

Des phénomènes collocationnels à leurs corrélats prosodiques : comparaisons de deux approches de la phraséologie prosodique dans le corpus *RHAPSODIE*

1

MARIA ZIMINA, NICOLAS BALLIER
EA 3967 CLILLAC-ARP

université
**PARIS
DIDEROT**
PARIS 7



Journée Collocations génériques
One-day Seminar on Generic Collocations

lundi 15 janvier 2018
Université Paris Diderot

Plan

2

- **Phraséologie et corpus oraux** (état de l'art)
- Corpus *Rhapsodie* : annotation prosodique
- **Liens** entre la **prosodie** et des **phénomènes collocationnels**
 - Analyse textométrique de données multiannotées
 - Récurrences observées en début des périodes intonatives
 - Spécificités par genre, besoins communicationnels
- Premières conclusions
- Perspectives

Phraséologie et prosodie : statut expérimental de la recherche

3

“If most of the formulaic expressions we know have been acquired from and are used in speech, the phonological representation of formulaic expressions should, in theory, play a fundamental role in the lexical storage and retrieval.” (Lin, 2013)

- **Lin, Ph. M.S. (2013).** The prosody of formulaic expression in the IBM/Lancaster Spoken English Corpus. *International Journal of Corpus Linguistics* 18 (4), 561–588 (2013).
- **Aston, G. (2015).** Learning phraseology from speech corpora. In: Leńko-Szymańska, A., Boulton, A. (eds.) *Multiple Affordances of Language Corpora for Data-driven Learning (Studies in Corpus Linguistics 69)*, pp. 63–84. John Benjamins, Amsterdam-Philadelphia (2015).

Corpus oraux annotés : ressources singulières

4

Projet *Rhapsodie* :
corpus de référence en français
<http://projet-rhapsodie.fr>

Laboratoires porteurs :

- MODYCO
- IRCAM
- LATTICE
- ERSS
- LPL

Partenaires

- ATILF
- CORAL
- CRDO Paris (Centre de ressources pour la description de l'oral)
- INRIA Bordeaux, Paris
- SYLED



Base de données *Rhapsodie*

5

- Un corpus constitué de 57 échantillons de français parlé
- (5 minutes en moyenne), soit 3 heures de parole (33 000 mots, 89 locuteurs) munies d'une transcription orthographique et phonétique.
- Annotations micro / macro syntaxique et prosodique (plus de 60 couches d'annotations)
- Modélisation de l'interface prosodie, syntaxe, discours en français parlé

Rhapsodie :

interface prosodie / syntaxe / discours

6

- La compréhension du rôle que jouent les indices intonosyntaxiques dans la segmentation du continuum sonore en unités informationnelles et discursives
- La modélisation de l'interface prosodie, syntaxe, discours en français parlé
- La base des données prosodiques alignées sur le temps et des données syntaxiques calibrées sur les tokens syntaxiques.
- ...

Annotation prosodique dans *Rhapsodie* :

Exemple (Lacheret *et al.*, 2014)

7

IPE	que vous soyez devenue une vedette vous étiez normalement entraînée																
IPA	que vous soyez devenue une vedette vous étiez normalement entraînée																
RG	que vous soyez devenue					une vedette			vous étiez			normalement			entraînée		
MF	kvuswajədəvny					ynvədət			vuzetje			nɔr	malmã		ãtrene		
syllable	kvu	swa	je	dəv	ny	yn	və	dət	vu	ze	tje	nɔr	mal	mã	ã	tre	ne
Prom	0	0	0	0	W	0	0	W	0	0	W	S	0	0	0	0	S

Méthodologie :

Combinaison de l'annotation manuelle (*proéminence* et *disfluences*) et automatique :

- 1) Annotation manuelle (critères formels acoustiques et perceptifs)
- 2) Caractérisation automatique des constituants prosodiques à partir de l'annotation manuelle
- 3) Stylisation automatique de contours mélodiques et l'annotation des tons liés aux constituants prosodiques

La hiérarchie prosodique dans *Rhapsodie*

8

Intonational Periods (IPE)

Intonational PAcKages (IPA)

Rhythmic Groups (RG)

Metrical Feet (MF)

Syllabes (avec la proéminence : non-proéminent : **o**, fort : **S**, faible : **W**)

Caractère exploratoire de la recherche

9

- Etablir des **liens** entre la **prosodie** et des phénomènes **phraséologiques**
- Plus de **60 couches** annotations dans *Rhapsodie*
 - morpho-syntaxiques
 - syntaxiques
 - macro-syntaxiques
 - prosodiques
- Point de départ ?
- Options théoriques, pratiques, outils

Approche traditionnelle (1)

10

- Les recherches antérieures sur la phraséologie de l'oral ne tenaient pas compte de la hiérarchie prosodique, c'est-à-dire des différentes tailles des constituants prosodiques (Nespor et Vogel, 2007 ; Lin, 2013)
- Des analyses préliminaires d'exemples de *segments répétés* du corpus Rhapsodie, tels que *jeune fille* (F=20) ou *je veux dire* (F=21) nous ont amenés à observer la **non-congruence** de la récurrence de traits prosodiques, telle que la *proéminence*, et des unités phraséologiques traditionnelles

Exemple :

11

- ... une **jeune**|**strong** **fille**|**weak** euh habillée tout en noir...
- ... c'est une **jeune**|**tail** **fille**|**tail** # pauvre et affamée ...
- § et on se retrouve dans la rue avec une **jeune**|**weak** **fille**|**filled-dis** #

- § vous voyez ce que **je**|**strong** **veux**|**strong** **dire**|**strong** #
- § euh vous voyez ce que **je**|**weak** **veux**|**weak** **dire**|**weak** §
- § **je**|**strong** **veux**|**strong** **dire**|**strong** là ça me laisse rêveur #

Emboîtement des structures

12

per	do~ k@ ja a y n2 z9n Fi j@ a bi je tu ta~ nwaR																												
pkg	do~ k@ ja a y n2 z9n Fi							j@ a bi je tu				ta~ nwaR																	
rhg	do~ k@		ja a y n2			z9n Fi			j@ a bi je tu				ta~ nwaR																
feet	do~ k@		ja a y n2			z9n Fi			j@ a		bi je tu			ta~ nwaR															
pree	W		W			S			W		W																		
syl	do~	k@	ja	y	n2	Z9n	Fi	j@	a	bi	je	tu	ta~	nwaR															
ph	d	o~	k	@	j	a	y	n	2	Z	9	n	F	i	j	@	a	b	i	j	e	t	u	t	a~	n	w	a	R
w/ort	donc		euh	il	y	a	une	jeune	fille	euh	habillée			tout	en	noir													
w/ph	do~k		@	j	a	yn2	Z9n	Fij	@	abije			tut	a~	nwaR														
w/syl	do~		k@	ja	yn2	Z9n	Fi	j@	abije			tu	ta~	nwaR															
<i>english</i>	so		uh	there is	a	young	girl	uh	dressed			all	in	black															

Notre approche (2) : options théoriques

13

- Analyse de séquences préfabriquées qui constituent des **objets phraséologiques plus larges que les expressions figées classiques**
- Mise en relation de ces « prêts à dire » avec les caractéristiques du **genre discursif** dans lequel ces séquences apparaissent

Sitri, F., Tutin, A. (dir.): Phraséologie et genres de discours. Patrons, motifs, routines. *LIDIL* 53 (2016).

Notre approche : outils méthodologiques

14

- **CLILLAC-ARP** : utilisation des corpus pour identifier les **schémas lexicogrammaticaux** (**schémas LG**) :
- composés d'un ou plusieurs élément(s) obligatoire(s) :
 - le « **pivot** »
- et d'un ou plusieurs éléments variables :
 - le « paradigme »
- associés à des **fonctions discursives** spécifiques
 - Les schémas représentent les unités de structure pour la construction des **discours spécialisés**
 - Les **imbrications de schémas LG** mènent à la création de chaînes discursives plus étendues (Gledhill *et al.*, 2017)

Illustration à partir de Rhapsodie

15

- § et donc euh **c'est pour ça qu'**aujourd'hui je suis en italien en XXX ...
- § c'est-à-dire § ouais § un mois < **c'est pour ça que ça** s'appelle radio Timsit ...
- § mais bien sûr donc **c'est pour ça** bien sûr bien sûr **que** je parlais oui XXX ...
- § **c'est pour cela que** je tenais à vous rencontrer la veille de notre fête ...

Données *Rhapsodie* (version initiale)

16

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	U
1	TextID	TreeID	TokenID	Token	Lemma	POS	Mode	Tense	Person	Number	Gender	Gov_rection	Type_rection	Gov_para	Type_para	Gov_inher	Type_inher	Gov_junc	Type_jun	Gov_junc-inhe	Type_junc-in
2	M2006	1	1	bonjour	bonjour	B_I						0	root								
3	M2006	1	2																		
4	M2006	1	3	Eric	Eric	B_N				sg	masc	0	root								
5	M2006																				
6	M2006	2	1	bonjour	bonjour	B_N				sg	masc	0	root								
7	M2006	2	2																		
8	M2006	2	3	à	à	B_Pre						1	dep								
9	M2006	2	4																		
10	M2006	2	5	tous	tous	B_Pro				3	pl	masc	3	dep							
11	M2006																				
12	M2006	3	1	nouvelle	nouveau	B_Adj				sg	fem	3	dep								
13	M2006	3	2																		
14	M2006	3	3	nuit	nuit	B_N				sg	fem	0	root								
15	M2006	3	4																		
16	M2006	3	5	de	de	B_Pre						3	dep								
17	M2006	3	6																		
18	M2006	3	7	pillage	pillage	B_N				sg	masc	5	dep								
19	M2006	3	8																		
20	M2006	3	9	et	et	B_J														5	junc
21	M2006	3	10																		
22	M2006	3	11	d	de	B_Pre								5	para_coord		3	dep_inherited		9	junc

Ces données sont constituées par un certain nombre de textes (*TextID* visible dans la première colonne), chacun d'eux est segmenté en *Unités Illocutoires* (seconde colonne), chacune d'elle est segmentée en *Tokens* (troisième colonne), chacun d'eux est annoté (les autres colonnes)

Base textométrique *Rhapsodie* (S. Fleury)

17

Transcodage des données : *Trame/Cadre*

```
<item type="delim" pos="46"><f> </f><c>BLANK</c><l>BLANK</l><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a></item>
<item type="forme" pos="47"><f>lance</f><c>B_V</c><l>lancer</l><a>indicative</a><a>present</a><a>3</a><a>sg</a><a>-</a><a>ROOT</a><a>-</a><a>-</a><a>-</a><a>-</a></item>
<item type="delim" pos="48"><f> </f><c>BLANK</c><l>BLANK</l><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a></item>
<item type="forme" pos="49"><f>un</f><c>B_D</c><l>un</l><a>-</a><a>-</a><a>-</a><a>sg</a><a>masc</a><a>DEP(51)</a><a>-</a><a>-</a><a>-</a><a>-</a></item>
<item type="delim" pos="50"><f> </f><c>BLANK</c><l>BLANK</l><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a></item>
<item type="forme" pos="51"><f>appel</f><c>B_N</c><l>appel</l><a>-</a><a>-</a><a>-</a><a>sg</a><a>masc</a><a>OBJ(47)</a><a>-</a><a>-</a><a>-</a><a>-</a></item>
<item type="delim" pos="52"><f> </f><c>BLANK</c><l>BLANK</l><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a><a>-</a></item>
```

Trame

```
<p n="D2013 " d="58359" f="59730" nd="85" nf="86"/>
<p n="M2006 " d="1" f="2086" nd="1" nf="2"/>
<p n="M0015 " d="40347" f="40514" nd="59" nf="60"/>
<p n="M0022 " d="60079" f="60554" nd="89" nf="90"/>
<p n="M0010 " d="33229" f="33372" nd="41" nf="42"/>
```

Cadre

Le processus de transcodage des données issues du projet *Rhapsodie* (extrait). Dans la dernière version de la base, chaque item de la *Trame* est associé à **61 niveaux d'annotation** (prosodie, micro et macro-syntaxe).

Base *Rhapsodie* importée dans *Le Trameur* (Fleury et Zimina, 2014)

The screenshot displays the Le Trameur software interface, which is used for text analysis. The interface is divided into several sections:

- Top Menu:** Includes options like Cadre, Ventilation, Section, Forme-Lemme, Catégorie-Tag, Segment, Coo, Stat, Conco, Conco, Patron, Graphe, Relation, Sélection, Rapport, and Param.
- Left Sidebar:**
 - Chargement de la Carte des sections :** A section for loading the section map.
 - Délimiteur de sections :** A section for defining section delimiters, including a text input field for '#' and a 'Partie' checkbox.
 - Parties:** A list of parts: SUBGENRE, INTERACTIVITY, and SOCIAL_CONTEXT.
 - Recherche Forme sur la carte :** A section for searching forms on the map, including a text input field for a regular expression and a 'RegExp' checkbox.
 - Spécificités sur Sections:** A section for specifying section characteristics, including a 'BI-TEXT' checkbox and two input fields for 'V1' and 'V2'.
 - Sélection Annotation :** A section for selecting annotations, including checkboxes for 'Forme', 'Lemme', and 'Catégorie', and a text input field for 'Forme' with the value '1'.
- Main Area:**
 - Grid:** A large grid of checkboxes for selecting subgenres. The subgenres listed are: PROCEDURAL, ORATORY, NARRATIVE, ARGUMENTATION, DESCRIPTION, and CIRCUMFLEXION. The grid shows various combinations of checked and unchecked boxes.
 - Text Editor:** A text editor showing the analyzed text with words highlighted in red. The text is:
1 euh bon pour aller du CRDT à la gare euh de Grenoble je euh ben je sors déjà du CRDT \$
2 je remonte euh l'avenue Général Champon \$
3 je traverse euh face à la euh MDE \$
4 et je euh je continue je continue jusqu'à une place qui est face à la grande poste #
- Bottom Panel:** A panel showing the number of selected sections (0), the number of sections (1), the annotation (1), and the aperçu (50).

Annotations croisées dans *Le Trameur*

The screenshot displays the Le Trameur software interface with three main panels:

- Left Panel (Section):** Contains controls for section management, including a grid for section delimiters, a search box for forms, and a list of parts (PARTIE, GENRE, SUBGENRE).
- Middle Panel:** Shows a list of morphological forms with their frequencies. A red dot highlights the form <I_V_strong> at position 10199. A red arrow points from this form to the corresponding annotation in the text below.
- Right Panel (Le Métier Lexicométrique):** Displays a grid of annotations for the text 'je comprends'. A red dot is placed on the annotation '1' above the 'c' in 'comprends', indicating the start of a section.

The text 'je comprends' is shown with a grid of annotations above it. A red dot is placed on the annotation '1' above the 'c' in 'comprends', indicating the start of a section. A red arrow points from the form <I_V_strong> in the middle panel to this annotation.

Annotations croisées: The text 'je comprends' is annotated with a grid of boxes. A red dot is placed on the annotation '1' above the 'c' in 'comprends', indicating the start of a section. A red arrow points from the form <I_V_strong> in the middle panel to this annotation.

Annotations croisées: The text 'je comprends' is annotated with a grid of boxes. A red dot is placed on the annotation '1' above the 'c' in 'comprends', indicating the start of a section. A red arrow points from the form <I_V_strong> in the middle panel to this annotation.

Nouvelle couche d'annotation qui fusionne la catégorie (POS), l'information concernant le début de la période intonative (IPE) et la prééminence (exemple)

Prise en compte des constituants prosodiques : *hypothèses*

20

- La segmentation en **constituants prosodiques** est liée à l'**organisation discursive** (en relation avec le genre)
- La hiérarchie des constituants prosodiques (IPE/IPA/RG/MF) a des régularités qui reflètent **la présence des unités de structures** (**schémas LG**) pour la construction de l'oral
- L'analyse de la **saillance initiale** d'un constituant prosodique peut constituer l'un des points d'entrée

Saillance initiale dans les IPEs

(Zimina et Ballier, 2017)

21

Catégorie	Saillance initiale (B_IPE)	Total dans le corpus
Cl (Clitic pronoun)	511 occ.	4 179 occ.
J (Coordinating conjunction)	443 occ.	1 142 occ.
I (Interjection)	439 occ.	1 984 occ.
Adv (Adverb)	287 occ.	2 789 occ.
Pre (Preposition)	238 occ.	3 443 occ.
D (Determiner)	209 occ.	4 080 occ.
V (Verb)	112 occ.	5 994 occ.
Qu (Relative pronoun)	97 occ.	799 occ.
CS (Subordinating conjunction)	74 occ.	729 occ.
N (Noun)	65 occ.	6 317 occ.

IPEs : segments avec saillance initiale

22

Segment répété (catégories)	Saillance initiale (B_IPE)	Total dans le corpus
CL + V	257 occ.	2 223 occ.
D + N	129 occ.	2 919 occ.
Pre + D	90 occ.	1 112 occ.
J + Cl	77 occ.	164 occ.
Cl + Cl	76 occ.	525 occ.
J + Adv	70 occ.	150 occ.
Cl + Cl + V	69 occ.	479 occ.
J + I	67 occ.	107 occ.
I + I	60 occ.	258 occ.
Pre + D + N	55 occ.	939 occ.

La méthode des **spécificités**

(Lebart et Salem, 1994. *Statistique textuelle*. Dunod.)

23

PARTIES

<i>Unités textuelles</i>			
		K_{ij}	F_i
		t_j	

Tableau lexical :

K_{ij} : fréquence de l'unité j dans la partie i

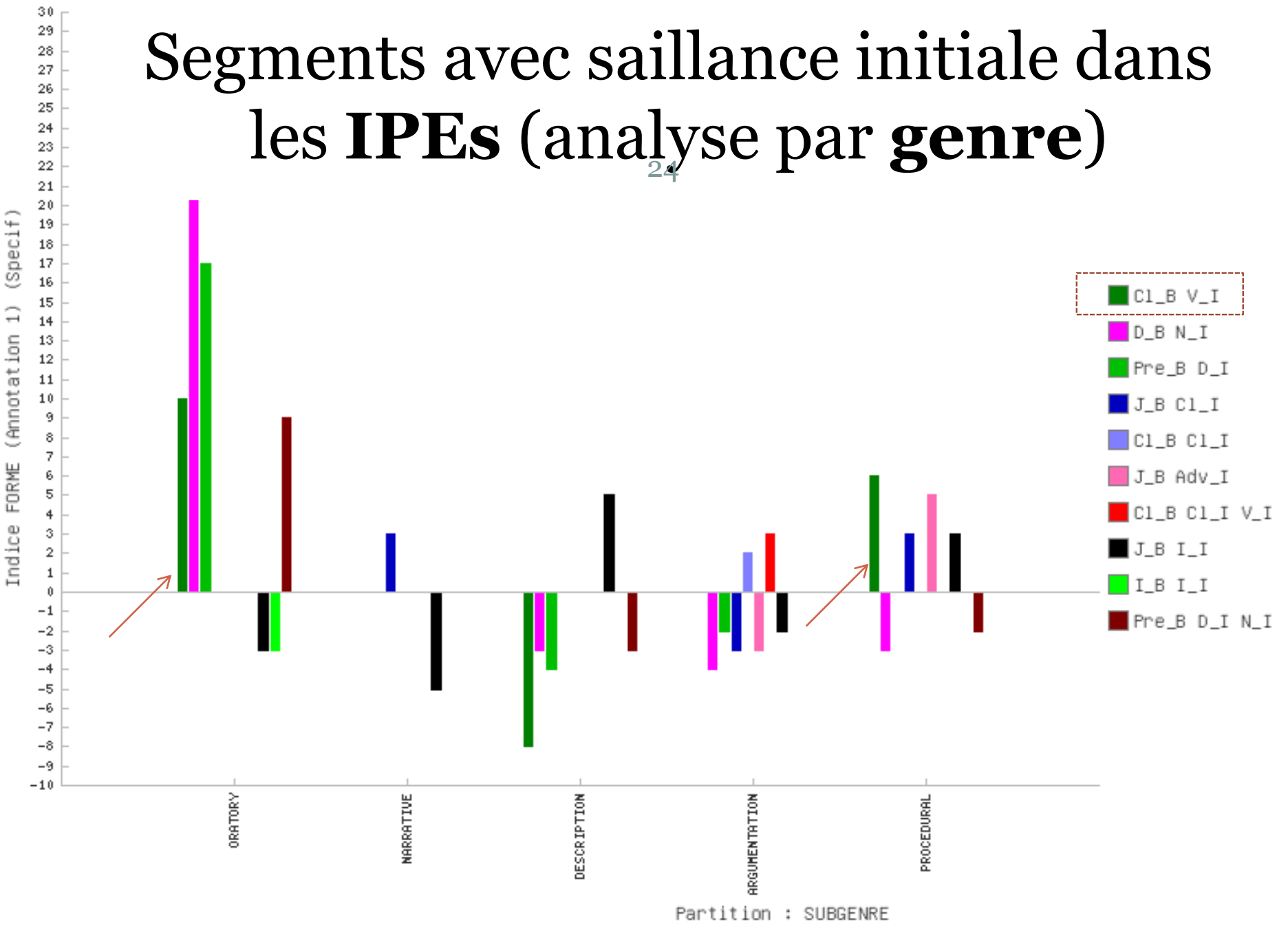
F_i : fréquence de i dans le corpus

t_j : taille de la partie j

Si l'effectif K_{ij} ne se situe pas dans les limites de ce que le **calcul probabiliste** permettait d'espérer, on calcule un *indice de spécificité* : **sur-emploi** ou **sous-emploi** de l'unité (spec. $\pm xx$)

Segments avec saillance initiale dans les IPEs (analyse par genre)

24



Contextes avec saillance initiale

25

Oratory : CL + V

- # **je suis** heureux de me retrouver ce soir #
- [...] est la nation entière qui vous rend hommage # **elle salue** la loyauté #
- # **il faut** les faire grandir #
- # **je souhaite** que l'Europe #

énoncés performatifs

Oratory : D + N

- # **la démocratie** politique et sociale #
- # **la France** sera ce que nous voudrons qu'elle soit # une nation unie #
- # **le droit** de grève # le droit à l'instruction #
- # **un moment** fort #
- # **l'exigence** de solidarité #

focus

Procedural : Cl + V

- # **on passe** devant le le kiosque à journaux #
- # **tu vas** tout droit #
- # **vous continuez** # vous prenez le rond-point tout droit #
- # **on traverse** la rue #
- # **tu descends** toute la pente #

instructions

Le caractère « # » marque le début de la période intonative (IPE)

Proéminence initiale (comparaison par genre pour **CL + V** en début d'IPEs)

26

Indice FORME (Annotation 1) (Specif)

Strong

Weak

Oratory

Procedural

B_C1_strong I_V_strong
B_C1_weak I_V_weak

ORATORY

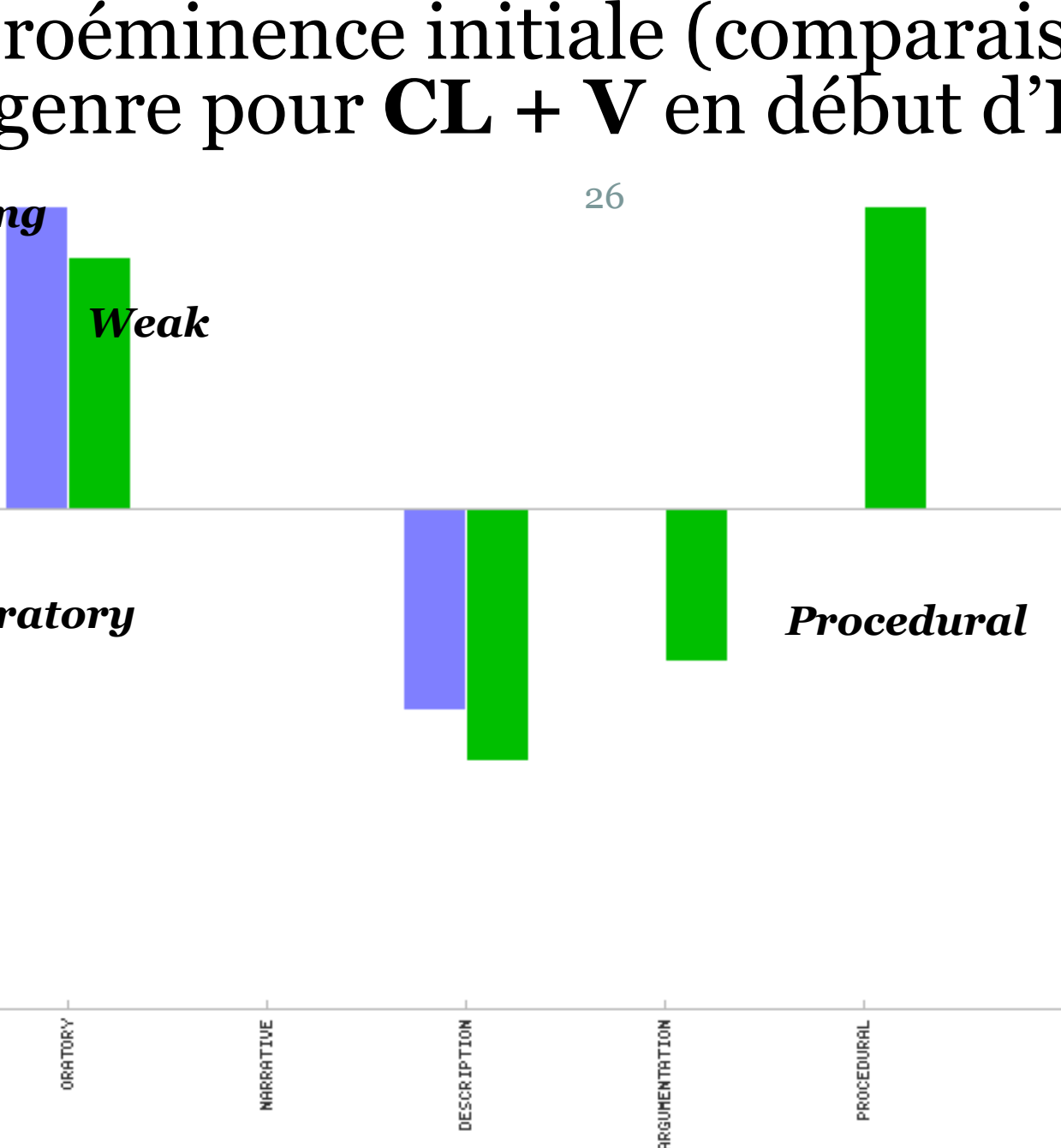
NARRATIVE

DESCRIPTION

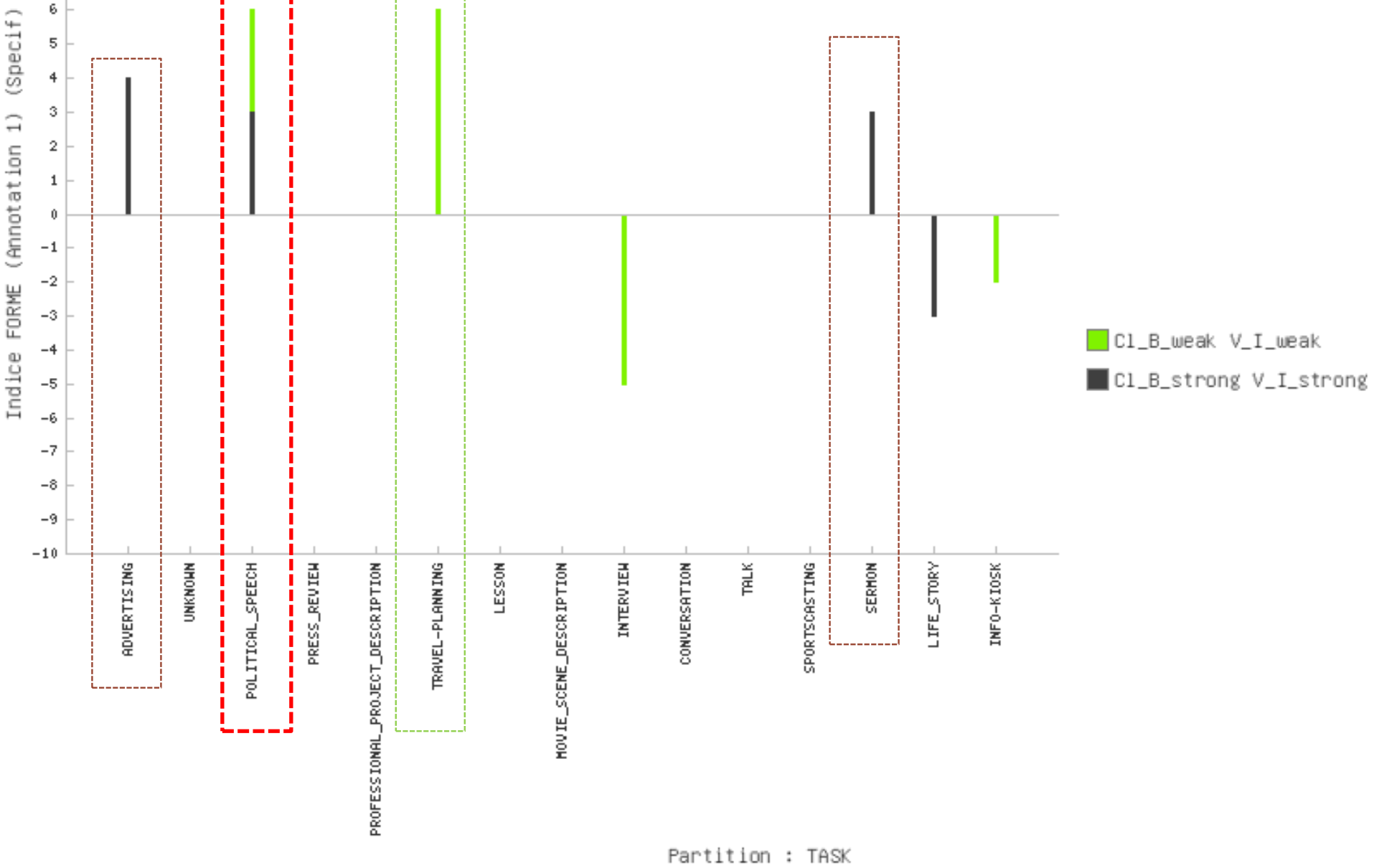
ARGUMENTATION

PROCEDURAL

Partition : SUBGENRE



système complexe des fonctions discursives du discours politique :
- > fonctions structurante, décisionnelle, pédagogique, thérapeutique
(A. Dorna, 1995)



Saillance initiale avec PROÉMINENCE : pivot CL + V (TASK : *political speech*)

28

proéminence FORTE (*Strong*)

----- TASK-POLITICAL_SPEECH -----

dix-neuf dans l'épreuve # \$ **je** pense aux nombreuses victimes de la tempête
ce qui nous semblait acquis # \$ **nous** voyons combien # tout peut être
les règles et les habitudes # \$ **je** comprends ces mouvements de l'âme #
et qui garantissent notre avenir # \$ **nous** avons choisi ensemble de faire grandir la
de faire reculer la pauvreté # \$ **ce** sera tout le sens du combat de
\$ mes chers compatriotes # **je** mesure l'honneur et la responsabilité #
-mer # de l'étranger # **je** souhaite très chaleureusement # une bonne

proéminence FAIBLE (*Weak*)

----- TASK-POLITICAL_SPEECH -----

mes chers compatriotes # **je** voudrais d'abord # exprimer ma
dont nous partageons la peine # \$ **je** pense à nos concitoyens # cruellement
coordination des moyens du pays # \$ **nous** mesurons surtout le prix de l'aide
est devenu # contemporain immédiat \$ **je** suis sûr que beaucoup d'entre vous
de l'éthique # \$ **je** sais que bien des tragédies aujourd'hui
ne doit pas nous diviser # \$ **elle** doit profiter # à chacun #
\$ en les faisant vivre # **nous** serons plus forts pour aborder les temps
\$ la France change # \$ **elle** doit le faire au rythme du monde
de passion # d'enthousiasme # **elle** continue comme hier # à ouvrir
plus fraternel # plus volontaire # **il** aura les couleurs # que nous

Premières conclusions

29

- Observations très partielles par rapport à la richesse du corpus : fondées exclusivement sur quelques niveaux d'annotation et sur un type de saillance
- L'analyse textométrique de constituants prosodiques (**IPE**) a fait émerger des **éléments caractéristiques** dont la **saillance initiale** et la **proéminence** varient en fonction de **genres discursifs**
- Les « **pivots** » recensés correspondraient aux éléments stables de **schémas LG** qui structurent l'oral ("*je salue*", "*elle souhaite*", "*il faut*", "*on continue*", etc.)
- La **saillance initiale** reflète les besoins communicationnels (interaction, tours de parole, ...)

Pistes de recherche

30

- **Co-occurrences de régularités** à d'autres échelles de la hiérarchie prosodique
- **Emboîtements des constituants** prosodiques ou frontières complexes entre niveaux d'analyse
- *Prise en compte de **plusieurs niveaux d'annotation** prosodiques (proéminences, pauses, tonalité, ...)*
- ***Critères perceptifs** utilisés lors de l'annotation manuelle des constituants prosodiques ?*
- Test d'identification de séquences prosodiques phraséologiques par re-synthèse (*humming*)
- Comparaison avec les travaux sur l'anglais (annotation 1984 du **MARSEC**)

MERCI



Journée Collocations génériques
One-day Seminar on Generic Collocations

lundi 15 janvier 2018
Université Paris Diderot

Références (extrait) /.../

Aston, G. (2015). Learning phraseology from speech corpora. In: Leńko-Szymańska, A., Boulton, A. (eds.) *Multiple Affordances of Language Corpora for Data-driven Learning (Studies in Corpus Linguistics 69)*, pp. 63–84. John Benjamins, Amsterdam-Philadelphia (2015).

Gledhill C., Patin S., Zimina M. (2017). « Lexico-grammaire et textométrie : identification et visualisation de schémas lexico-grammaticaux caractéristiques dans deux corpus juridiques comparables en français. » *Corpus 17* .

Fleury, S., Zimina, M. (2014). Trameur: A Framework for Annotated Text Corpora Exploration. *COLING 2014 the 25th International Conference on Computational Linguistics: System Demonstrations*, August 2014, Dublin, Ireland, pp. 57-61.

Lacheret, A., Kahane, S., Beliao, J., Dister, A., Gerdes, K., Goldman, J-P., Obin, N., Pietrandrea, P., Tchobanov, A. (2014). Rhapsodie: a Prosodic-Syntactic Treebank for Spoken French. In: *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, Reykjavik, Iceland, May 2014.

Lin, Ph. M.S. (2013). The prosody of formulaic expression in the IBM/Lancaster Spoken English Corpus. *International Journal of Corpus Linguistics* 18 (4), 561–588 (2013).

Sitri, F., Tutin, A. (dir.) (2016). Phraséologie et genres de discours. Patrons, motifs, routines. *LIDIL 53* (2016).

Zimina M., Ballier N. (2017). Intonational PERiods (IPE) and Formulaic Language: A Genre-based Analysis of a French Speech Database. *Proceedings of Europhras 2017 Conference of 13-14 November 2017*, London: *Computational and Corpus-based Phraseology: Recent Advances and Interdisciplinary Approaches. Volume II.*