

Lexique et analyse sémantique de textes

structures, acquisitions, calculs, et jeux de mots

M. Lafourcade

Christian Boitet	Université Joseph Fourier - Grenoble 1 - LIG	Père spirituel
Vladimir Fomichov	Higher School of Economics, Moscou	Rapporteur
Marianne Huchard	Université Montpellier 2 - LIRMM	Modèle
Violaine Prince	Université Montpellier 2 - LIRMM	Mère spirituelle
Christian Retoré	Université Bordeaux 1 - LaBRI	Rapporteur
Eric Wehrli	Université de Genève - LATL	Examineur
Michael Zock	Université de Marseille - LIF	Examineur
Pierre Zweigenbaum	LIMSI - CNRS	Rapporteur

Contexte

... et problématique

Analyse de textes

structures de représentation du sens

thématique : vecteurs ? mots-clés ?

sémantique : quelles relations entre syntagmes ?

désambiguïisation lexicale - inventaire de sens ?

fonctions lexicales ?

magn(fièvre) = forte fièvre (?) magn(maison) = villa

référence et substitution ?

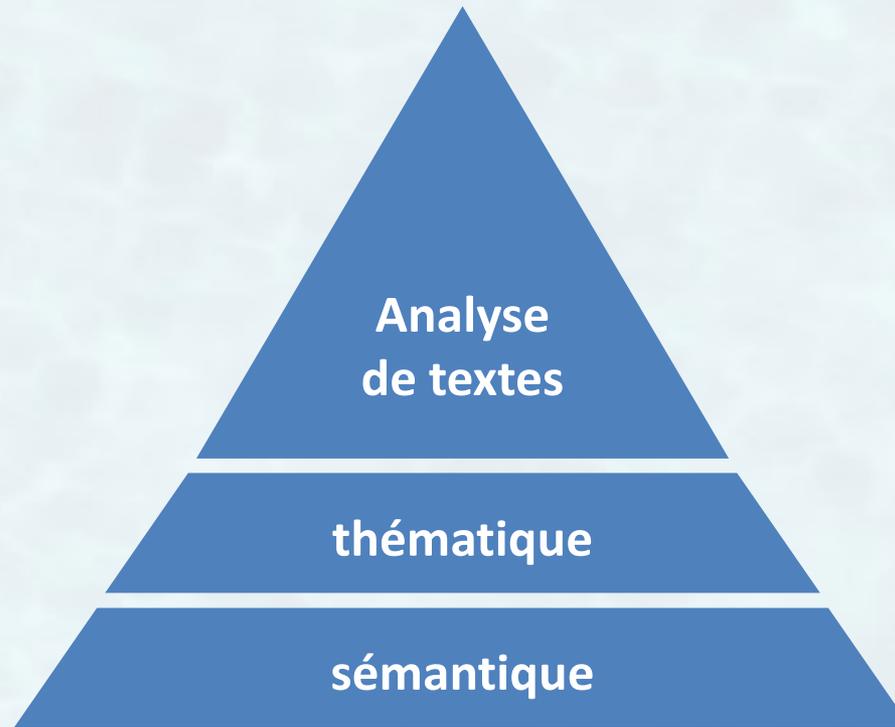
chat \Rightarrow matou / animal (?) mouche \Rightarrow drosophile / insecte

Calcul / acquisition des structures

vecteurs / réseaux lexico-sémantiques

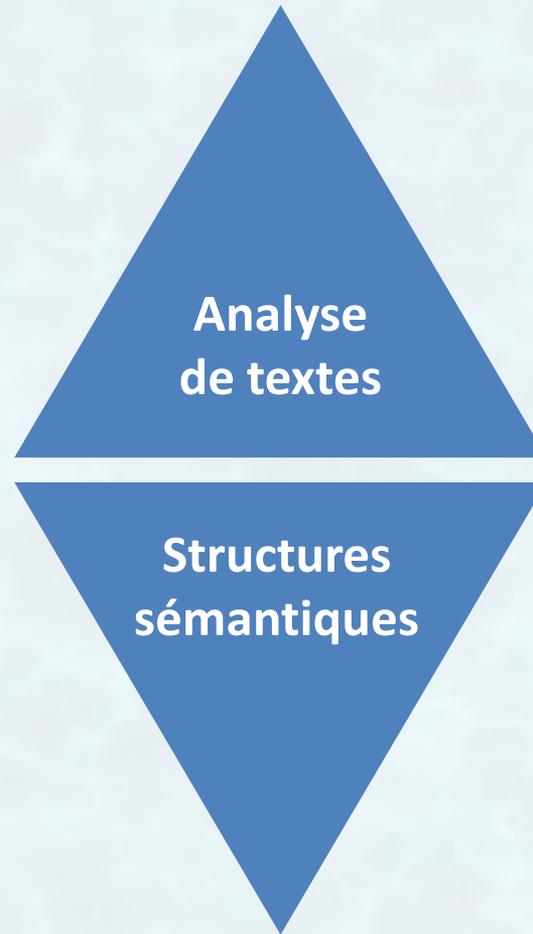
Contexte

structures, acquisitions, calculs, et jeux de mots



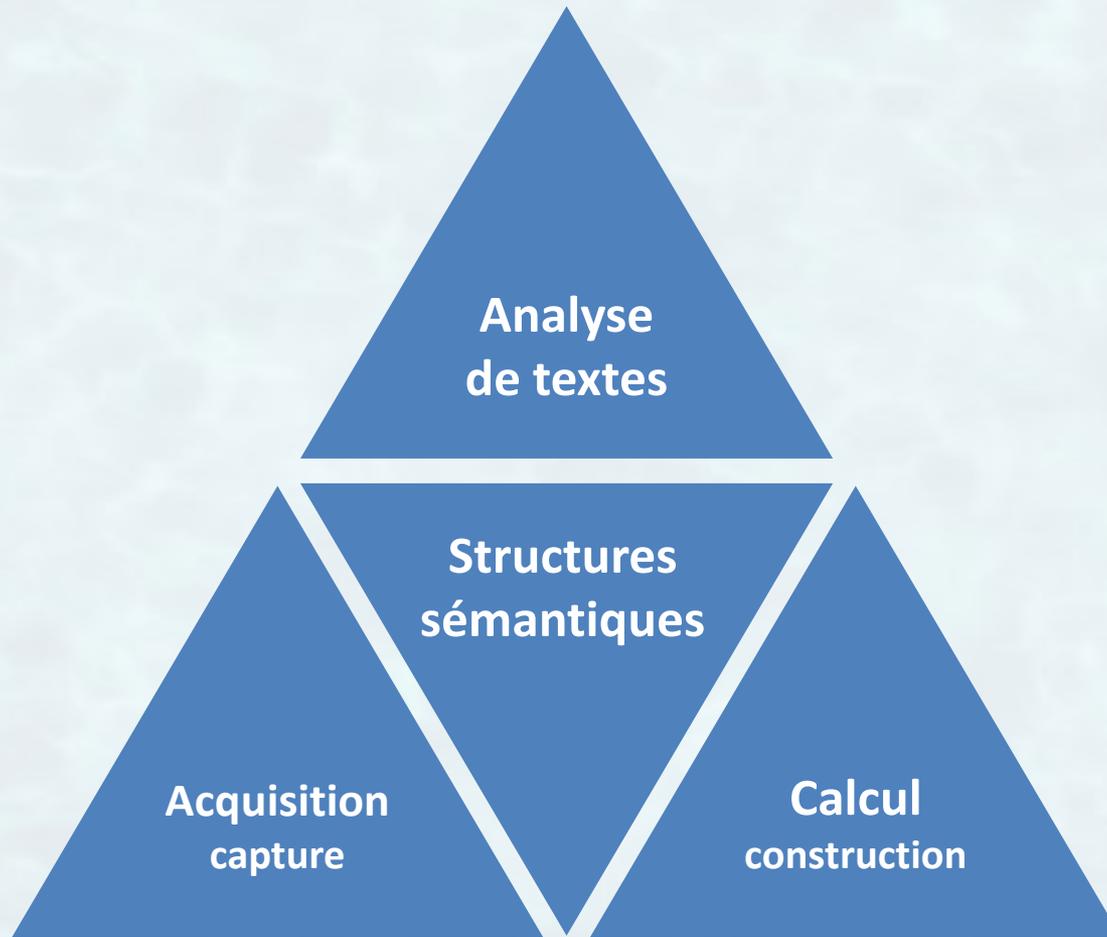
Contexte

structures, acquisitions, calculs, et jeux de mots



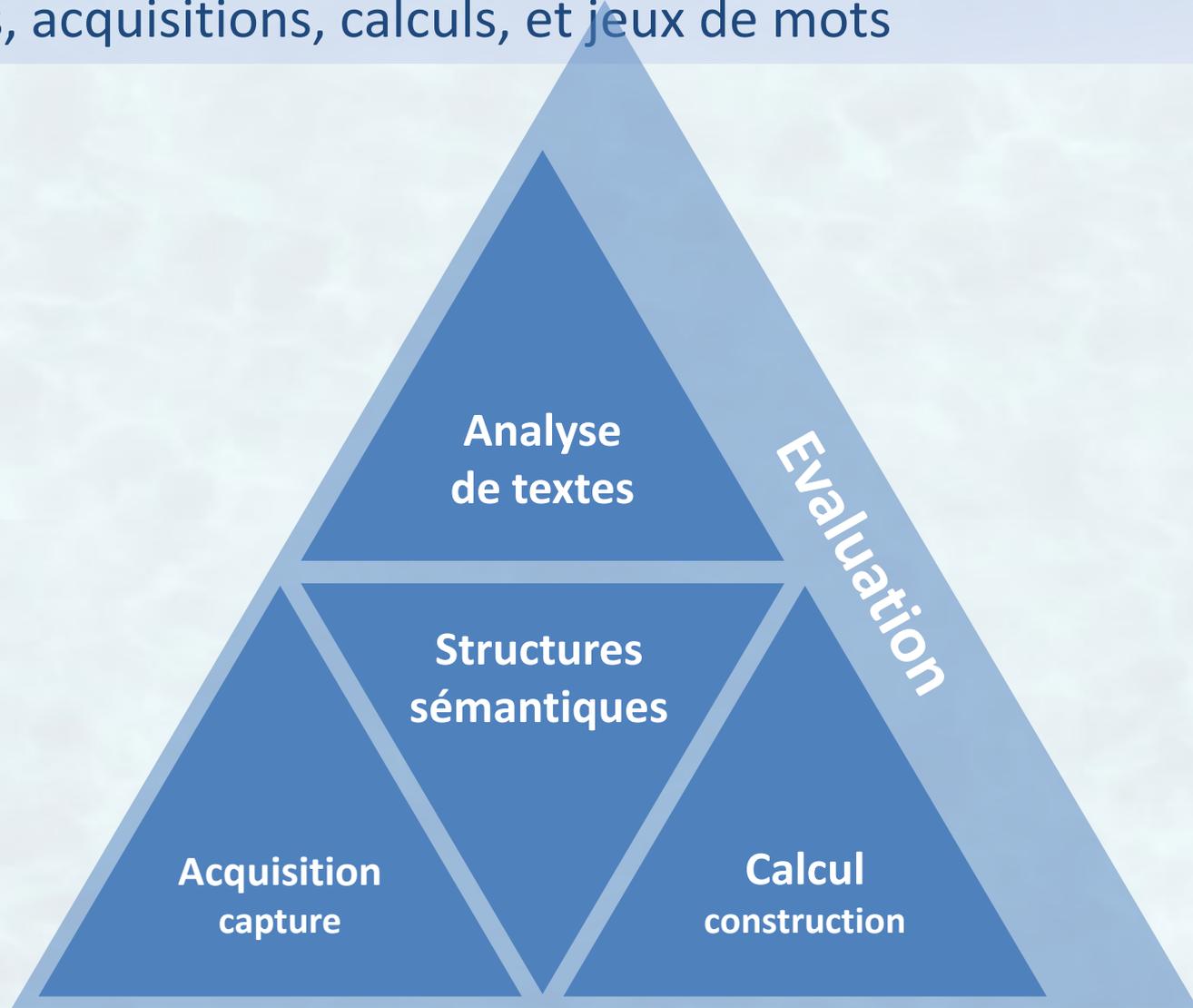
Contexte

structures, acquisitions, calculs, et jeux de mots



Contexte

structures, acquisitions, calculs, et jeux de mots



Calcul de vecteurs d'idées

Principes, variétés et construction
Fonctions lexicales
Pliage et dépliage

Acquisition de relations lexicales

JeuxDeMots : principe et résultats
Inventaire de sens / usages
Tentative d'évaluation qualitative : AKI

Analyse de textes

Thématique par vecteurs
Thématique par mot-clés
Métaheuristique bioinspirée

bouclage

permanence de l'acquisition

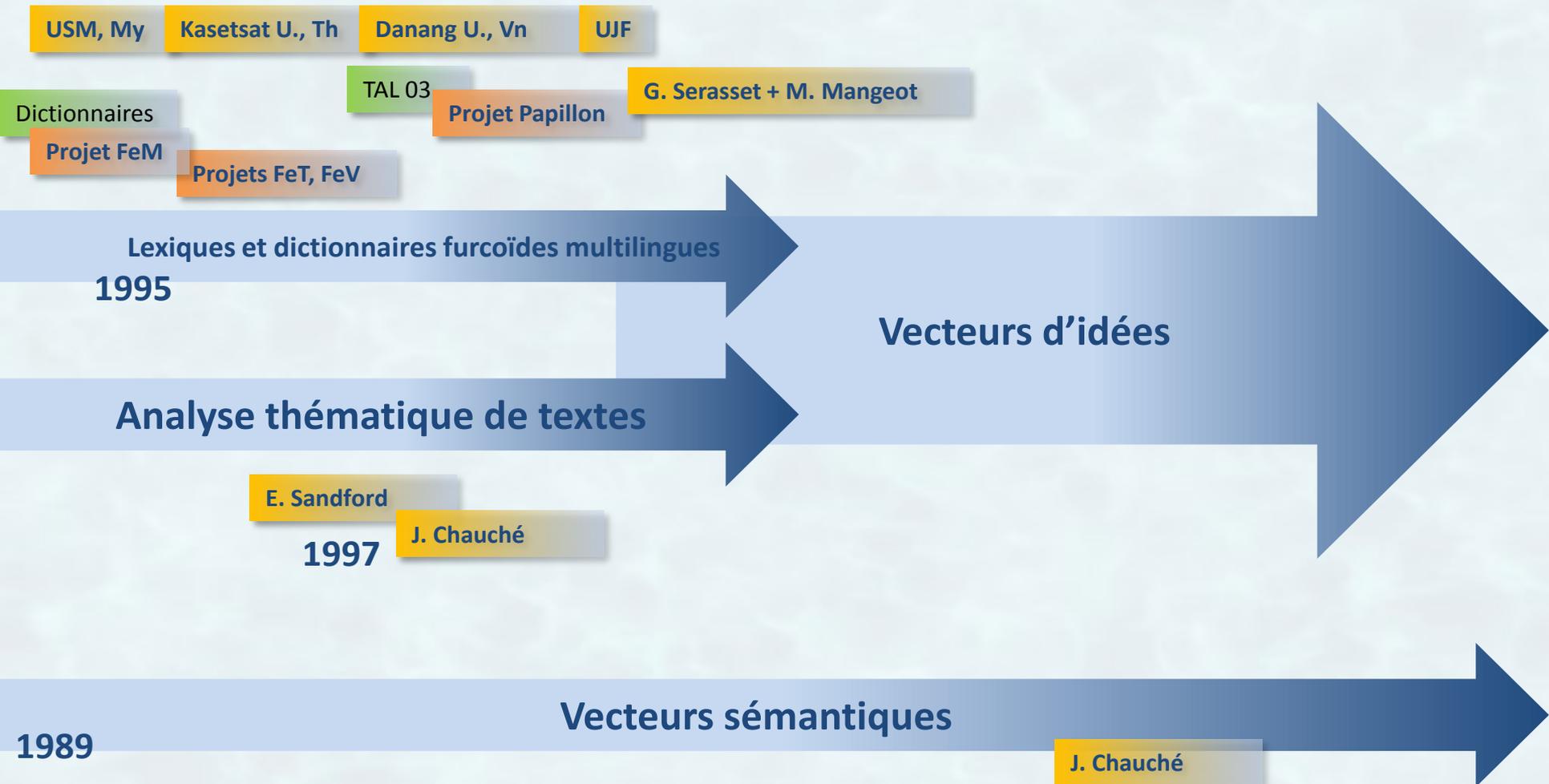
consensus populaire

raffinement des structures

propagation et diffusion

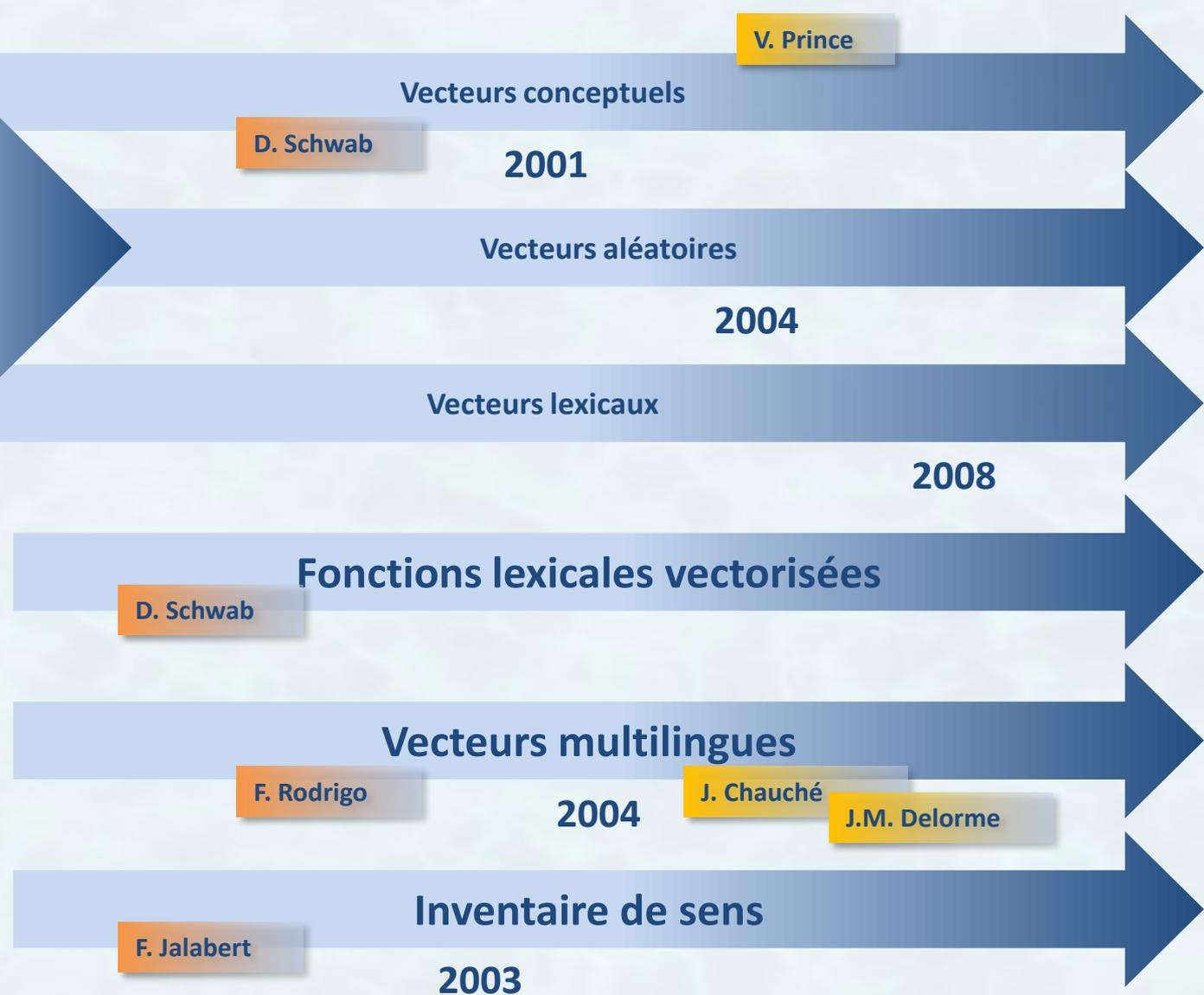
activation et inhibition

Chronologie



Chronologie

Vecteurs d'idées



Chronologie

Vecteurs d'idées

Vecteurs conceptuels

V. Prince

2001

Vecteurs aléatoires

2004

Vecteurs lexicaux

2008

Fonctions lexicales vectorisées

D. Schwab

Vecteurs multilingues

F. Rodrigo

2004

J. Chauché

J.M. Delorme

Inventaire de sens

F. Jalabert

2003

Vecteurs d'idées

vecteurs conceptuels

Une idée

= une combinaison de concepts = un vecteur

⇒ un espace vectoriel (en fait, une famille génératrice)

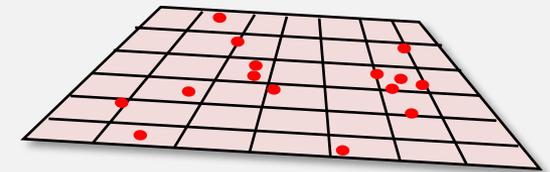
Un concept

= une idée importante/récurrente/élémentaire = un vecteur (ancrage)

= combinaison de lui-même + voisinage

Un Espace des sens

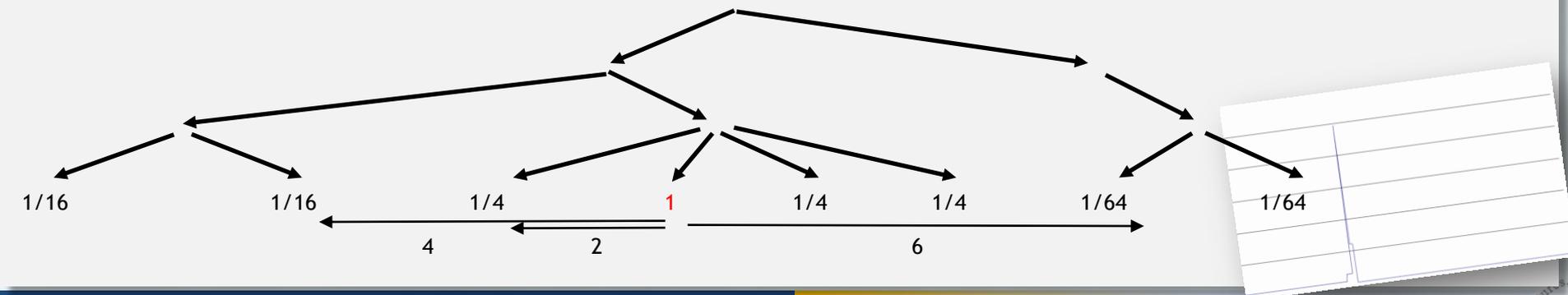
= une collection de vecteurs (un espace peuplé)



Exemple : thésaurus Larousse = 873 concepts

⇒ des vecteurs de dimension 873

Exploiter la hiérarchie du thésaurus pour représenter le voisinage



Vecteurs conceptuels

Construction à partir d'une ressource

Comment peupler un espace?

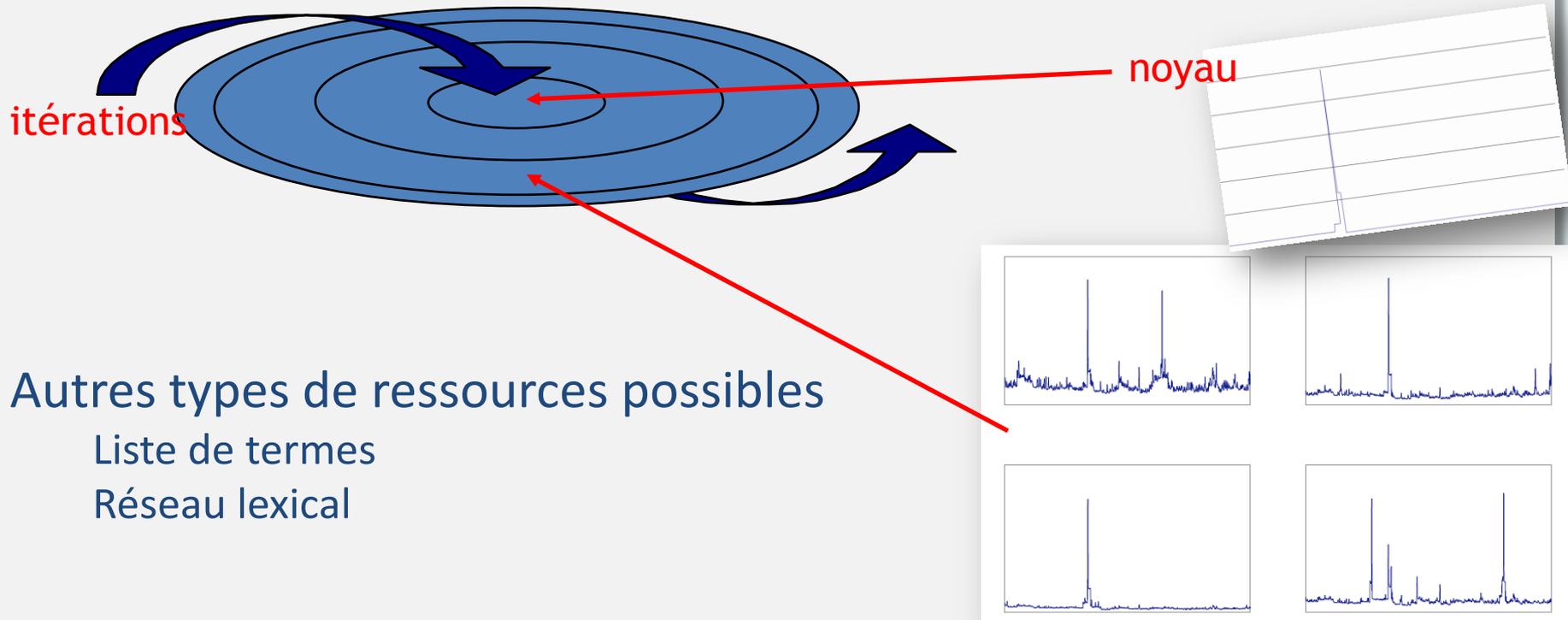
Par exemple, par analyse thématique

Définitions dictionnairiques (Larousse, Hachette, etc.)

Définition = genre + différence

Vecteur nul pour un terme non encore indexé

permanence de l'acquisition



Vecteurs d'idées

opération de comparaison

Distance angulaire $D_A(X, Y) = \text{angle}(X, Y)$

$$0 \leq D_A(X, Y) \leq \pi/2$$

si 0 alors colinéaire - même idée

si $\pi/2$ alors rien en commun

La norme n'intervient pas !

$$D_A(X, X) = 0$$

$$D_A(X, Y) = D_A(Y, X)$$

$$D_A(X, Y) + D_A(Y, Z) \geq D_A(X, Z)$$

$$D_A(X \oplus X, Y \oplus X) = D_A(X, Y \oplus X) \leq D_A(X, Y)$$

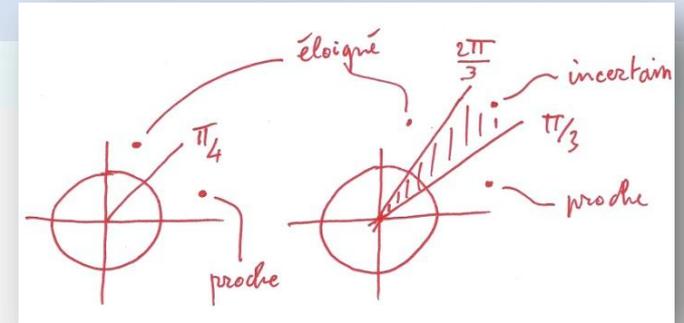
$$D_A(X \oplus Z, Y \oplus Z) \leq D_A(X, Y)$$

Interprétation géométrique

$$\text{Sim}(X, Y) = \cos(D_A(X, Y))$$

$$\text{Disim}^2(X, Y) = 1 - \text{Sim}^2(X, Y)$$

$$\text{Ponderation}(X, Y) = \text{Sim}(X, Y) / \text{Disim}(X, Y) = \cot(D_A(X, Y))$$



Vecteurs d'idées

opérations

Somme

$$V = X \oplus Y \Rightarrow v_i = x_i + y_i$$

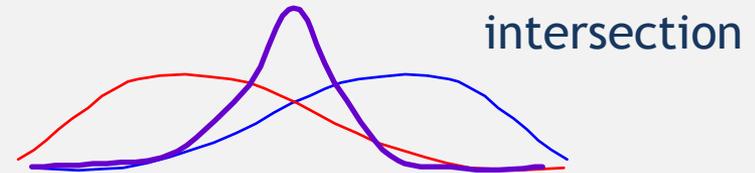
Élément neutre : **0** (vecteur nul)



Produit terme à terme

$$V = X \otimes Y \Rightarrow v_i = x_i * y_i$$

Élément neutre : **1** (vecteur unitaire)



Contextualisation faible

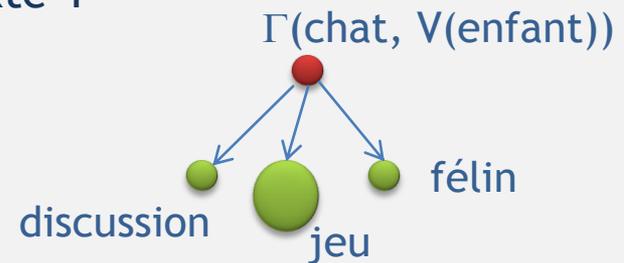
$$V = \gamma(X, Y) = X \oplus (X \otimes Y)$$

Sélection des idées de X présentes dans le contexte Y

Contextualisation forte du terme w avec le contexte Y

$$V(w) = \Gamma(w, Y) = \beta_1 V(w_1) \oplus \dots \oplus \beta_k V(w_k)$$

et $\beta_i = \text{Ponderation}(Y, V(w_k))$



Vecteurs d'idées

fonctions lexicales vectorielles

Distance de synonymie relative

$\text{Syn}_R(X, Y, Z)$ — Z est une référence

$$\text{Syn}_R(X, Y, Z) = D_A(\gamma(X, Z), \gamma(Y, Z))$$

$D_A(\text{charbon}, \text{nuit})$	= 0.9
$\text{Syn}_R(\text{charbon}, \text{nuit}, \text{couleur})$	= 0.4
$\text{Syn}_R(\text{charbon}, \text{nuit}, \text{noir})$	= 0.35

D. Schwab

V. Prince

TAL 02

NLPRS

utile en structuration terminologique
et analyse de métaphore

Travail sur l'antonymie

complémentaire
scalaire
duale

lexicalisation croissante

D. Schwab (DEA)

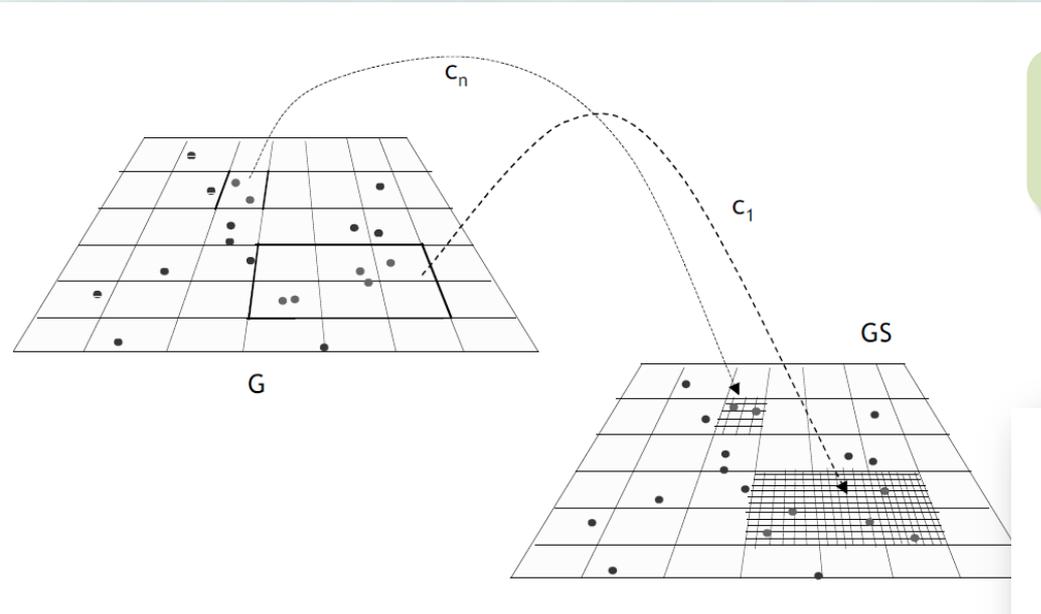
COLING

utile en gestion de la négation
dans l'analyse thématique

Vecteurs d'idées

pliage et dépliage

Comment étendre un espace à un domaine de spécialité ?

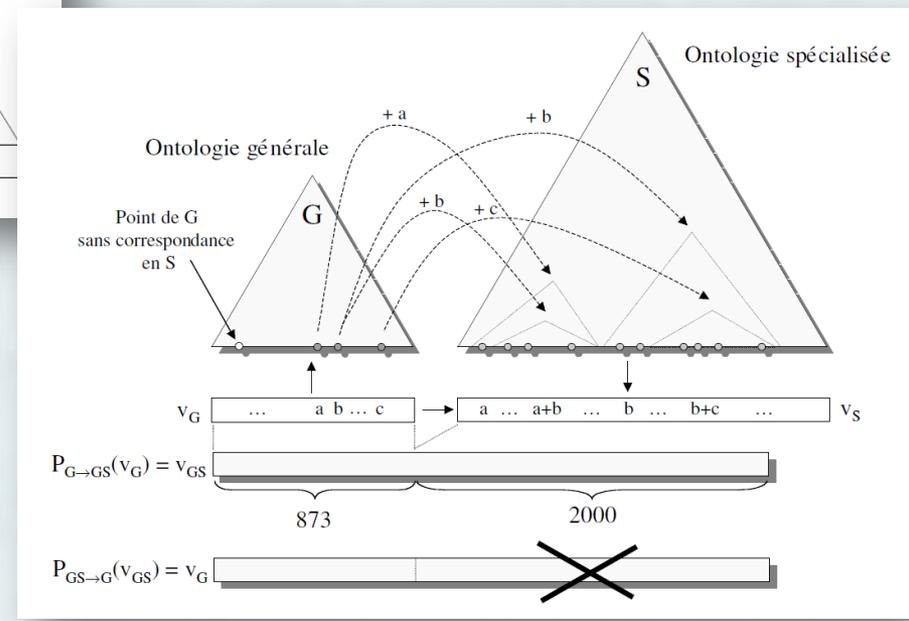


séparation dans S
 si $X, Y \in GS : D_{A,G}(X, Y) > D_{A,GS}(X, Y)$

D. Schwab

V. Prince

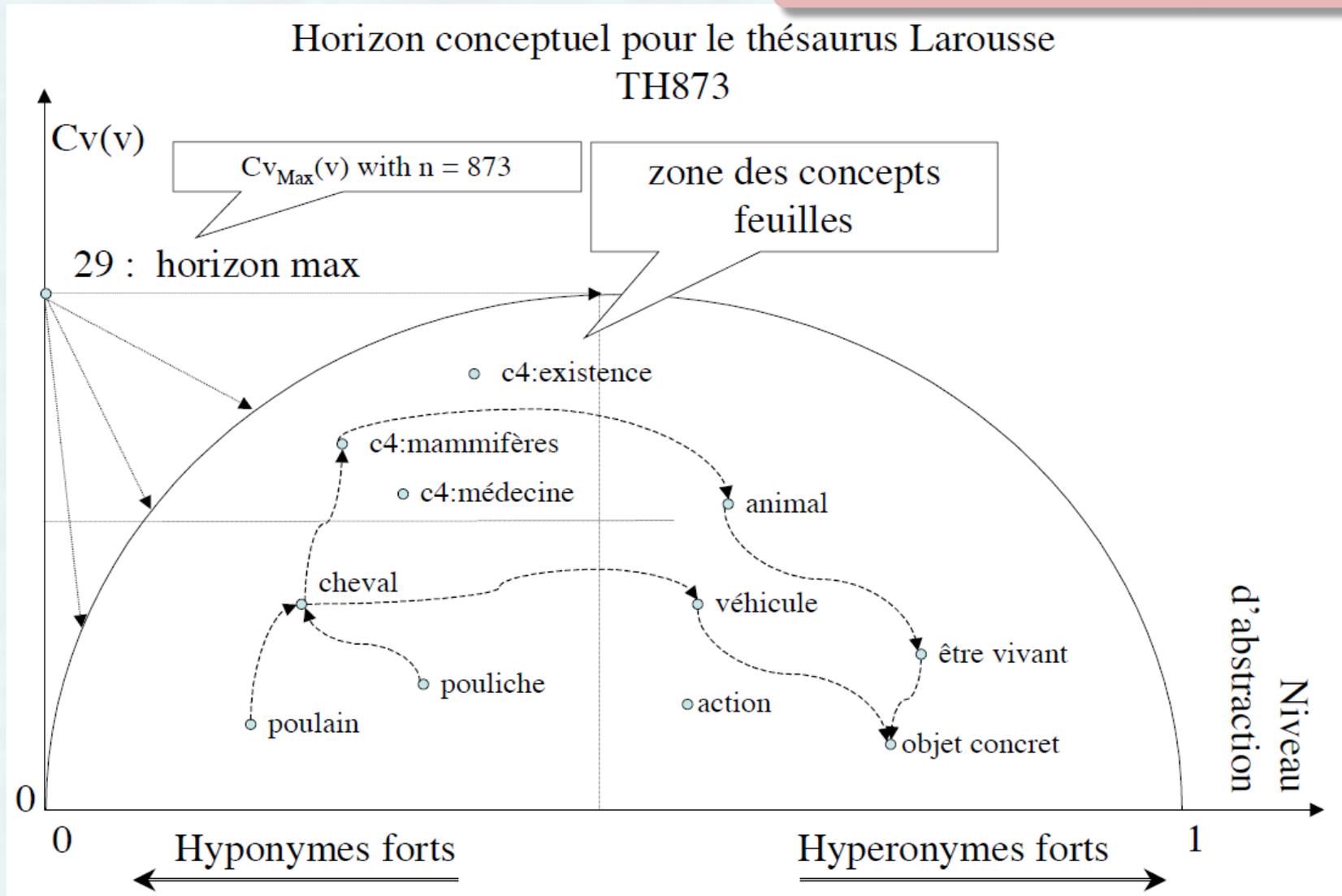
myopie hors S ?
 si $X, Y \notin GS : D_{A,G}(X, Y) < D_{A,GS}(X, Y)$



Vecteurs d'idées

horizon conceptuel

une limite indépassable ?



Vecteurs d'idées - autres formes

vecteurs aléatoires/émergents

vecteurs lexicaux

Choix d'une taille de vecteurs

Placement initial aléatoire

LREC

Random indexing - Sahlgren

Fonction d'agglomération

⇒ analyse de définitions

⇒ moyenne pondérée de termes

Fonction de répulsion nécessaire
car pas d'ancrage

Bonne répartition sur l'espace

Ajustable avec recalcul

Composantes non décodables

Compilation locale d'un réseau lexical

⇒ liste de mots pondérée

≈ Salton ?

Forte précision mais rappel moindre

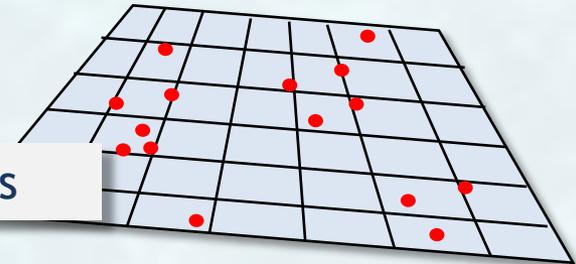
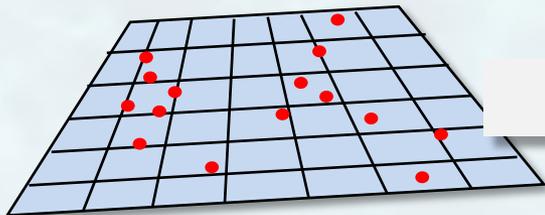
Ajustable sans recalcul

Fortement décodables mais
déconceptualisés

⇒ chat, matou, minet

invariance globale des voisinages

J. Chauché – TAL89



Chronologie



Chronologie

Vecteurs d'idées

Réseaux lexico-sémantiques

2004

D. Schwab

JeuxDeMots

2007

PtiClic

TER L2, LM1

TER L2, LM1

TER L3, LM1

2009

V. Zampa

AKI

2010

M. Zock

Inventaire de sens

A. Joubert

Pour quoi faire ?

Utiles pour des applications nécessitant
des informations terminologiques
des représentations du monde
des informations linguistiques

... et en TAL : Analyse de textes, TAO, résumé, RI, ToT, correction
vers la compréhension automatique du discours, de modèles, etc.
⇒ Analyse sémantique

Comment ?

À la main...

Automatiquement (extraction à partir de textes)

mais les 'savoirs' ne sont pas toujours dans les textes

Ni exclusivement

Ni totalement

Réseaux lexico-sémantiques

toutes les relations potentiellement utiles

Associations libres : chat \Rightarrow animal, chien, souris, ...

Relations ontologiques : hyperonymes – hyponymes – partie de – tout – substance/matière ...

Relations lexicales : synonymie – contraire – locutions – même famille lexicale – magn/antimagn ...

Relations de typicalité/rôles sémantiques : agent – patient – instrument – lieu – rôle téléique (fonction/but) – rôle agentif (mode de création)...

Usages : plus général qu'acception/sens

sapin (arbre), **sapin (sapin de Noël)**,
sapin (bois), sapin(cercueil), sapin (fiacre)

Réseaux lexico-sémantiques

autres réseaux existants

Quelques exemples

WordNet

(25 ans) Miller et al.- *The 2006 version contains 155,287 words organized in 117,659 synsets for a total of 206,941 word-sense pairs.*

coût = 6 M\$

HowNet

(début fin des années 80) Dong et Dong – Anglais-Chinois – environ 100k termes dans chaque langue – 2600 concepts

Wolf

(2006) Sagot et al. – Adaptation au **Français** de WordNet (et EuroWN)

BabelNet et WordNet++

(2010) Navigli et al. – croisement automatique de WordNet et Wikipedia

Construction manuelle et/ou par extraction
+ vérification manuelle

Acquisition de relations

par consensus populaire – jeux sérieux



Peut-on construire un tel réseau autrement ?

Avec des jeux proposés à des non-spécialistes

Hypothèse : rapide – gratuit – efficace

Von Ahn

TAL

TALN

Expérience de faisabilité : projet JeuxDeMots

Principe : demander à des joueurs de renseigner une fonction f pour un terme x

Rendre la tâche

amusante en permettant la comparaison avec les propositions d'un autre joueur

communautaire /compétitive par le gain de points – classement - forum

appropriable par les joueurs avec des possibilités de contrôle partiel sur le jeu

auto-évaluation à l'aune des autres joueurs – satisfaction/frustration

JeuxDeMots - exemple de partie

DONNER DES IDEES ASSOCIEES AU TERME QUI SUIT :

kaput
Crédits : 934050
Honneur : 140849
870



Pearl Harbor

30s

OK



< flotte (navire)
île
raid aérien
7 décembre
base navale
Japon
attaque

7/13



Dernier terme proposé : **flotte**

Raffinements possibles :

1. **flotte** (flotte aérienne)
2. **flotte** (eau)
3. **flotte** (navire)

Ce terme a plusieurs sens ou il en manque ? [Demandez](#) de l'aide à vos amis

JeuxDeMots

exemple de partie : résultat

DONNER DES IDEES ASSOCIEES AU TERME QUI SUIT :

Pearl Harbor

Réponses données par kaput : flotte (navire) - île - raid aérien - 7 décembre - base navale - Japon - attaque

Réponses données par azrael : États-Unis - Oahu - flotte - deuxième guerre mondiale - île - flotte (navire) - Japon - Honolulu - américain - guerre - USA - attaque - base navale

flotte (navire) 🌟 - île - base navale - Japon - attaque

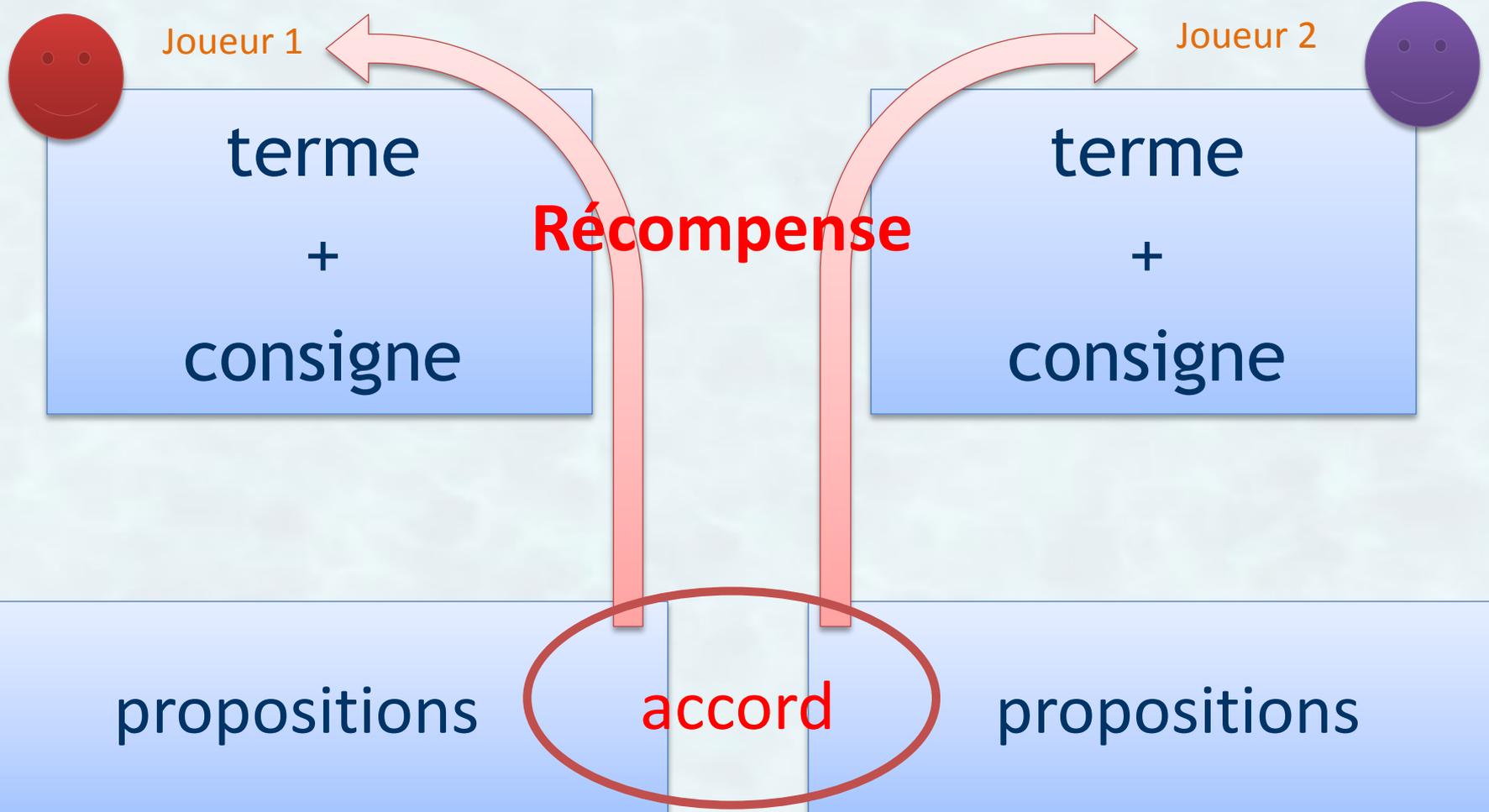
Vous gagnez 376 crédits et 8 point(s) d'honneur



👍 Soyez le premier de vos amis à indiquer que vous aimez ça.

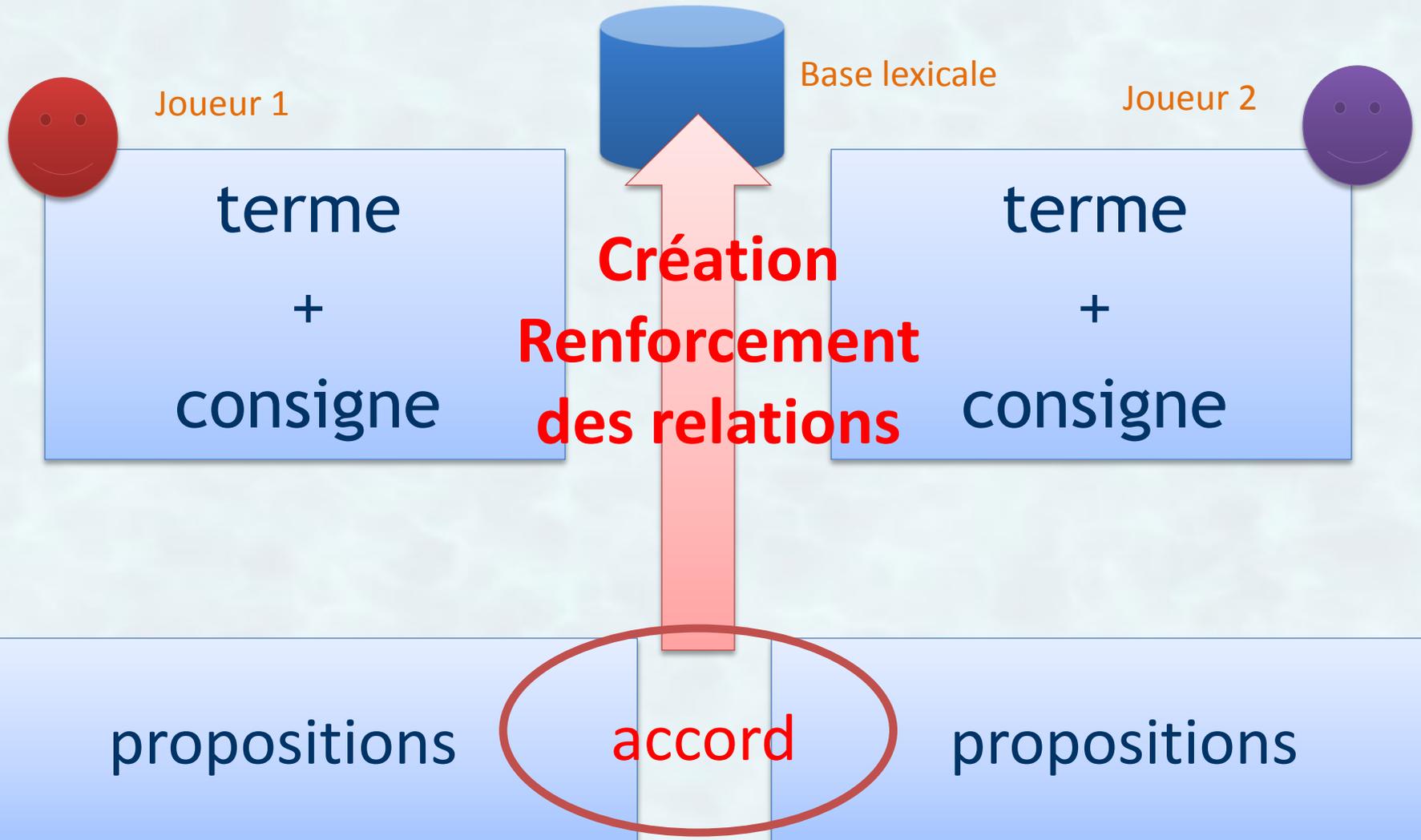
JeuxDeMots

modèle d'interaction



JeuxDeMots

modèle d'interaction



Réseau lexical

modèle d'interaction – filtrage

Accord par paires

Consensus minimal itéré

Minimiser le bruit, maximiser le rappel

(pondération)

(longue traîne)

Caractéristiques

Mot tiré au hasard

Partenaire inconnu durant la partie

Partie asynchrone

Points

⇒ calibrés pour développer la longue traîne (augmenter le rappel)

Bcp

si relation peu
pondérée

peu

si relation très
pondérée

rien

si relation taboue

JeuxDeMots

résultats quantitatifs

Depuis septembre 2007

Plus de 2500 inscrits

Environ 1 million de parties jouées - 16574h de temps de jeu cumulé (690 j)

1 250 000 relations entre 230 000 termes

(amorçage avec 0 relation et 150 000 termes)

632835 r_associated

83349 r_domain

152463 r_syn

41674 r_isa

12111 r_anto

10098 r_hypo

12158 r_has_part

9322 r_holo

9346 r_locution

11306 r_agent

...

coût = 1 mcf/année

Ressource accessible librement

sortie complète mensuelle (pseudo rdf)

sortie terme à terme (pseudo XML)

Des JeuxDeMots dans d'autres langues

English (anglais)

ภาษาไทย (thaï)

ភាសាខ្មែរ (khmer)

العربية (arabe)

日本語 (japonais)

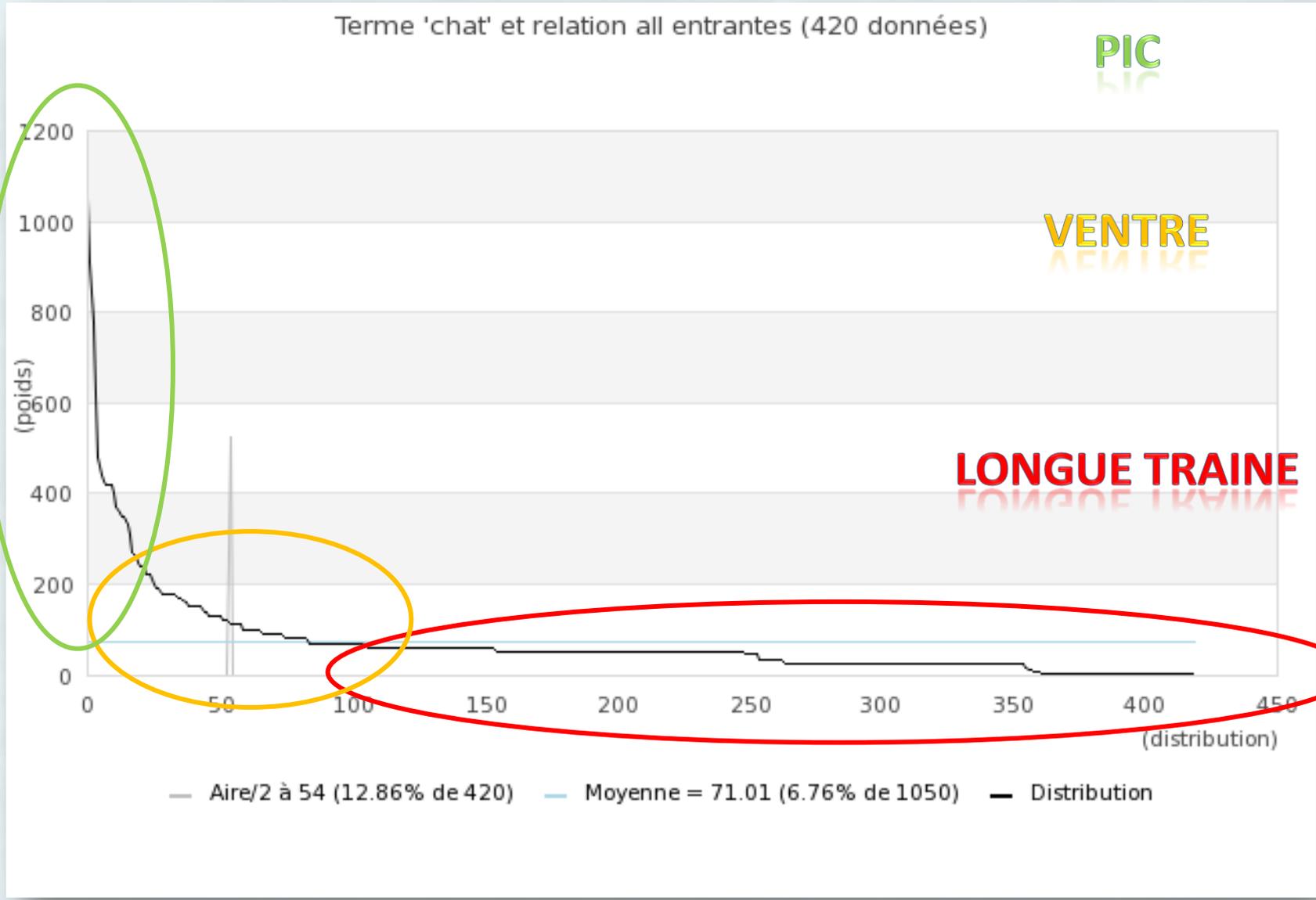
português (portugais)

español (espagnol)

shikomori (comorien)

Tiếng Việt (vietnamien)

Réseau lexical – distribution pour *chat* - zipfien



Réseau lexical

autres jeux

Jeu de renforcement / raffinement des relations

prochaine
longue
nuit
calendrier
jours
période
mars
faire le pont
chargée
vendredi
semaine
week-end
jour
lundi
mardi
mercredi
365 jours
dimanche
année
durée
temps

... est une partie de 'semaine'

Une caractéristique de 'semaine' est ...

'semaine' est une sorte de ...

J'ai fini !

PtiClic

V. Zampa

STICEF

raffinement des structures

Réseau lexical

autres jeux

Jeu de propagation sur les usages
et création de relations négatives / inhibition

prochaine
longue
calendrier
mars
vendredi
jour
lundi
année
durée

Est-ce que

homme (mâle)

peut

travailler

OUI

NON

... est une partie de 'semaine'

Askit

activation et inhibition

Réseau lexical

identification d'usages

A. Joubert

TAL

JADT

Peut-on construire un inventaire de sens / usages ?

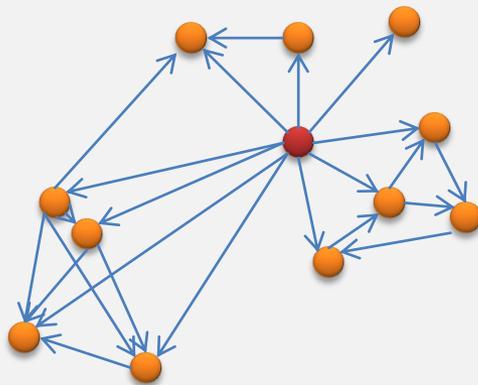
Identification de cliques et quasi-cliques

Classification hiérarchique ascendante

Nommage/glosage descendant

S . Ploux & B Victorri -
synonymes

● chat



Réseau lexical

identification d'usages

A. Joubert

TAL

JADT

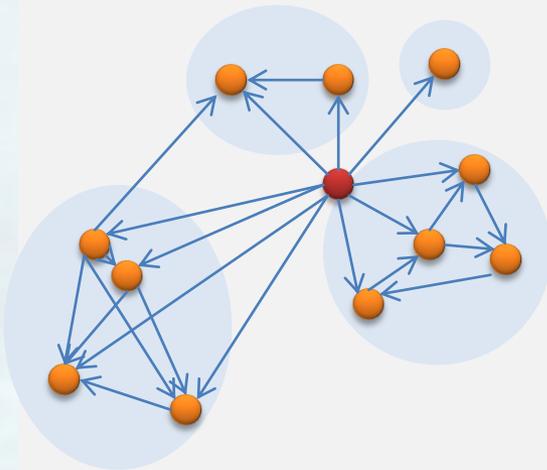
Peut-on construire un inventaire de sens / usages ?

Identification de cliques et quasi-cliques

Classification hiérarchique ascendante

Nommage/glosage descendant

S . Ploux & B Victorri -
synonymes



Réseau lexical

identification d'usages

A. Joubert

TAL

JADT

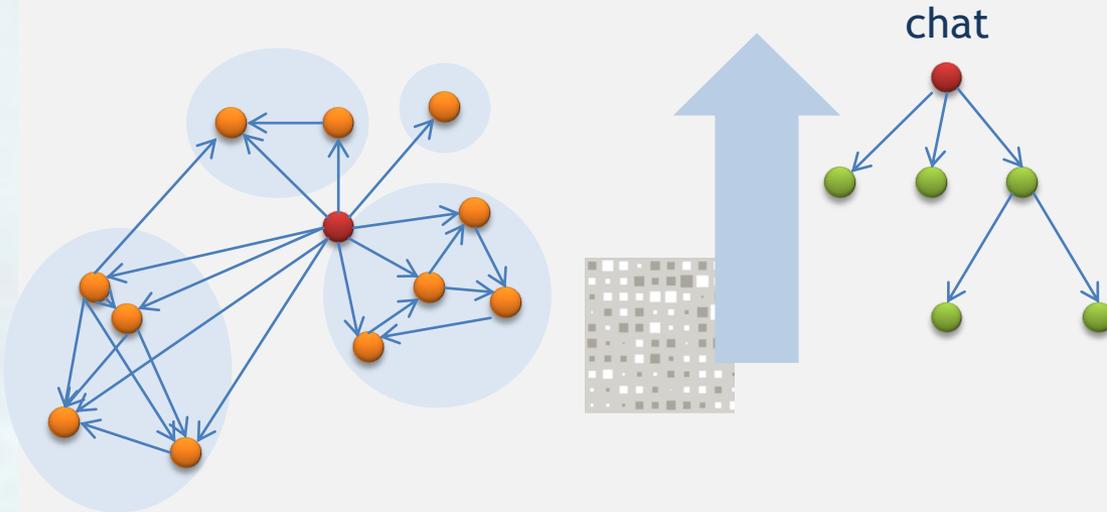
Peut-on construire un inventaire de sens / usages ?

Identification de cliques et quasi-cliques

Classification hiérarchique ascendante

Nommage/glosage descendant

S. Ploux & B. Victorri -
synonymes



Réseau lexical

identification d'usages

A. Joubert

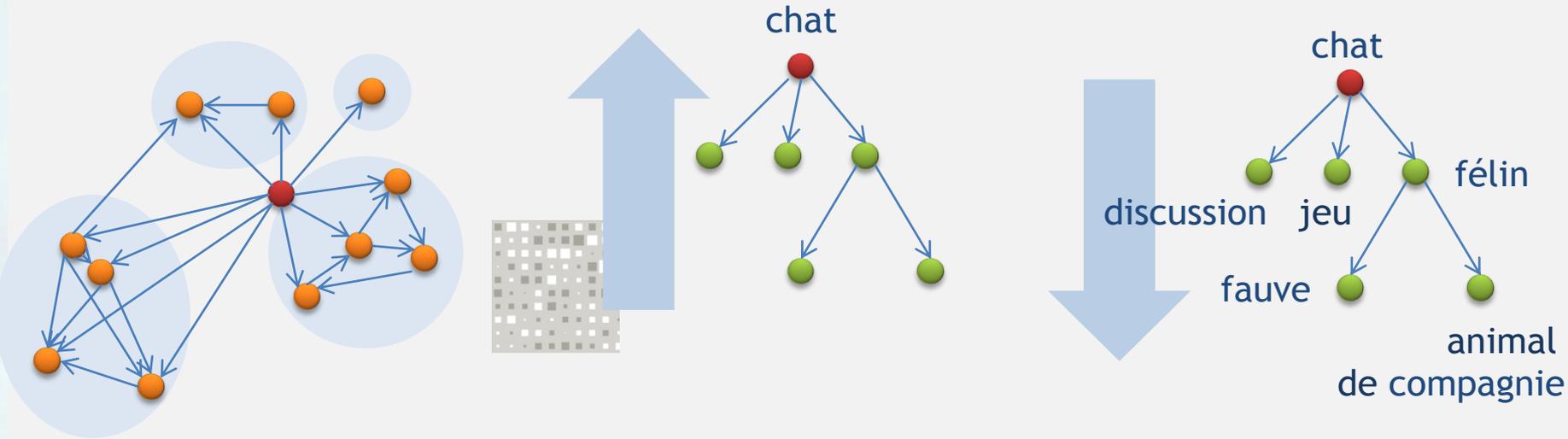
TAL

JADT

Peut-on construire un inventaire de sens / usages ?

Identification de cliques et quasi-cliques
Classification hiérarchique ascendante
Nommage/glosage descendant

S . Ploux & B Victorri -
synonymes



Injection dans le réseau lexical
raffinements utilisables / jouables par les joueurs

bouclage

AKI – un jeu de devinette : faire trouver un terme à partir d'indices

Variantes avec plusieurs groupes de TER

Soumettre l'indice

Après 3 indices, il s'agit sûrement de :

maladie de Parkinson

C'est la bonne réponse !
si ce n'est pas ça vous pouvez proposer un nouvel indice...

Vos indices	Mes propositions
maladie vieillesse tremblement	grippe mort maladie de Parkinson

Après 4 indices, il s'agit sûrement de :

cothurne

C'est la bonne réponse !
si ce n'est pas ça vous pouvez proposer un nouvel indice...

Vos indices	Mes propositions
chaussure théâtre antiquité montante	ped clown lieu cothurne

Après 4 indice(s), je suis perdu, désolé

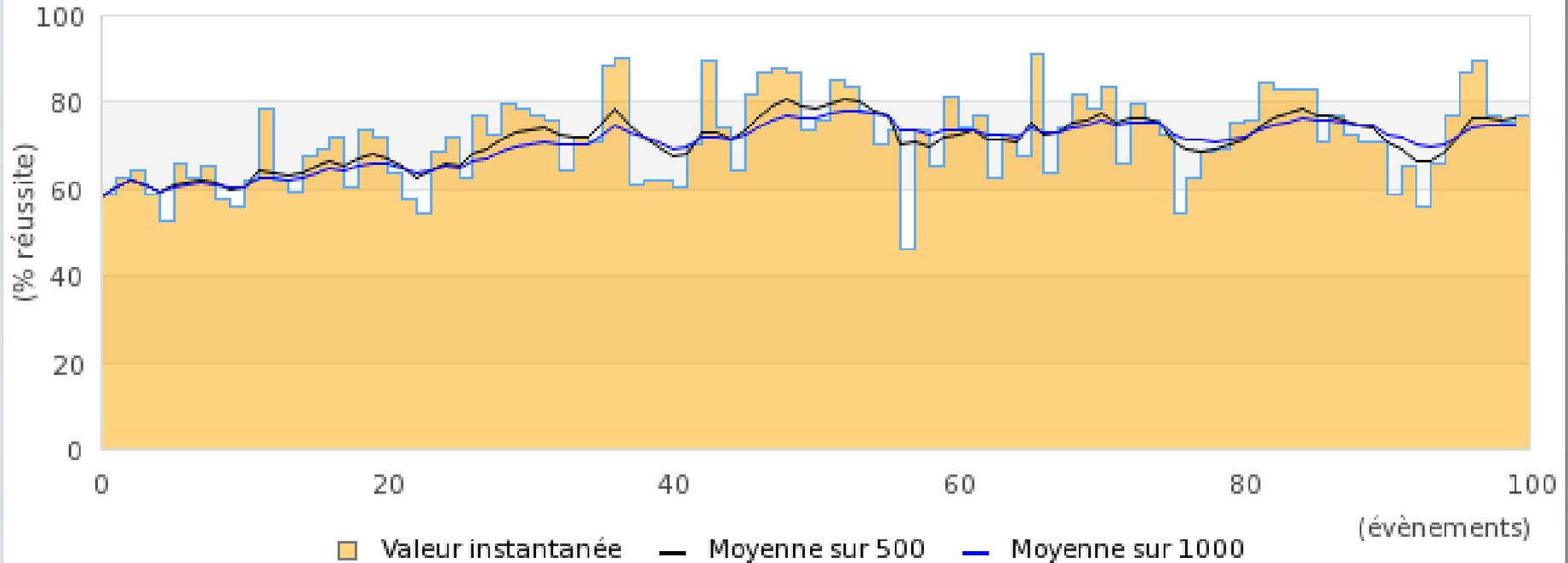
Il s'agissait de

(faites attention aux accents et aux majuscules/n)

Vos indices	Mes propositions
chemin de fer voie de garage wagon wagon de marchandise	rail sous-station électrique transport ferroviaire

Hypothèse : si le terme est trouvé alors il est bien indexé

Données AKI (12103 données - segments de taille 121)



TALN 2011

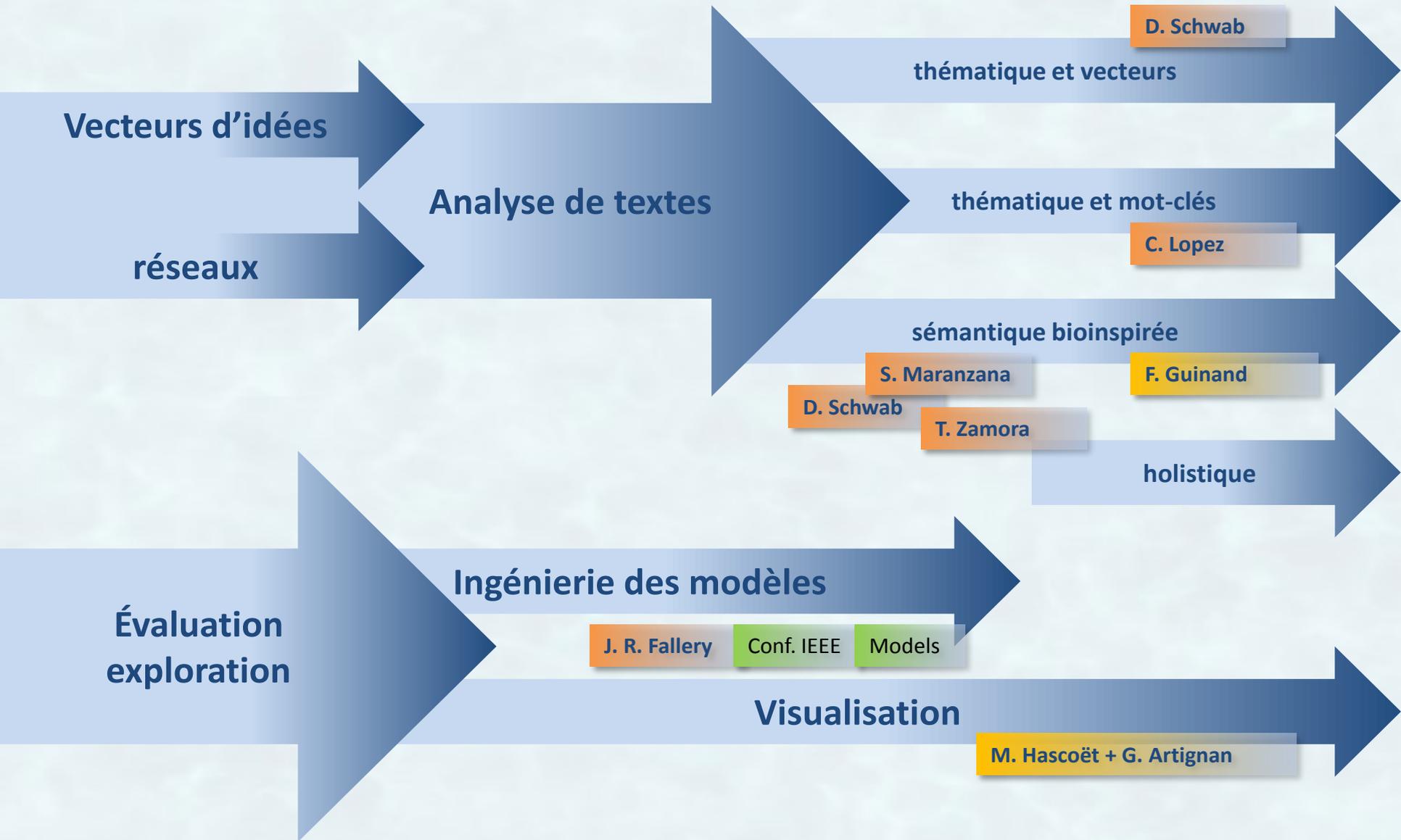
AKI ~ 75%

Test sur > 10000 parties dont les termes sont choisis par les utilisateurs

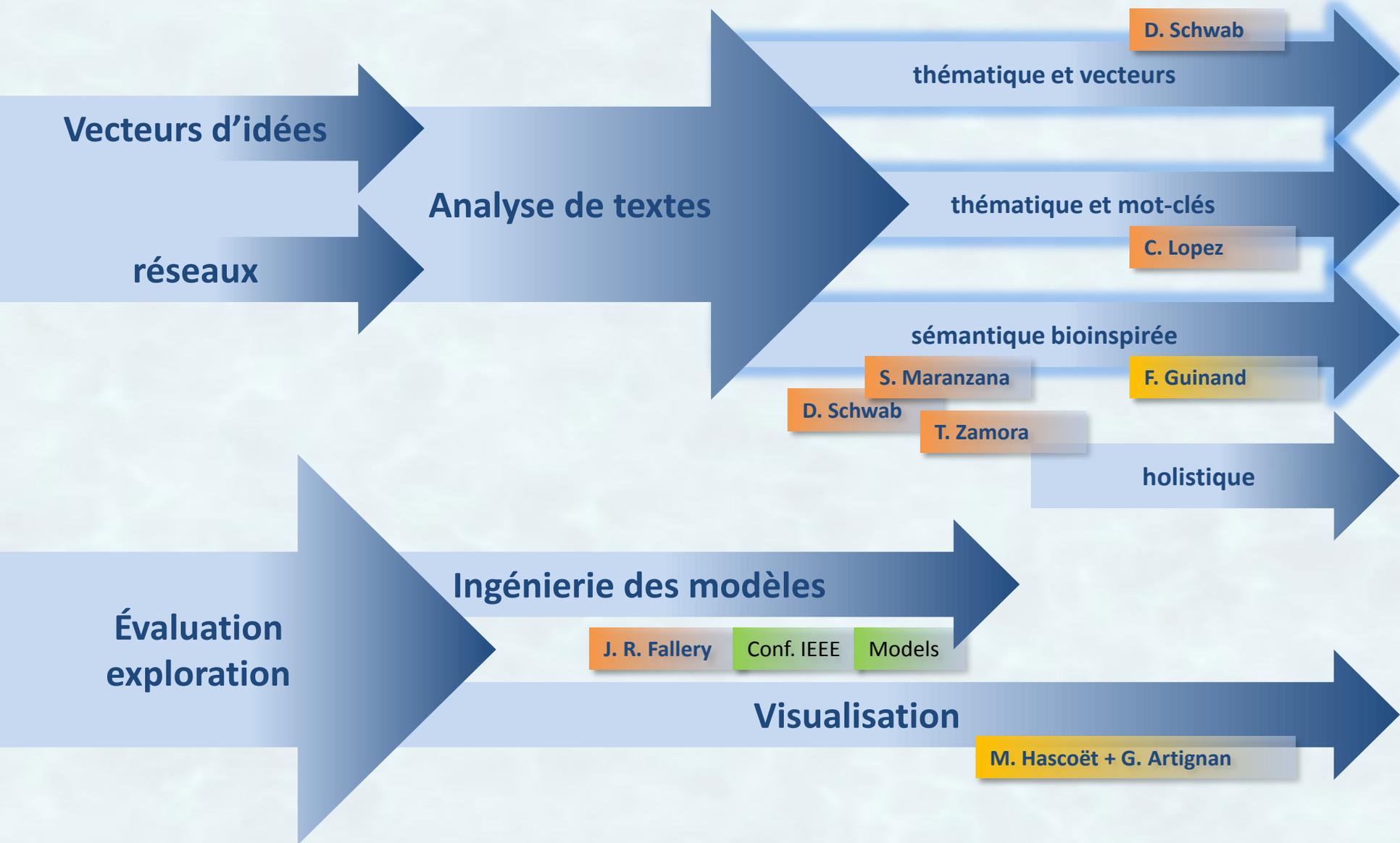
Utilisateur ~ 46%

Test sur environ 200 entrées et les termes tabous de AKI donnés comme indices

Chronologie



Chronologie



Analyse de texte

thématique et vectorielle - propagation

Construire des vecteurs
pour des textes et des définitions ?

Support : arbre d'analyse morpho-syntaxique

⇒ Sygfran

⇒ syntagmes

⇒ gouverneurs/têtes

⇒ fonctions syntaxiques



différence de pondération des vecteurs pour \oplus

J. Chauché

Principe : propagation de vecteurs

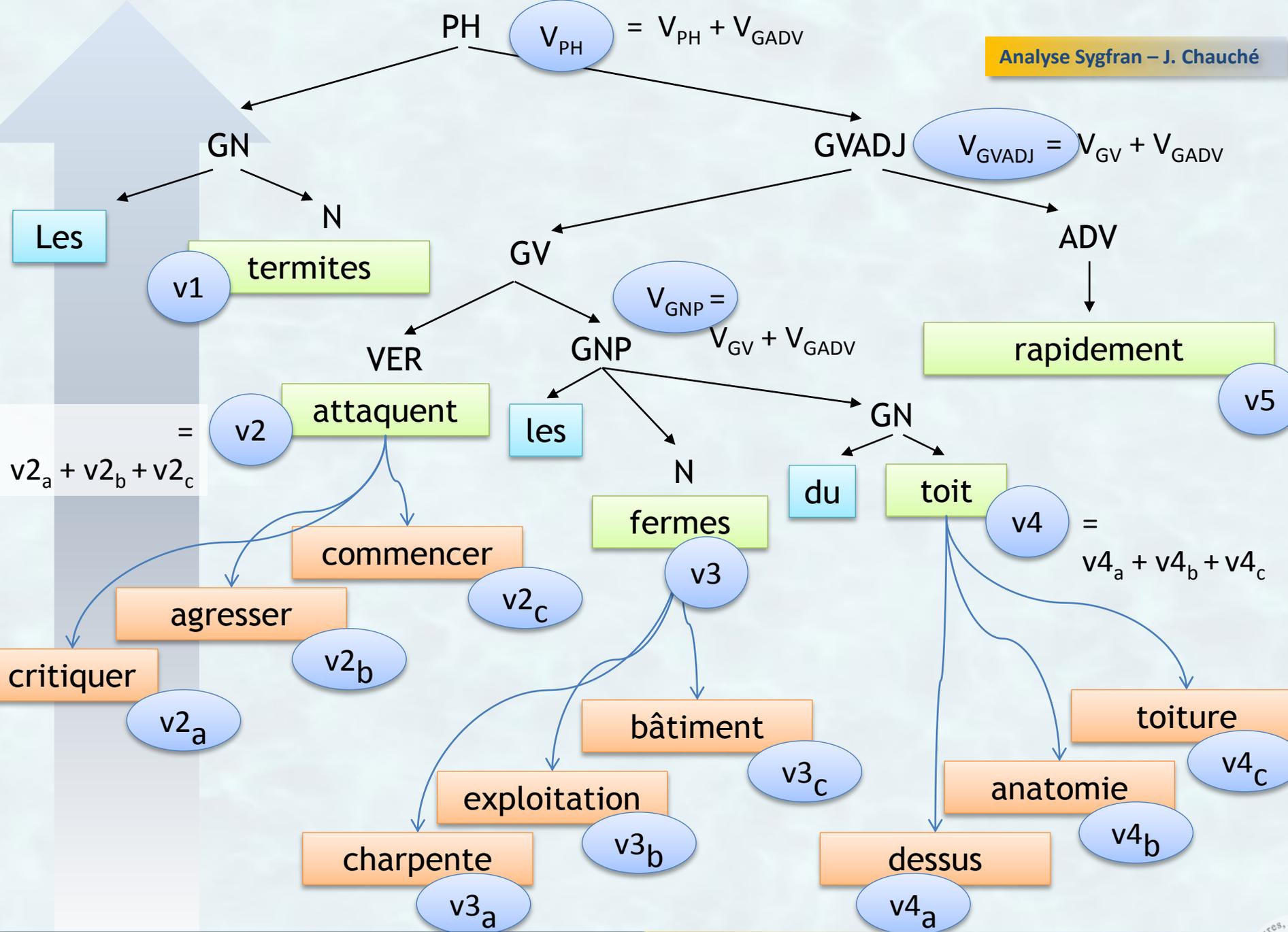
D. Schwab

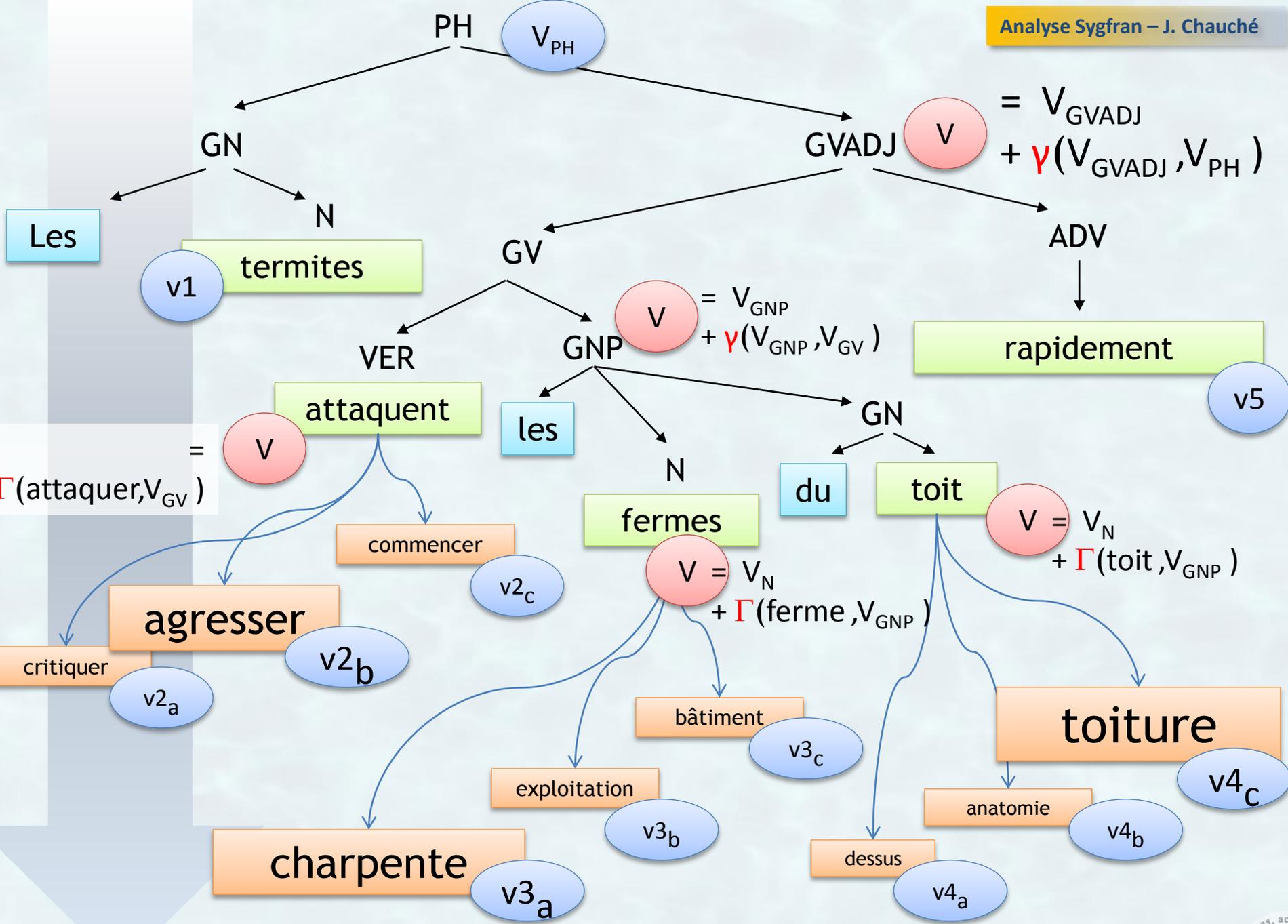
TAL

itération (3-4 cycles de remontée/descente suffisent en général)

convergence globale (vecteur de la racine)

contextualisations faible γ et forte Γ





Analyse de texte

thématique et lexicale - diffusion

C. Lopez

Extraire les termes *pertinents* d'un texte

Amorçage d'un noyau

mot-clés centraux

par TF-IDF

Itération

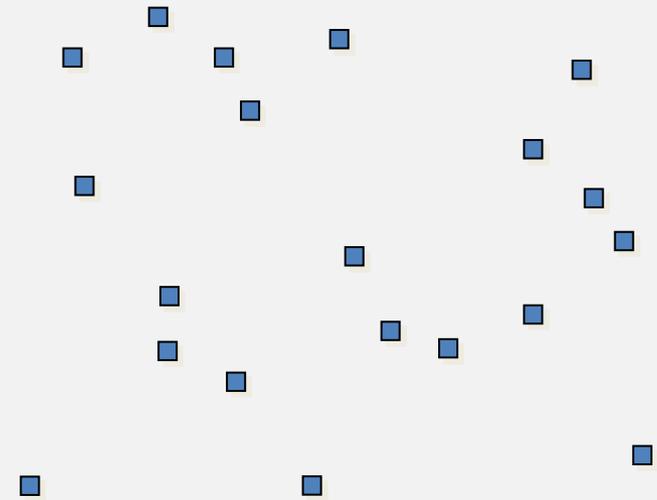
mot-clés périphériques

réduction de la distance seuil

Capture de mot-clés proches dans le réseau

proche au sens de DA sur le vecteur lexical construit

Exemple : *carambolage sur l'A7*



Analyse de texte

thématique et lexicale - diffusion

C. Lopez

Extraire les termes *pertinents* d'un texte

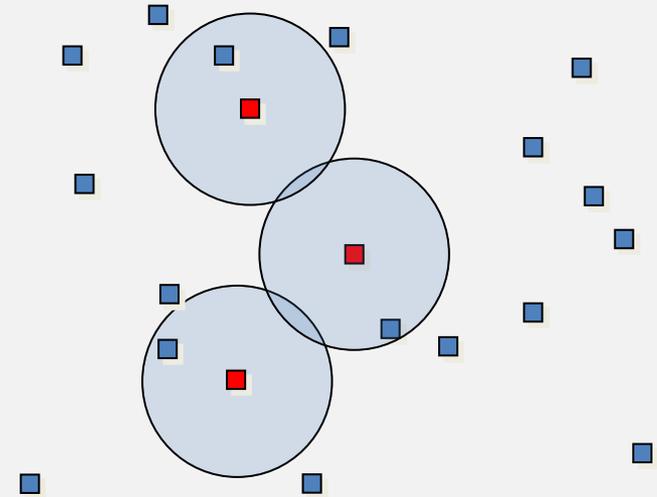
Amorçage d'un noyau
mot-clés centraux
par TF-IDF

Itération

mot-clés périphériques
réduction de la distance seuil

Capture de mot-clés proches dans le réseau
proche au sens de DA sur le vecteur lexical construit

Exemple : *carambolage sur l'A7*



Analyse de texte

thématique et lexicale - diffusion

C. Lopez

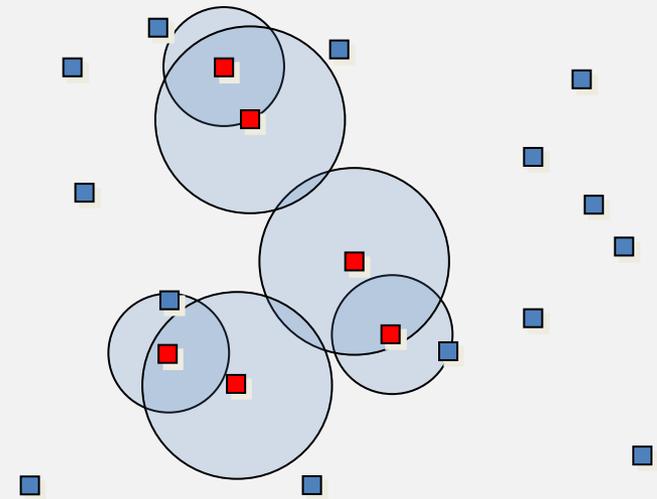
Extraire les termes *pertinents* d'un texte

Amorçage d'un noyau
mot-clés centraux
par TF-IDF

Itération
mot-clés périphériques
réduction de la distance seuil

Capture de mot-clés proches dans le réseau
proche au sens de DA sur le vecteur lexical construit

Exemple : *carambolage sur l'A7*



Analyse de texte

thématique et lexicale - diffusion

C. Lopez

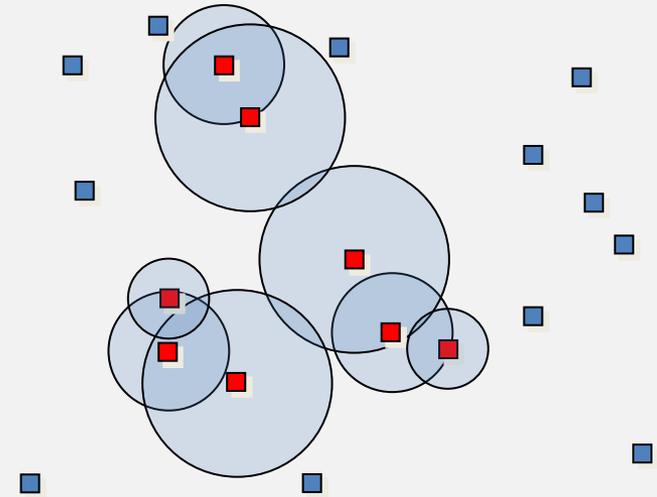
Extraire les termes *pertinents* d'un texte

Amorçage d'un noyau
mot-clés centraux
par TF-IDF

Itération
mot-clés périphériques
réduction de la distance seuil

Capture de mot-clés proches dans le réseau
proche au sens de DA sur le vecteur lexical construit

Exemple : *carambolage sur l'A7*



Analyse de texte

thématique et lexicale - diffusion

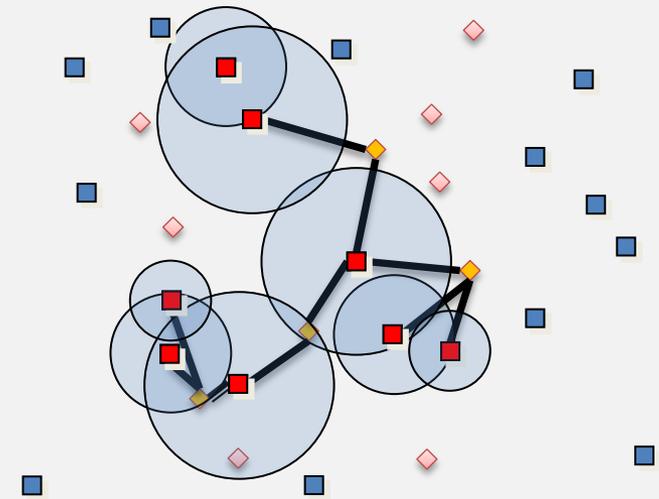
C. Lopez

Extraire les termes *pertinents* d'un texte

Amorçage d'un noyau
mot-clés centraux
par TF-IDF

Itération
mot-clés périphériques
réduction de la distance seuil

Capture de mot-clés proches dans le réseau
proche au sens de DA sur le vecteur lexical construit



Exemple : *carambolage sur l'A7*

carambolage, A7, automobile, transport par route, accident de la route, autoroute, voiture > automobile

Diffusion dans un espace

Propagation dans le réseau

Analyse de texte

sémantique, relationnelle et bioinspirée

Comment être moins localiste ?
Expliciter l'interprétation ?

Multi-agents réactifs : transporteurs d'informations

Support : réseau (arbres, graphes UNL)

Déplacements pseudo-aléatoires

stygmérie – le contrôle par le marquage temporaire de l'environnement

Création de ponts : court-circuits

typés pour rendre explicite la (les) relation(s) entre les segments

Utilisation (recopie locale) du réseau lexico-sémantique

Rattachement de groupes prépositionnels

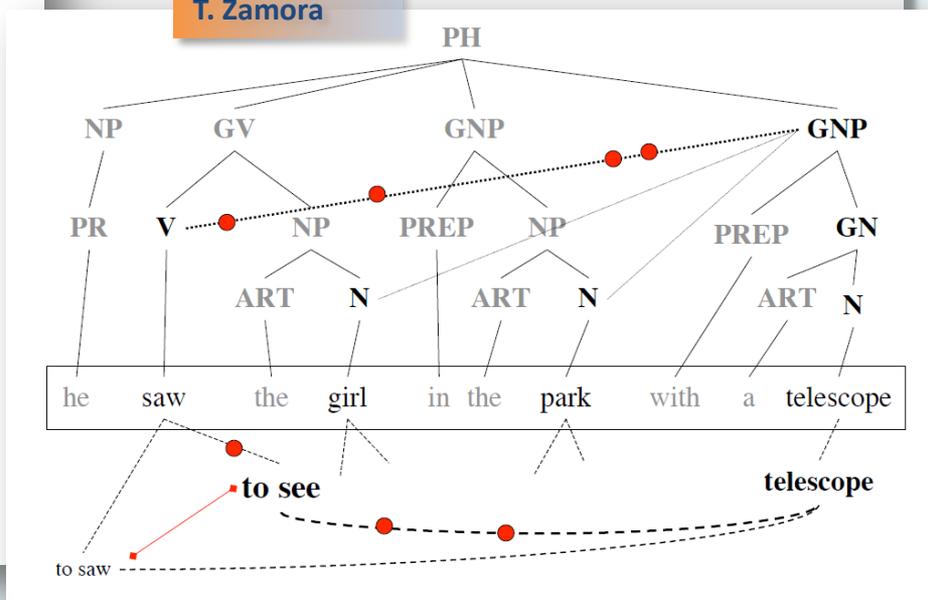
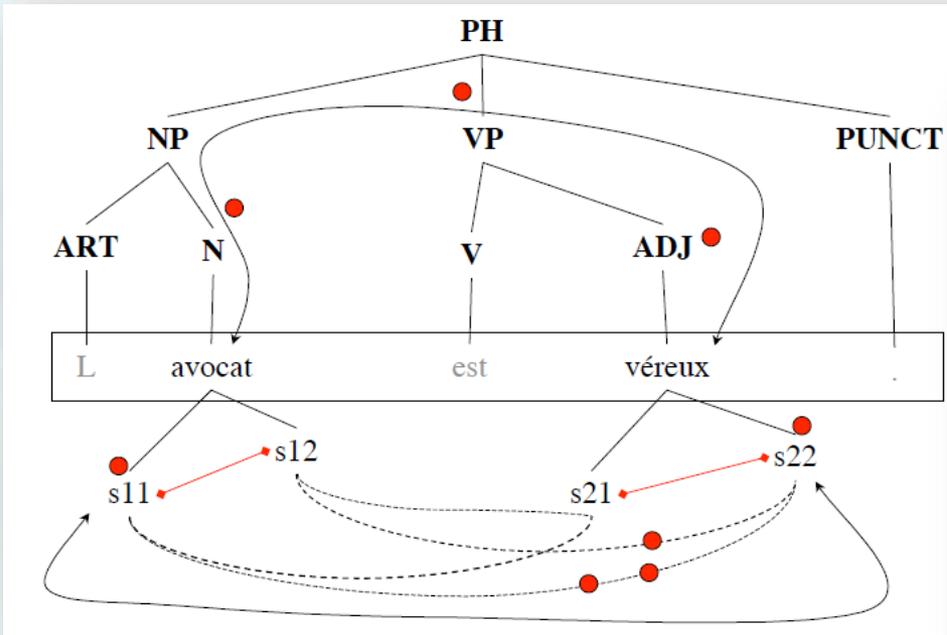
LREC + chap. livre Ch. Boitet

2 chap. de livre F. Guinand

ECTI N. Gala

S. Maranzana

T. Zamora



Conclusions - apports

Analyse de textes

Thématique et sémantique
abordée de façon généraliste

Construction de ressources lexico-sémantiques

Vecteurs et relations : complémentaires
rappel (concepts) et précision (vocables)
Acquisition par consensus populaire – vers du *sens commun*
les rendre disponibles et contribuables

Calcul

Propagation (et diffusion) - Bouclage
Activation (et inhibition)
mise en concurrence d'objets lexicaux : sélection d'usages en contexte

Application et évaluation

Sortir du TAL (par exemple, ingénierie des modèles)
Exploration des ressources
Oracle lexical (AKI) : un outil ludique

Valorisation :

Lingua & Machina
Namae Concepts
Succeed Together
Cemagref

...

M. Huchard

B. Gaume

M. Hascoët

M. Zock

Conclusions - perspectives

Vers des mécanismes d'inférence pour les réseaux sémantiques

Endogène & Exogène (le contenu des textes comme amorçage)

Renforcer / créer certaines relations

Découvrir des règles d'inférences

Renforcer le couplage avec l'analyse

Objets lexicaux agrégés & meta information

(embarcation \Rightarrow **carac** \Rightarrow surchargée) \Rightarrow **conséquence** \Rightarrow chavirer

Réifier les occurrences de relations

$r \Rightarrow$ **info** \Rightarrow raciste, sexiste, humoristique, etc.

$r \Rightarrow$ **info** \Rightarrow tjrs vrai, probable, possible, déduite

Le commis / le vagabond / l'enfant

a pris

une tarte.



Vers du raisonnement et de l'explication

Restaurer au TAL son
désir d'intelligence

Merci