

Acquisition et validation de connaissances à partir de sorties syntaxiques

Éric de la Clergerie

INRIA Paris-Rocquencourt
<http://alpage.inria.fr>



Séminaire TEXTE
Montpellier, 28 Mai 2015

- Développement de **FRMG**, une grammaire TAG à large couverture
 - ▶ s'appuyant sur une **méta-grammaire** [Candito, Crabbé]
 - ▶ initiée en 2004 sur quelques mois pour la campagne EASy
 - ▶ largement étendue et améliorée depuis
- Conçue pour valider analyse syntaxique par charte pour les TAGs

2 grandes questions:

- ❶ À quoi sert un analyseur syntaxique ?
 - ▶ Extraction d'information (<http://passage.inria.fr/SAPIENS>)
 - ▶ Questions-Réponses
 - ▶ **Acquisition de connaissances sur corpus**
 - ▶ ...
- ❷ Comment continuer à améliorer un analyseur syntaxique ?
par **injection de connaissances** (attachement prépositionnel)

- 1 Une brève présentation de FRMG
- 2 Acquisition de connaissances
- 3 Aller plus loin: quelques pistes
- 4 Conclusion

FRMG: issue d'une description modulaire d'une hiérarchie de classes

- exprimant des contraintes sur les noeuds (précédence, dominance, ...) et leurs traits
- pouvant fournir ou requérir des fonctionnalités (*ressource*)

Lors de la phase de compilation:

- Des classes neutres sont obtenues par croisement des classes productrices/consommatrices
- Des arbres TAGs (minimaux) sont produits à partir des contraintes des classes neutres

Modularité \Rightarrow facilite l'ajout ou l'enrichissement de phénomènes syntaxiques

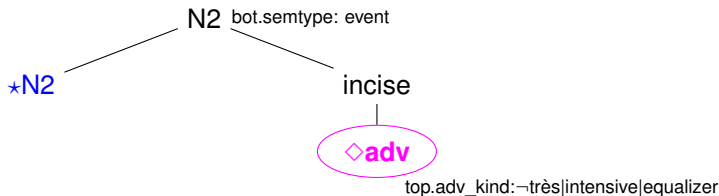
- gestion des arguments verbaux
- *qu'est-ce qu'une pomme ?* : ajout d'une classe, enrichissant les clivées
- gestion des "adverbiaux" flottants
type de flottants \times type d'incise \times positions \times nature modifiés

Exemple: noms événementiels

Accroche de modifieurs flottants (adverbiaux) après les noms événementiels:

*L'annonce, **ce matin**, d'un remaniment a surpris tous les commentateurs*

```
class mod_on_N2 {  
  + s_modifier;  
  - x_modifier;  
  - s_adj_pos;  
  node(Root).cat = value(N2);  
  node(Root).bot.semtype = value(event); %% semantic property  
}
```



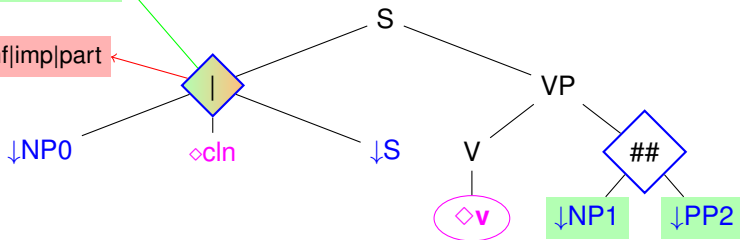
Note: Comment renseigner les noms événementiels !

Arbres factorisés

Le croisement des classes peut normalement produire beaucoup d'arbres
⇒ utilisation d'opérateurs réguliers pour des arbres factorisés

V.top.mode = \neg inf|imp|part

V.top.mode = inf|imp|part



Défactorisation: $(1_{\text{no subj}} + 3_{\text{subj}}) * (1_{\text{no arg}} + 2_{1 \text{ arg}} + 2_{2 \text{ args}}) = 20$ arbres

actuellement (2015) : 368 arbres factorisés

expansion FRMG (2007) ~ 200 arbres $\rightsquigarrow \sim 2$ millions d'arbres !

- Utilisation **SXPIPE** (tokenisation, EN) et **LEFFF** (lexique) **Sagot**
⇒ treillis de mots (DAG)
- analyse par charte, hybride TAG/TIG, relation coin-gauche, filtrage par lexicalisation, ...
- Analyse complète des phrases
si échec, essai de corrections sur la phrase (eg, accord sujet-verbe)
si échec, analyses partielles couvrant la phrase
⇒ très peu de phrases sans analyse (<1% timeout, 5 à 20% partielles)
- Retourne l'ensemble des analyses (totales ou partielles)
sous forme de forêts partagées
- Utilisation de clusters de machines pour traiter de gros corpus

Par défaut, l'analyseur rend une forêt partagée de toutes les dérivations mais nécessaire de pouvoir sortir la ou les meilleures analyses !

Conversion des arbres de dérivations en arbres de dépendances

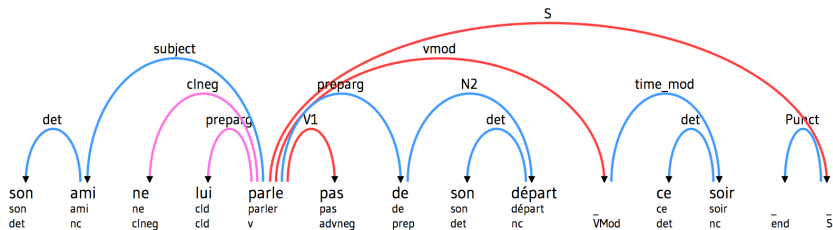
Utilisation de règles de désambiguïisation sur les dépendances pour rechercher la meilleure analyse

- initialement, règles et poids manuellement définis
exemple: favoriser les arguments verbaux, pénaliser inversion sujet, ...
- maintenant, apprentissage des poids appris sur FrenchTreeBank (FTB)

Conversion vers divers schémas d'annotations:

DepXML, **FTB/CONLL** & **PASSAGE** (+ variantes)

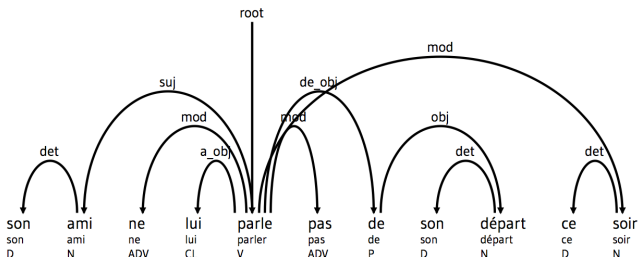
Schéma natif FRMG pour les arbres (et forêts) de dépendances



Derrière la vue graphique, un format XML riche en information

- tokens
- noeuds avec formes, lemmes, catégories et arbre TAG ancré
pseudo-noeuds pour les arbres non lexicalisés
- dépendances entre noeuds (label, type, dérivation)
- constituants “maximaux” (couverts par les arbres) avec traits
- info syntaxiques ([hypertag](#)): sous-catégorisation, contrôle, diathèse, ...

Schéma d'annotation utilisé pour la version en dépendance du French TreeBank (FTB) [Candito & Seddah]

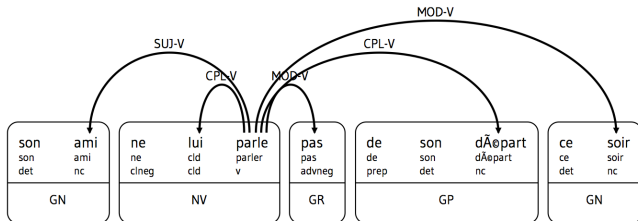


Représentation tabulaire (CONLL)

1	son	son	D	DET	n=s s=poss	2	det
2	ami	ami	N	NC	g=m n=s s=c	5	suj
3	ne	ne	ADV	ADV	s=neg	5	mod
4	lui	lui	CL	CLO	n=s p=3	5	a_obj
5	parle	parler	V	V	g=m m=ind n=s p=3 t=pst	0	root
...							

Un schéma issu des campagnes d'évaluation EASy/Passage avec

- 6 types of *chunks* (GN, NV, GA, GR, GP, PV)
- 14 types de dépendances (SUJ-V, COD-V, CPL-V, ATB-SO, MOD-N, ...)
- format XML



Moins riche en information que DepXML mais “standard”
⇒ utilisé pour les xp d'acquisition

Quelques résultats

Sur FTB et Sequoia (FTB/CONLL, LAS sans ponctuation)
et EASy (PASSAGE, F-mesure sur les relations [90.5% sur les chunks])

système	dev	test	sequoia	easy
FRMG init	80.85	82.08	81.13	65.92
FRMG +tuning	86.76	87.95	86.41	69.89
BKY	86.50	86.80	–	–
MALT	86.90	87.30	–	–
MST	87.50	88.20	–	–
TALISMANE (b=20)	88.10	88.50	–	–
Le Roux (2012)	–	89.20	–	–
DYALOG-SR (b=8)	88.17	89.01	85.02	–
DYALOG-SR+FRMG (b=10)	89.02	90.25	87.14	–

Constat:

- Performances FRMG similaires à celles des parseurs statistiques
- plus stable sur du hors domaine (SEQUOIA)
- meilleur sur phénomènes profonds
(Ribeyre, article RECITAL: contrôle et ellipses)

- 1 Une brève présentation de FRMG
- 2 Acquisition de connaissances
- 3 Aller plus loin: quelques pistes
- 4 Conclusion



*“Beware the Jabberwock, my son!
The jaws that bite, the claws that catch!
Beware the Jubjub bird, and shun
The frumious Bandersnatch!”*

*Il était grilheure; les slictueux toves
Gyraient sur l’alloinde et vriblaient:
Tout flivoreux allaient les borogoves;
Les verchons fourgus bourniflaient.*

Paul s’est cassé la **binti**.
Sa fracture à la **binti** a été correctement réduite.
Il a des douleurs dans la **binti**.

un tremblement de terre de forte magnitude
il observe un homme avec un télescope

Point de départ: des expérimentation menées dans le cadre du projet SCRIBO

- Traitement syntaxique de gros corpus (avec FRMG)
- Exploitation des sorties syntaxiques
- Acquisition de connaissances lexico-sémantiques, dont terminologie, réseau de mots, ...
 - ▶ ré-injection pour la désambiguïsation syntaxique
tremblement de terre de forte magnitude
 - ▶ aide à la construction de glossaires, thésaurus, ontologies lexicalisées (domaines spécialisés)

Obtention de grosses ressources (entre 1K à 100K *entrées*)

⇒ problématiques de visualisation, évaluation, et validation

Mise en place d'une interface de visualisation et validation collaborative

Un large corpus “*général*” hétérogène (CPL+AFP)

Corpus	#Mphrases	#Mmots	Description
Wikipedia	18.0	178.9	504K pages encyclopédie
Wikisource	4.4	64.0	12.8K textes littéraires
EstRepublicain	10.5	144.9	journalistique
JRC	3.5	66.5	directives européennes
EuroParl	1.6	41.5	débats parlementaires
Total CPL	38.0	495.8	
AFP	14.0	248.3	400K dépêches
Total ALL	52.0	744.2	

Mais aussi des corpus spécialisés plus petits (certains en Juridique)

Corpus	#Mphrases	#Mmots
fiscal	7.2	145.2
social	6.8	127.5
civil	2.6	40.9
affaires	7.2	133.8

Et de nombreux autres: botanique (1Mmots), automobile, récits de voyage, ...

Phase0 Analyser les corpus, sur cluster

Phase1: Collecter et compter

- Algorithme **MapReduce** (popularisé par Google et **HADOOP**)
version distribuée de `grep|sort|uniq -c|sort -nr`
- collecter des éléments d'information, étant donné un ou des motifs
- trier les éléments
- et les compter

Phase2: utiliser les décomptes pour prendre des décisions
clustering, classification, ranking, ...

(Phase3): Valider !

(Phase4): Injecter dans les phases 0, 1, ou 2

Extraction Terminologique

Un ensemble plus ou moins standard de critères guidant l'extraction de termes

- **structure** (multi-mots) : les chunks Passage sont bien adaptés
(GN) (GR*GA | GP | PV | NV) +
base: Chunks nominaux GN (et GP/Prep), avec modifieurs syntaxiques
- forte **fréquence**
note: pas (encore) d'information de contraste pour des domaines spécialisés
- forte **cohésion** interne (**information mutuelle**)
- (**nouveau**) **autonomie**: un terme doit apparaître nu
i.e., sans modifieurs, dans des fonctions sujets ou objets, ...
- (**nouveau**) **diversité**: un terme apparaît dans des contextes diversifiés
(évite collocations, phrases dupliquées, ...)
- **variants** : un peu de variation (mais pas trop) sur même lemmes
élimination de certaines entités nommées

↪ plus de 100k termes potentiels (ordonnés) extraits des dépêches AFP
pas de seuil (favorise rappel)
bruit: entités nommées (organisations, événements), VP mal coupés, ...

Quelques exemples

Diversité des motifs (et existence de long termes)

dioxyde de carbone
carbon dioxid

[dioxyde/nc]_{GN} [de/prep carbone/nc]_{GP}

hockey sur glace
ice hockey

[hockey/nc]_{GN} [sur/prep glace/nc]_{GP}

téléphone portable
mobile phone

[téléphone/nc]_{GN} [portable/adj]_{GA}

lait écrémé
skimmed milk

[lait/nc]_{GN} [écrémer/v]_{NV}

permis de conduire
driving license

[permis/nc]_{GN} [de/prep conduire/v]_{PV}

procréation médicalement assistée
medically assisted procreation

[procréation/nc]_{GN} [médicalement/adv]_{GR}
[assisté/adj]_{GA}

implant chirurgical non actif
non active chirurgical implant

[implant/nc]_{GN} [chirurgical/adj]_{GA}
[non/adv]_{GR} [actif/adj]_{GA}

Limites: *filis et filles de*

Collationner et compter les dépendances

<gouverneur>	<rel>	<gouverné>	<freq>
-----	-----	-----	-----
chaise_nc	et	table_nc	235
asseoir_v	sur	chaise_nc	227
chaise_nc	modifieur	long_adj	168
chaise_nc	de=	poste_nc	115
tomber_v	sur	chaise_nc	103
chaise_nc	modifieur	musical_adj	102
se_asseoir_v	sur	chaise_nc	93
prendre_v	cod	chaise_nc	87
chaise_nc	modifieur	électrique_adj	82
chaise_nc	modifieur	vide_adj	80
chaise_nc	à=	porteur_nc	80
dossier_nc	de	chaise_nc	78
avoir_v	cod	chaise_nc	71
table_nc	et	chaise_nc	62
chaise_nc	de=	paille_nc	56

À partir des dépendances Passage, avec un peu de correction et d'abstraction:

- rectification des passifs
- ajout de *se* pour les verbes pronominaux
- relation directe entre un sujet et un attribut
(*pomme rouge* dans *la pomme est rouge*)
- abstraction du verbe dans les arguments sententiels
pouvoir cod sentence \rightsquigarrow *pouvoir cod manger*
- distribution entre éléments coordonnées

- ajout des attachements ambiguës

<code>terre_nc</code>	<code>de=</code>	<code>magnitude_nc</code>	344
<code>tremblement_nc</code>	<code>de=*</code>	<code>magnitude_nc</code>	357

- injection des termes candidats

<code>qualité_nc</code>	<code>de=</code>	<code>président_du_conseil</code>	189
<code>tremblement_de_terre</code>	<code>de=*</code>	<code>magnitude_nc</code>	

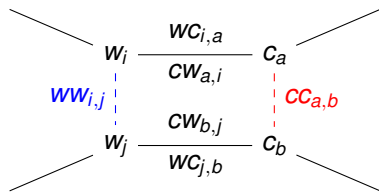
Une dépendance (*se_asseoir*, *sur*, *chaise*) fournit l'association

- d'un **contexte syntaxique** $\langle se_asseoir\ sur\ \bullet \rangle$ avec le **mot** *chaise*
- et, symétriquement, $\langle \bullet\ sur\ chaise \rangle$ avec *se_asseoir*

Corpus	#dep (millions)	#formes distinctes (milliers)	# contextes distincts (millions)
CPL	170	1149	4
AFP	93	378	2
Total ALL	263	1366	5

Inspiré de **Markov clustering** [MCL, **van Dongen**] cherchant, dans un graphe,

- à favoriser les fortes densités de chemins courts
- à pénaliser les chemins longs



$$WW_{i,j} = \frac{1}{Z_i} \left(\sum_{a,b} (WC_{i,a})(CC_{a,b})(WC_{j,b}) \right)^\alpha$$

$$CC_{a,b} = \frac{1}{Z_a} \left(\sum_{i,j} (CW_{a,i})(WW_{i,j})(CW_{b,j}) \right)^\alpha$$

avec **inflation** $\alpha > 1$ (défaut: 2) et normalisation $\frac{1}{Z}$
 \Rightarrow renforce les forts coefficients, réduit les faibles !

Les équations précédentes s'expriment sous forme matricielle compacte:

$$\begin{cases} W = \Gamma_\alpha(F^t C F) \\ C = \Gamma_\alpha(G^t W G) \end{cases}$$

avec l'opérateur d'inflation et normalisation Γ_α

où:

- $W = (w w_{i,j})$ et $C = (c c_{a,b})$ sont les matrices de similarités à calculer
- $F = (w c_{j,a})$ et $G = (c w_{a,i})$ matrices de paramètres
 - ▶ $w c_{j,a}$: importance (pondérée) du contexte c_a pour le mot w_j
 - ▶ $c w_{a,i}$: importance (pondérée) du mot w_i pour le contexte c_a

Formulation récursive, attaquable par un algorithme itératif de point-fixe en démarrant d'une matrice initiale $W^{(0)}$

Variation de tf.idf, avec un ratio de normalisation:

$$w_{c_{i,a}} = \frac{\ln(u_{ai}) * \eta_a}{Z_i} \quad \text{with} \quad \eta_a = \ln \left(\frac{\#\text{mots distincts}}{\sqrt{|\{w_j | u_{aj} > 0\}|}} \right) \quad (1)$$

$$c_{w_{a,i}} = \frac{\ln(u_{ai}) * \eta_i}{Z_a} \quad \text{with} \quad \eta_i = \ln \left(\frac{\#\text{contextes distincts}}{\sqrt{|\{c_b | u_{bi} > 0\}|}} \right) \quad (2)$$

où u_{ai} nombre de co-occurrences de c_a avec w_i

Intuition: si w_1 et w_2 sont similaires, alors le sont aussi les contextes:

• $\langle w_1 r \bullet \rangle$ et $\langle w_2 r \bullet \rangle$

• $\langle \bullet r w_1 \rangle$ et $\langle \bullet r w_2 \rangle$

e.g. **chaise** \sim **tabouret** $\rightsquigarrow \langle \bullet \text{ sur chaise} \rangle \sim \langle \bullet \text{ sur tabouret} \rangle$

Idem des contextes vers les mots

Formalisation par des matrices de transfert:

$$\begin{cases} W = \Gamma_\alpha(F^t C F + \tau(C)) \\ C = \Gamma_\alpha(G^t W G + \rho(W)) \end{cases} \quad (3)$$

avec

$$\tau(C) = \sum_r \beta_r T_r^t C T_r \text{ and } \rho(W) = \sum_r \beta_r T_r W T_r^t \quad (4)$$

et $T_r = (t_{r,ia})_{ia}$ avec $t_{r,ia} = 1$ si $c_a = r.w_i$ et 0 autrement

Par défaut, $\beta = 0.2$

Possibilité d'injecter à chaque itération une matrice L de bonus/malus:

$$W = \Gamma_1((I + L) \circ \Gamma_\alpha(F^t C F + \tau(C))) \quad \text{avec } (A \circ B)_{ij} = A_{ij} \cdot B_{ij}$$

Utilisation:

- forcer l'**auto-similarité** d'un mot avec lui-même
- Coordination (en fonction du nombre d'occurrences)
- Distance d'édition: des mots proches (si assez longs) tendent à être proches sémantiquement
- différences de productivités et fréquences
⇒ émergence d'une hiérarchie sur les mots avec des *hubs*
- **Random indexing** : distance cosinus entre vecteurs de taille 1000
projection des vecteurs de contextes
élargissement *qualitatif* des vecteurs de contextes
- similarités germes, par exemple venant de WordNet
(peu testé pour l'instant)

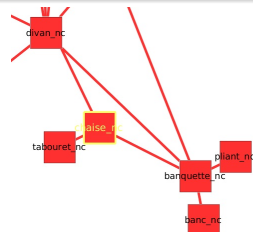
On peut estimer la contribution d'un contexte c_a pour grouper w_i et w_j :

$$\langle \mathbf{w}w_{i,j}, \mathbf{c}_a \rangle = \sum_b (\mathbf{w}c_{i,a})(\mathbf{c}c_{a,b})(\mathbf{w}c_{j,b})$$

⇒ permet de mesurer la corrélation entre les rangs des contextes groupant w_i et w_j , avec ceux pour w_i (et pour w_j):

$$\Delta_{ij} = \sum_{c_a, \mathbf{w}c_{i,a} > 0} \left(1 + \left| 1 - \frac{\text{rank}_{aj}}{\text{rank}_{ai}} \right| \right)^{-1} \left(1 + \left| 1 - \frac{\langle \mathbf{w}w_{i,j}, \mathbf{c}_a \rangle}{\mathbf{w}c_{i,a}} \right| \right)^{-1} \quad (5)$$

À quoi sert une chaise ?



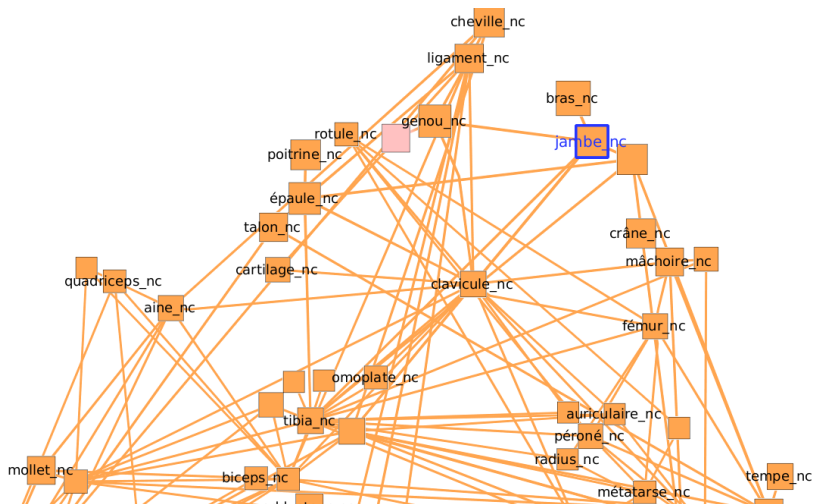
Les regroupements sont motivés par des contextes syntaxiques $\langle ww_{i,j}, C_b \rangle = \sum_a (WC_{i,a})(CC_{a,b})(WC_{j,b})$

	chaise divan	chaise tabouret	banquette divan	banquette canapé	banquette chaise
se asseoir sur [•]	●	●	●	●	●
asseoir sur [•]	●	●	●	●	●
allonger sur [•]	●		●	●	
dormir sur [•]	●		●	●	●
tomber sur [•]	●		●	●	●
monter sur [•]		●			●
place sur [•]					
grimper sur [•]		●			●
installer sur [•]		●			●
poser sur [•]		●			●

Visualisation: que d'os, que d'os !

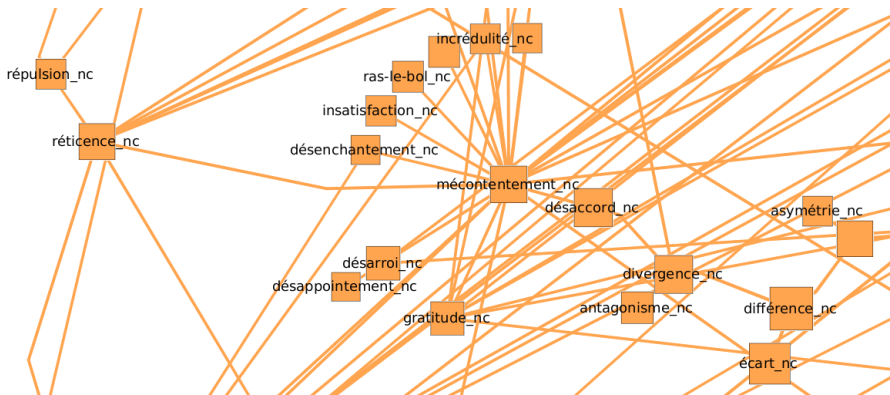
Grphe d'environ 40K connections

Visualisation avec TULIP (<http://tulip.labri.fr/>), layout *BubbleTree*

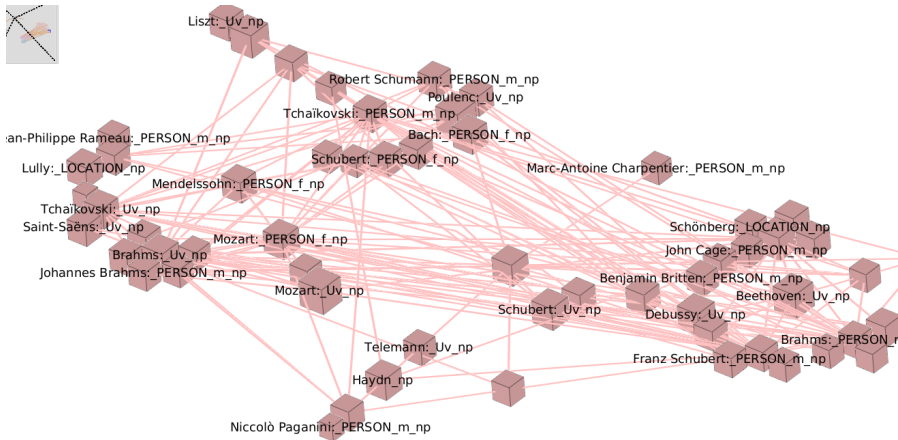


d'autres fragments sur <http://alpage.inria.fr/~clerger/wnet/wnet.html>

Visualisation: sentiments négatifs



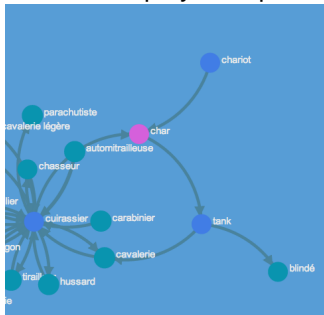
Visualisation: compositeurs



Structures topologiques

Vue à gros grain déjà utile pour repérer certaines structures topologiques:

- **buissons** fortement connectés: très proches de classes sémantiques
- **filaments**: glissement de sens
- **étoiles**: un centre avec de nombreux satellites parfois pertinents, souvent mauvais !
- des mots polysémiques à la jonction de buissons



char et **chariot**

<• *modifieur atteler*>, <*promenade en* •>

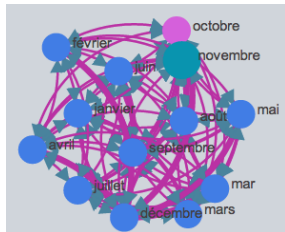
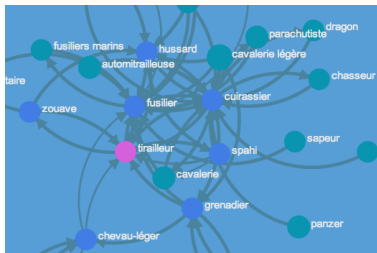
char et **tank**

<• *de combat*>, <*régiment de* •>

Quelques classes topologiques

Les *buissons* permettent l'extraction d'environ 4000 classes (ALL)

- <79> *sulky malinois fox-terrier setter cocker colley chiot fox labrador ratier griffon caniche teckel épagneul*
- <80> *arrière-garde canonnier cavalerie carabinier tirailleur hussard panzer voltigeur blindé grenadier cuirassier avant-garde zouave lancier*
- <83> *pneumonie paludisme diphtérie pneumopathie variole dysenterie malaria botulisme poliomyélite septicémie varicelle polio rougeole méningite*



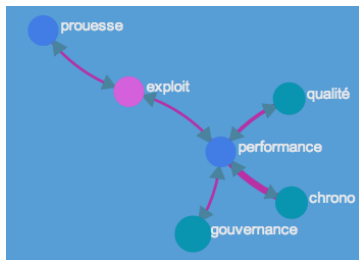
Évaluation

Plusieurs types d'évaluation, dont une par des tests TOEFL construits à partir de wordnets du français (FrenchWordNet, Wolf, JAWS)

Génération aléatoire de tests à partir des synsets

toutefois	<i>néanmoins</i>	complètement	progressivement	sensiblement
exploit	<i>prouesse</i>	offset	plie	bit

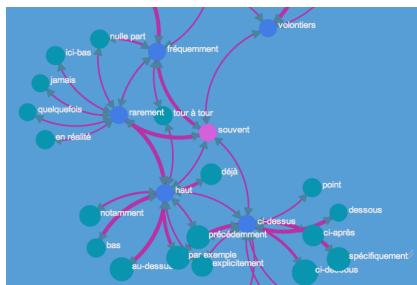
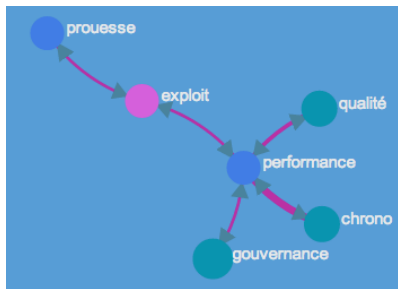
Sélection du candidat le plus proche ($d < 10$) pour répondre
+ choix aléatoire autrement (*baseline hasard=25%*)



network	fwn		wolf	
	%ok	#tests	%ok	#tests
all	51,5	4,121	42,1	7,674
fiscal	46,1	104	37,0	493
business	35,1	248	43,2	1,055
social	39,4	274	37,7	1,345
wolf	64,5	1,076	-	-

Forte différence en fonction des catégories syntaxiques et pb de couverture.

wordnet	pos	#tests	%ok	%bad	%missing	%ok/(ok + bad)
wolf	v	3,876	35,5	30,9	33,6	53,4
	nc	1,078	33,5	2,1	64,4	94,0
	adj	2,085	22,3	11,3	66,4	66,3
	adv	1,533	36,9	41,9	21,7	46,8



Analyse de la connectivité (POS)

corpus	partie du discours	#mots	#conservés	%c/m
cpl	np	1,779,848	8,749	0,5
	nc	35,417	5,782	16,3
	v	11,224	2,480	22,1
	adj	10,198	3,108	30,5
	adv	2,693	776	28,8
fiscal	terme	50,479	4,981	9,9
	nc	18,760	1,976	10,5
	v	4,783	593	12,4
affaires	terme	65,138	5,142	7,9
	nc	23,506	2,095	8,9

Principe:

- les ressources (même WordNet) sont vus comme des graphes.
- étant donné une connection (a, b) dans le graphe G_1 , vérifier si (a, b) 10-connectable dans G_2
baseline hasard: $\sim 1\%$; résultats 33% (prec. 80%, rappel: 21%)

Faible couverture de notre réseau

mais en fait également faible couverture des ressources de référence !

Exemples de paires non 10-connectables dans French WordNet:

écrivain pianiste

canne betterave

dirigeant activiste

éminence promontoire

tendon auriculaire

informatique biologie

cancer hémorragie

provocation agression

haschisch cocaïne

...

caverne abîme

karaté yoga

typhus leucémie

harpe violoncelle

banderole pancarte

nationalisation privatisation

électricien ajusteur

phoque baleine

informatique génétique

caverne gouffre

facteur enjeu

conjuraison conspiration

cigare cigarette

opportunité intention

scooter fourgonnette

mise application

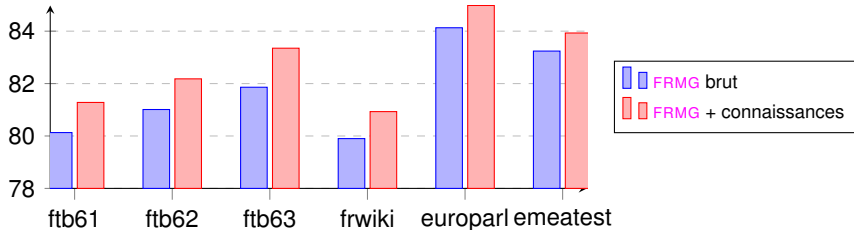
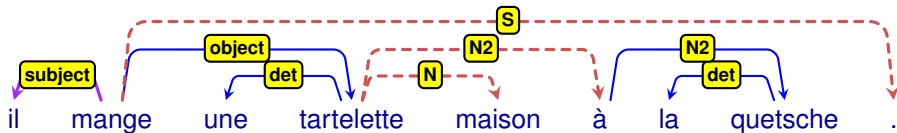
militant dirigeant

cancer infarctus

Évaluation par la tâche

Injection des connaissances acquises dans FRMG (similarité + contextes)
il mange une tartelette maison à la quetsche.

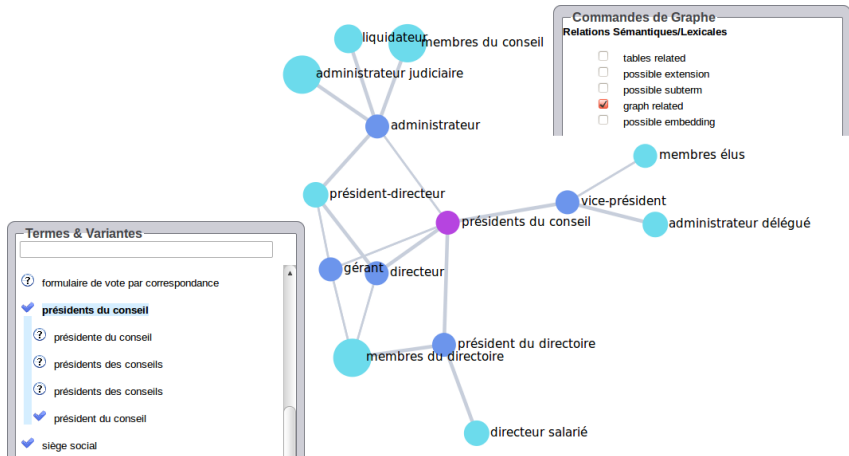
tartelette proche de tarte
quetsche sorte de fruit
aux_fruits contexte fréquent sur tarte } ⇒ tartelette à la quetsche



Vrai besoin de:

- visualisation riche, locale (zoom), mais sans surcharge
- mécanismes simples de navigation et de recherche
- accès à des explications et exemples
- validation collaborative
 - ▶ ressources imparfaites
 - ▶ maintenance et évolutions
 - ▶ effort de validation important devant être distribué
 - ▶ échanges et discussions

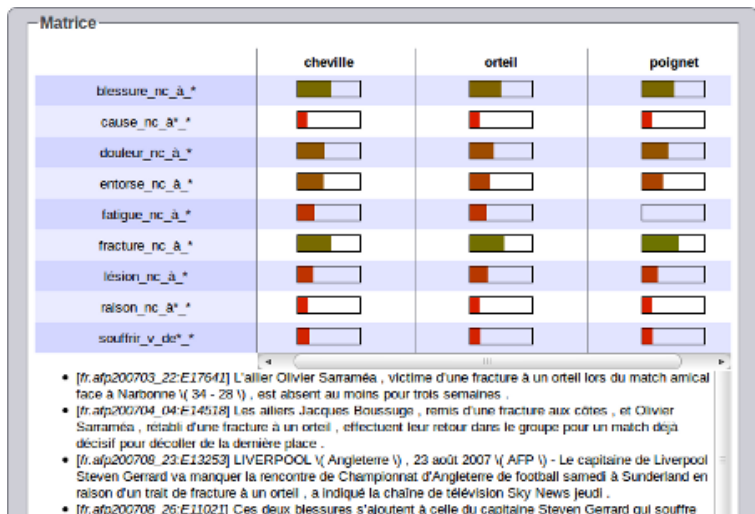
⇒ développement d'une interface WEB sur la plateforme **LIBELLEX** dans le cadre d'une collaboration avec Lingua et Machina.



Matrices explicatives

Les rapprochements explicables par des contextes dont l'importance fournie par

$$\langle ww_{i,j}, c_b \rangle = \sum_a (wc_{i,a})(cc_{a,b})(wc_{j,b})$$



- 1 Une brève présentation de FRMG
- 2 Acquisition de connaissances
- 3 Aller plus loin: quelques pistes**
- 4 Conclusion

D'autres expériences préliminaires menées dans le cadre de SCRIBO sur des sorties d'analyse de corpus:

Concepts

- Extraction de terminologie
garde à vue
- Construction de réseau de mots
- Regroupement de mots en cluster (*synset*), plus regroupement hiérarchique
- extraction de relations ontologiques (par ex. hypéronymie)
navire de guerre: destroyer, aviso

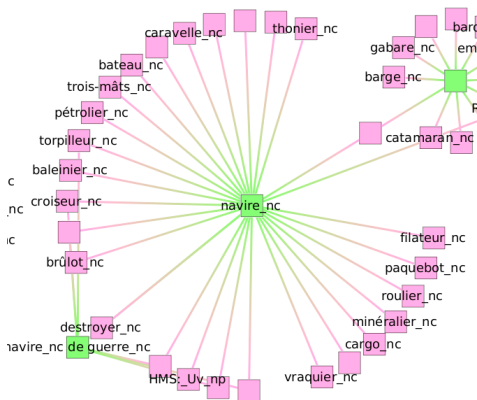
Évènements

- Regroupement de verbes, dénotant un type d'évènement
 - ▶ /transfer/ *donner, offrir, céder*
 - ▶ /communication act/ *annoncer, indiquer, affirmer*
- Identification de paires reliées verbe-nom
 - ▶ *déclarer/déclaration* ;
 - ▶ *identifier/identification* ;
 - ▶ *commencer/commencement/début*
- Découverte de chemins de dépendances caractéristiques entre des paires d'entités d'un certain type

Assigner des labels aux classes [van de Cruys]

Retrouver et analyser les définitions des mots (dans Wikipédia), avec des patrons d'extraction: *X est (une sorte de/une forme de/...) Y* _{genus} . . . differentia

Un *porte-avions* est un *navire de guerre* permettant de transporter et de mettre en oeuvre des avions de combat



Motivation: identifier des verbes correspondant au mêmes types d'événements

Collecte et décompte des *frames remplies*

```
12 accepter_v_active subject:cln_cln object:cla_cla
 9 accepter_v_active object:poste_nc
 7 accepter_v_active subject:_LOCATION_np de_xcomp:reprendre_v
 5 accepter_v_active subject:cln_cln object:quelque chose_pro
```

- Selection des 1000 verbes les plus fréquents (AFP)
- **algorithme:** unsupervised hierarchical clustering
- **measure de similarité:** Jaccard index

Résultats: groupes de verbes proches sémantiquement (max. 4 éléments) :

- partageant des cadres de sous-catégorisation
- avec des arguments proches sémantiquement dans ces cadres
- **pb principal:** la dispersion des données et les arguments non essentiels

- **synonymes**

accepter approuver adopter voter
attirer séduire ravir impressionner
clore clôturer boucler couronner

- **synonyms/antonyms**

chuter diminuer baisser augmenter
rassurer décevoir choquer surprendre

- **"scenarii"**

poser aborder résoudre régler (un problème)
déposer examiner étudier rejeter (une loi)
blesser tuer interpellier libérer (violence action)

- **même événement, autre point de vue**

subir entraîner provoquer causer
vendre acheter racheter acquérir

- **objectif:** identifier des paires verb/noms, où le nom nominalise le verbe:

X *manifester* contre Y – *manifestation* de X contre Y
? les *manifestants* contre Y

- **algorithme:** combiner
 - ▶ similarité distributionnelle: arguments partagés entre verbes et noms (événementiels), Jaccard Index avec PMI
 - ▶ distance d'éditions: forme faible de morphologie dérivationnelle

20 gérer_v / gestion_nc
21 fermer_v / fermeture_nc
22 envisager_v / possibilité_nc
23 reporter_v / report_nc
24 instaurer_v / instauration_nc
25 assigner_v / assignation_nc
26 réagir_v / réaction_nc
27 suspendre_v / suspension_nc
28 déployer_v / déploiement_nc
29 geler_v / gel_nc
30 protéger_v / protection_nc
31 verser_v / montant_nc
32 racheter_v / rachat_nc
33 verser_v / versement_nc
34 dépenser_v / enveloppe_nc
35 circuler_v / circulation_nc
36 durcir_v / durcissement_nc
37 débloquer_v / déblocage_nc
38 octroyer_v / octroi_nc
39 diffuser_v / diffusion_nc
40 contracter_v / propagation_nc
41 renforcer_v / renforcement_nc
42 exécuter_v / exécution_nc

Ticket [38] octroyer_v / octroi_nc

Status:

- Confidence 0.3019
- Rank 1

Morphology

- noun = +

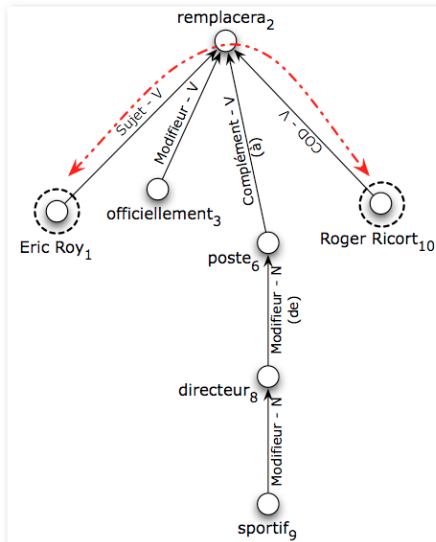
Roles

Verb role	Noun role	strength	Lemma
cod	de	0.21	asile prêt enveloppe prime
à	à	0.12	

Sentences

- [afp200901_14:E12198] La Constitution brésilienne **octroie l'asile** et empêche l'extradition d'étrangers pour des d'opinion, sans distinction idéologique.
- [afp200704_03:E4676] LONDRES, 3 avr 2007 (AFP) - Les nouveaux propriétaires américains de Liverpool ont o 60 millions d'euros à l'entraîneur Rafael Benitez pour recruter l'été prochain, affirme mardi la presse britannique.
- [afp200908_20:E17113] Un opposant au président du Venezuela Hugo Chavez a demandé **l'asile** politique au Pé Premier ministre péruvien jeudi, alors que les deux pays tentent de normaliser leurs relations, après des frictions **d'asile** à des opposants vénézuéliens.
- [afp200710_19:E19690] Ils veulent également que la direction s'engage à entamer des négociations sur **l'octroi** c financière qui servirait à contribuer aux conséquences sociales de la cession, avant de rentrer en négociation ex repreneur.

- Extraction of dependency paths between NEs
- **Hypothesis:** the paths denote the relation between NEs
~> patterns for the relation
- Step 1: **Clustering**
 - ▶ NE pairs given paths
 - ▶ paths given NE pairs
- Step 2: **Acquisition** (2 methods)
 - ▶ partially supervised: acquisition by induction
 - ▶ non-supervised: classification based on the shared contexts



NE pairs \rightsquigarrow paths

http://alpage.inria.fr/~nakamura/afp0520_0729_5_ch
http://alpage.inria.fr/~nakamura/afp0520_0729_5_chemins.xml

Chemin : [41] ==> (MOD-N) président (nc) <== (MOD-N)(de)

- [27] François Bayrou - MoDem
- [2] Jacques Rogge - Comité international olympique
- [3] Jacques Delors - Commission
- [1] Philippe de Villiers - MPF
- [9] Jean-Claude Trichet - BCE
- [1] Fahey - de l'AMA John
- [1] Elio Di Rupo - PS
- [2] Allain Bougrain-Dubourg - LPO
- [6] Hervé Morin - Nouveau Centre
- [4] François Bayrou - Modem
- [1] Mouammar Kadhafi - l'Union
- [1] François Bayrou - Mouvement Démocrate
- [4] Jean-Claude Trichet - Banque centrale
- [1] Jérôme Lejeune - Fondation
- [14] Jean-Paul Bailly - La Poste
- [4] Poul Nyrup Rasmussen - Parti socialiste
- [27] Dominique Sopo - SOS Racisme
- [1] Pascal Colombani - conseil d'administration
- [2] Jean-Luc Mélenchon - Parti de gauche
- [2] Romano Prodi - Commission
- [6] Jean-Paul Huchon - conseil régional
- [3] José Sarney - Sénat
- [2] Roger Karoutchi - l'Assemblée

paths \rightsquigarrow NE pairs

http://alpage.inria.fr/~nakamura/afp0520_0729_5_ENs.xml
http://alpage.inria.fr/~nakamura/afp0520_0729_5_ENs.xml

Relation : Ban Ki-moon - l'ONIL (254)

organization

[244] ==>

[1] <==

[4] ==>coréen<==

[2] ==>sud<==

[1] ==>appelle<==Conseil<==

[1] ==>lettre==>secrétaire<==

[1] ==>fait<==unanimité<==patron<==

Relation : José Manuel Barroso - Commission (168)

[3] <==

[41] <==président<==

[14] ==>président<==

[1] ==>président<==

Semi-supervised acquisition by induction

- Acquisition by induction given a few seed examples
- experiments for the membership relation
- 2 month of AFP news
- Results:
 - ▶ extraction of 136 paths, for 1469 NE pairs
 - ▶ manual validation of 178 pairs
⇒ 149 OK; 29 wrong

http://alpage.inria.fr/~nakamura/afp0520_0729_resultat.xml

Couples EN en relation d'appartenance : (rouge = individual, bleu = organization)

- Xavier Bertrand - UMP
- Nicolas Sarkozy - UMP
- Nicolas Sarkozy - New Delhi
- Mahamane Ousmane - CDS
- Calderon - Parti d'action
- Lazarus Murendo - tribunal de première instance
- Fredrik Reinfeldt - Commission
- Roland Koch - Deutschlandfunk
- Franco Frattini - Rai Uno
- Devedjian - BFM
- Sébastien Delahaye - CFDT
- M. Berlusconi - Rai
- Ballack - Bild
- M. Berlusconi - Rai Uno
- Mailly - BFM
- Batho - UMP
- M. Roche - BBC
- M. Obama - CBS News
- Moscovici - BFM
- Pascal Baudouin - CGT
- Fred Irwin - Handelsblatt
- Dubus - Paris
- Huchon - Paris
- René Raimondi - PS
- Biden - ABC

- 1 Une brève présentation de FRMG
- 2 Acquisition de connaissances
- 3 Aller plus loin: quelques pistes
- 4 Conclusion

Conclusions

- Confirmation de l'intérêt de données syntaxiques pour l'acquisition de connaissances lexicales
mais la problématique de l'évaluation reste !
- Un environnement assez complet avec une interface de visualisation utile
Nombreux composants échangeables
(chaîne linguistique, schéma d'annotation, algo. acquisition, ...)
⇒ nombreuses expériences à mener
- Validation des termes faites (domaine juridique)
⇒ (prévu:) amélioration outils et apprentissage classification/reranking
- évolution vers la gestion de ressources plus riches (**WOLF**)
mais aussi extraction de relations plus riches
- Système consultable en ligne
<http://alpage.inria.fr/Lbx> (*guest/guest*)
ok pour des traitements de corpus / chargement de ressources

