)$
 $\rightarrow Tax(?d1)$



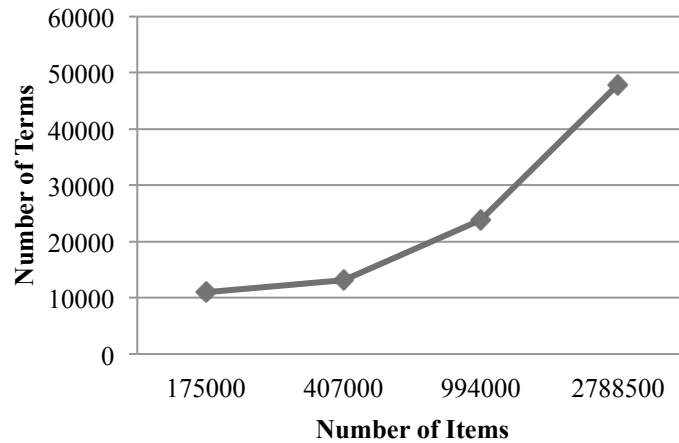
Classés dans Stardog

	L ₁	L ₂	L ₃	L ₄	L ₅
Doc ₁	X	-	-	-	X
Doc ₂	-	X	X	-	X
Doc ₃	X	-	-	X	-
Doc ₄	X	-	-	-	-
Doc ₅	-	X	-	X	-
Doc ₆	-	-	X	-	X
Doc ₇	X	-	-	X	-

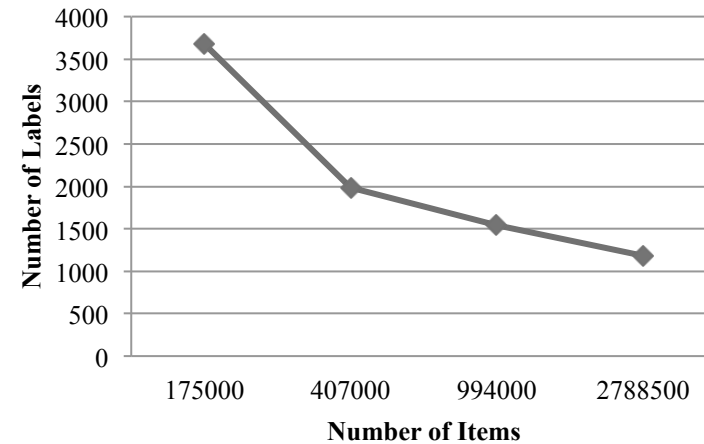
Vecteurs « document »

Evaluation de l'apprentissage

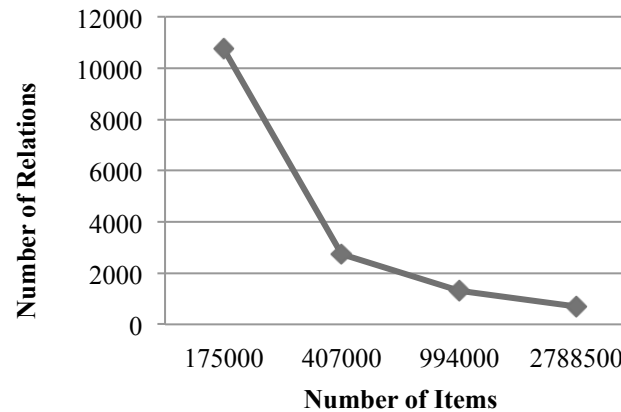
Learned Terms



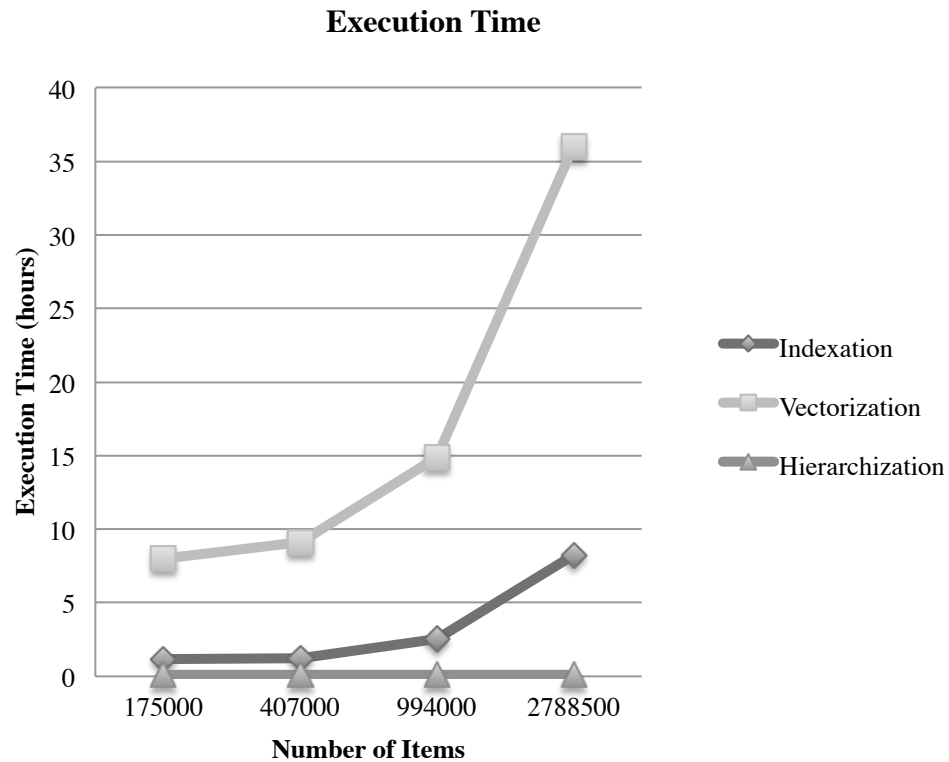
Learned Labels



Learned Subsumption relations



Evaluation de la performance du processus (4 nœuds):



Sémantique et connaissances

- La sémantique
- Problématique
- Application au Big Data

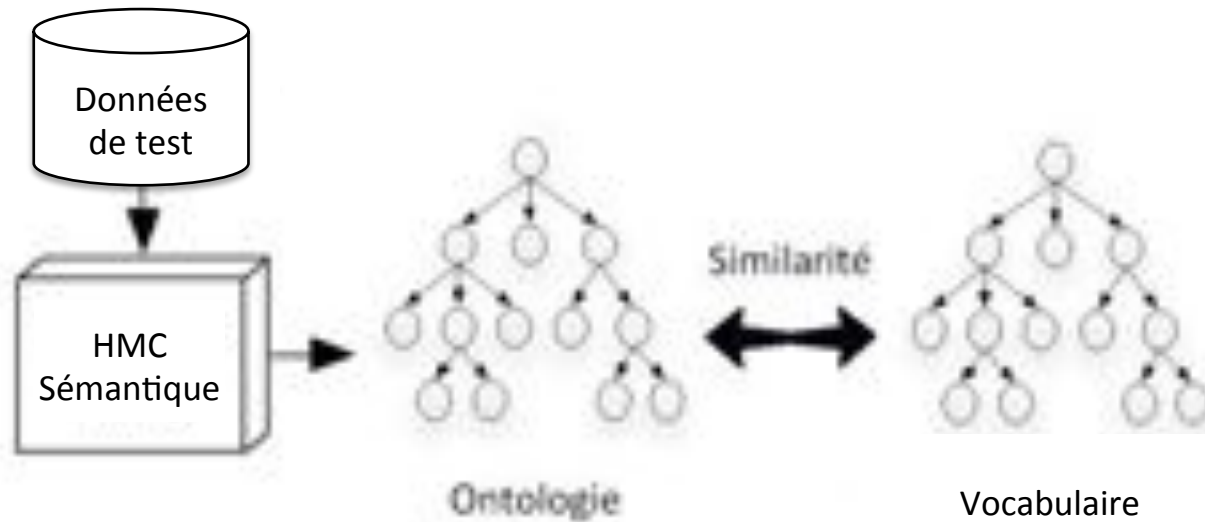


Merci ...

Questions ?

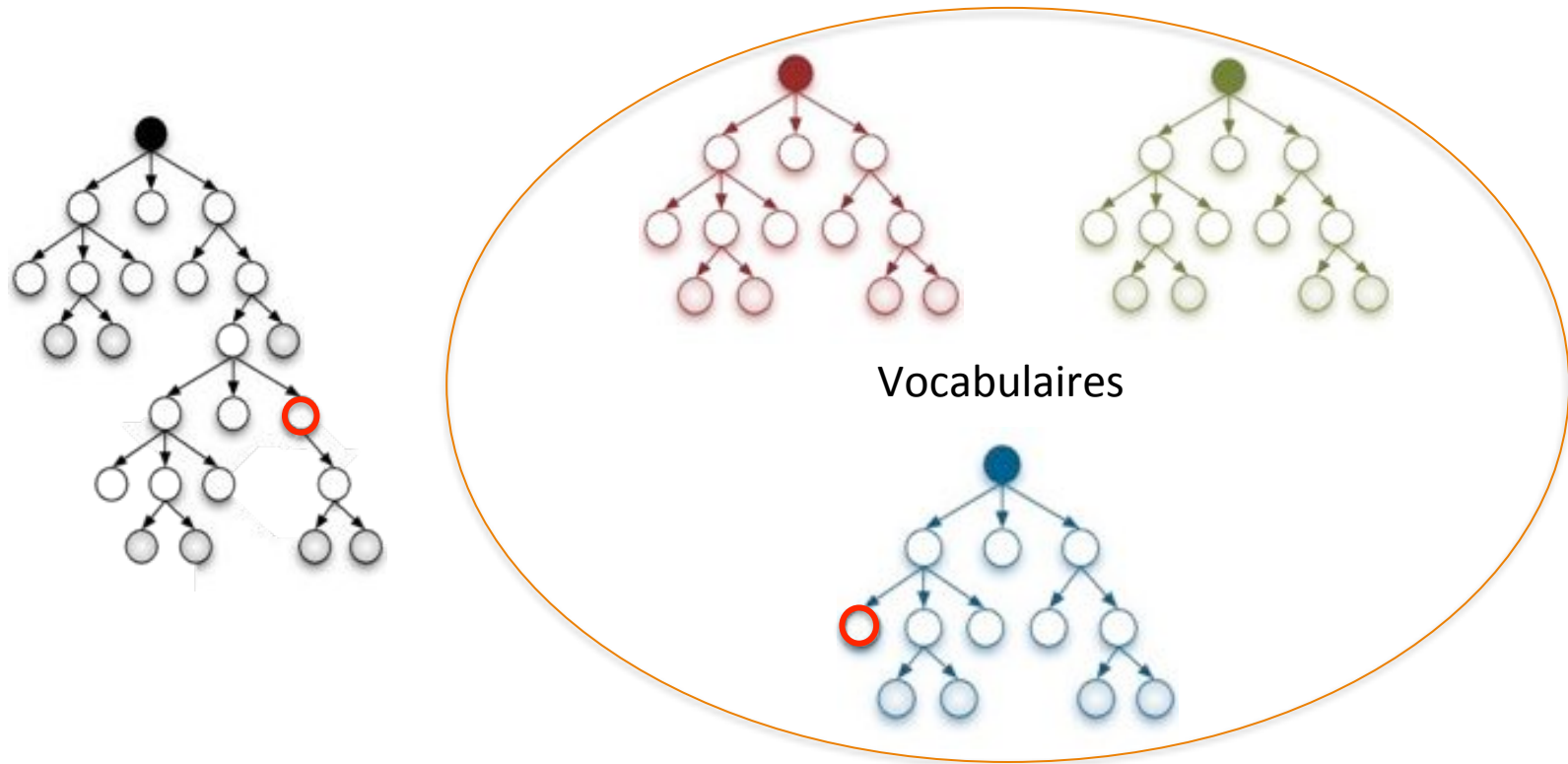


Evaluation de la qualité de la hiérarchie :



Profondeur, similarité des noeuds, topologie

Evaluation de la qualité de la classification :



Précision, rappel, F1

L'architecture ...

Ensemble d'outils logiciels

Experimentale



Fondamentale

Classification
Hiérarchique Multi-label Sémantique



Technologique

Semantic



Indexing



Distributed



Streaming

