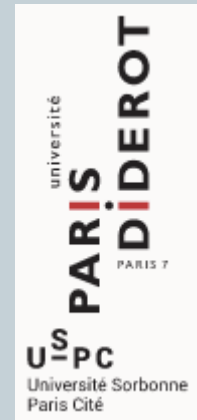


Analyse des seuils de constituance de la phraséologie prosodique dans le corpus *Rhapsodie* : *période intonative (IPE) et saillance initiale*

1

MARIA ZIMINA, NICOLAS BALLIER
EA 3967 CLILLAC-ARP



Colloque
PHRASÉOLOGIE FRANÇAISE
Phraséologie française 2017, Université d'Artois, Arras, France.
21-22 sept. 2017 Arras (France)



je te dis quoi

Plan

2

- **Phraséologie et corpus oraux** (état de l'art)
- Corpus *Rhapsodie* : annotation prosodique
- **Liens** entre la **prosodie** et des **phénomènes phraséologiques** :
 - Analyse textométrique de données multiannotées
 - Récurrences observées en début des périodes intonatives
 - Spécificités par genre
- 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.:** The prosody of formulaic expression in the IBM/Lancaster Spoken English Corpus. *International Journal of Corpus Linguistics* 18 (4), 561–588 (2013).
- **Aston, G.:** 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

Structures prosodiques 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**)

Nature 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

Options théoriques (1/2)

10

- 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).

Options théoriques (2/2)

11

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

Exemple :

<Il y a lieu de **procéder** à une évaluation des différentes méthodes de recyclage>
<il convient de **procéder** à un second examen de toutes les régions du système nerveux qui présentent ces altérations>

Données *Rhapsodie*

12

	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 (l'identifiant du texte est 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)

13

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*

The screenshot displays the 'Le Trameur' software interface. At the top, a menu bar includes 'Cadre', 'Ventilation', 'Section', 'Forme-Lemme', 'Catégorie-Tag', 'Segment', 'Coo', 'Stat', 'Conce', 'Patron', 'Graphe', 'Relation', 'Sélection', 'Rapport', and 'Param'. The 'Section' menu is active, showing a grid of checkboxes for various subgenres: SUBGENRE, PROCEDURAL, ORATORY, NARRATIVE, ARGUMENTATION, and DESCRIPTION. Below the grid, there are several rows of checkboxes, some of which are checked. A status bar at the top right indicates 'seuillage : 1 5 10 ++ | Modifier le seuillage :'. Below the grid, there is a status bar showing 'Nb 1. Sections sélectionnées : 0 N° Sect. : 1:(1,119) Annotation : 1 Aperçu : 50'. The main text area contains the following text: 'euh bon pour aller du CRDT à la gare euh de Grenoble je euh ben je sors déjà du CRDT \$ je remonte euh l'avenue Général Champon \$ je traverse euh face à la euh MDE \$ et je euh je continue je continue jusqu'à une place qui est face à la grande poste #'. The text is displayed in a monospaced font, with some words highlighted in red and others in yellow. At the bottom, there is a status bar with 'Annotations : 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24'.

Annotations croisées dans *Le Trameur*

The screenshot displays the Le Trameur software interface. On the left, there are control panels for 'Section' (with a 'Chargement de la Carte des sections' message), 'Recherche Forme' (with a search box containing '^[\W]_Adj_strong' and 'RegExp' checked), and 'Spécificités sur Sections' (with 'BI-TEXT' selected and 'VI' set to 1, 'V2' set to 2). The main window shows a text editor with the text 'je compren' and a list of morphological annotations on the right. A red arrow points from the annotation '<I_V_strong>' to a new annotation layer at the bottom of the interface, which contains the text 'ne # Nouvelle couche d'annotation qui fusionne la catégorie (POS), l'information concernant le début de la période intonative (IPE) et la prééminence'.

Position: <10199>
Forme: <comprends>|Freq: 3
Lemme: <comprendre>|Freq: 17
Cat: <V>|Freq: 5994
a-00004: |Freq: 34451
a-00005: <indicative>|Freq: 4313
a-00006: <present>|Freq: 3700
a-00007: <1>|Freq: 1939
a-00008: <sg>|Freq: 17661
a-00009: <->|Freq: 22506
a-00010: <ROOT>|Freq: 6169
a-00011: <ROOT>|Freq: 6169
a-00012: <->|Freq: 37506
a-00013: <->|Freq: 36350
a-00014: <->|Freq: 35841
a-00015: <->|Freq: 38394
a-00016: <0>|Freq: 25855
a-00017: <I>|Freq: 31763
a-00018: <I>|Freq: 25322
a-00019: <0>|Freq: 35163
a-00020: <0>|Freq: 36100
a-00021: <0>|Freq: 37537
a-00022: <0>|Freq: 37880
a-00023: <0>|Freq: 37278
a-00024: <0>|Freq: 37682
a-00025: <0>|Freq: 36644
a-00026: <0>|Freq: 36783
a-00027: <0>|Freq: 36531
a-00028: <0>|Freq: 35917
a-00029: <0>|Freq: 36976
a-00030: <S>|Freq: 11651
a-00031: <0>|Freq: 22720
a-00032: <->|Freq: 34124
a-00033: <80.91585028476825>|Freq: 1
a-00034: <89.71624864605332>|Freq: 1
a-00035: <U>|Freq: 23747
a-00036: <V_hOH1>|Freq: 3
a-00037: <hOH1>|Freq: 3
a-00038: <259.3999999999909>|Freq: 143
a-00039: <198.750000000000398>|Freq: 16
a-00040: <\$L1>|Freq: 24847
a-00041: <->|Freq: 34494
a-00042: <I>|Freq: 31417
a-00043: <->|Freq: 35575
a-00044: <hLH1>|Freq: 31
a-00045: <L>|Freq: 7984
a-00046: <->|Freq: 27570
a-00047: <->|Freq: 27570
a-00048: <lone>|Freq: 14262
a-00049: <hOH2>|Freq: 14
a-00050: <L>|Freq: 9688
a-00051: <->|Freq: 23983
a-00052: <->|Freq: 23983
a-00053: |Freq: 20944
a-00054: <hOH2>|Freq: 7
a-00055: <L>|Freq: 9379
a-00056: <->|Freq: 21603
a-00057: <->|Freq: 21603
a-00058: |Freq: 19125
a-00059: <hOH2>|Freq: 6
a-00060: <200.702971>|Freq: 1
a-00061: <201.112971>|Freq: 1
a-xxxxx54: <_>|Freq: 38423
a-xxxxx55: <comprends>|Freq: 3
a-xxxxx56: <V_strong>|Freq: 2856
<I_V_strong>|Freq: 2444
<strong_L>|Freq: 16897

Le Métier Lexicométrique @CLA2T-P3 V. 12.148
Patron Graphe Relation Sélection Rapport Param
Ctrl-clic sur carré : sélection | Shift-Control-clic sur sélection : désélection
de page : sélection 5 sections | Shift-control-clic sur marqueur de page : sélecti
1 Aperçu : 50
ne # Nouvelle couche d'annotation qui fusionne la catégorie (POS), l'information concernant le début de la période intonative (IPE) et la prééminence

Hypothèses

16

- La segmentation en **constituants prosodiques** est liée à l'**organisation discursive** (en relation avec le genre).
- La constitution des périodes intonatives (IPEs) 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** de l'IPE est le point d'entrée.

IPEs : saillance initiale

17

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

18

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*.)

19

PARTIES

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

Tableau lexical :

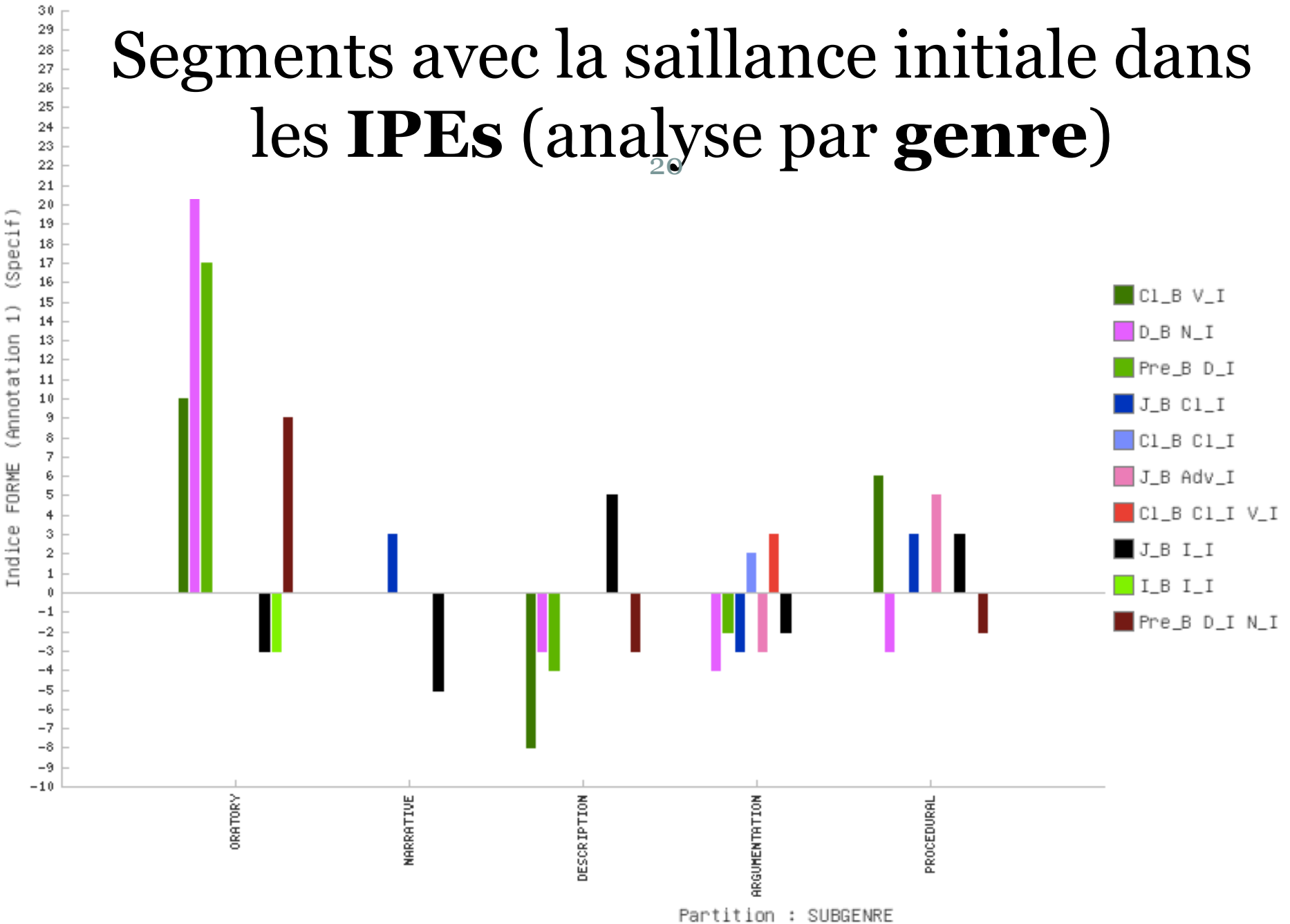
K_{ij} : fréquence de l'unité j dans la partie i

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** permettrait d'espérer, on calcule un *indice de spécificité* : **sur-emploi** ou **sous-emploi** de l'unité (spec. $\pm xx$)

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



Contextes avec saillance initiale

21

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 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)

Premières conclusions

22

- L'analyse textométrique de constituants prosodiques (**IPEs**) a fait émerger des **éléments caractéristiques** dont la **saillance initiale** varie 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).

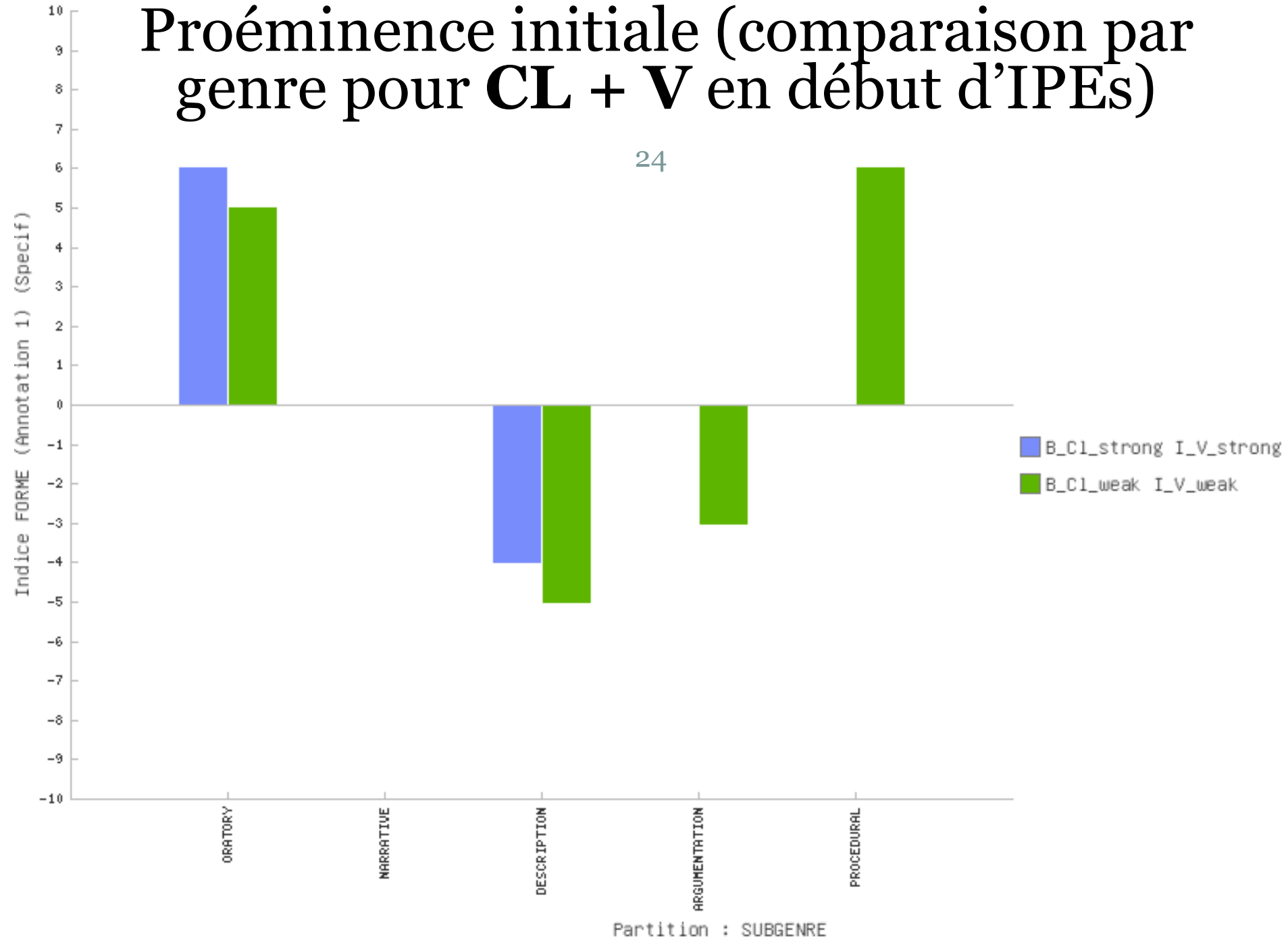
Perspectives

23

- *Prise en compte de plusieurs niveaux d'annotation prosodiques (proéminences, pauses, tonalité, ...)*

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

24



Saillance initiale avec PROÉMINENCE : pivot CL + V (genre : *Oratory*)

25

proéminence FORTE

quatre-vingt-dix-neuf dans l'épreuve # \$ je pense aux nombreuses victimes de la tempête # et
endeuillées dont nous partageons la peine # \$ je pense à nos concitoyens # cruellement touchés dans leur
précarité de ce qui nous semblait acquis # \$ nous voyons combien # tout peut être parfois remis en
dans l'épreuve # faire parler leur coeur # je voudrais dire # merci à tous les Français #
on maîtrisait les règles et les habitudes # \$ je comprends ces mouvements de l'âme # \$ pourtant
engagent # et qui garantissent notre avenir # \$ nous avons choisi ensemble de faire grandir la France dans
et capable de faire reculer la pauvreté # \$ ce sera tout le sens du combat de la France
toute sa place # pour dire le droit # le faire respecter # avec autorité # avec justice
bien vivre ensemble # \$ mes chers compatriotes # je mesure l'honneur et la responsabilité # qui m
d'outre-mer # de l'étranger # je souhaite très chaleureusement # une bonne # &

proéminence FAIBLE

mes chers compatriotes # je voudrais d'abord # exprimer ma sympathie #
'assurer la coordination des moyens du pays # \$ nous mesurons surtout le prix de l'aide fraternelle #
deux mille # est devenu # contemporain immédiat \$ je suis sûr que beaucoup d'entre vous # vont
an deux mille cent # \$ ces progrès # ne prendront tout leur sens # que s'ils bénéficient
être le siècle # de l'éthique # \$ je sais que bien des tragédies aujourd'hui font douter
l'intérieur de chaque nation # une exigence # se fait entendre # toujours plus forte # pour
ces valeurs # \$ en les faisant vivre # nous serons plus forts pour aborder les temps qui viennent
qui viennent # \$ la France change # \$ elle doit le faire au rythme du monde # \$
\$ en étant fidèle à son génie propre # elle saura conjuguer le changement et la cohésion sociale #
riche de fièvre de passion # d'enthousiasme # elle continue comme hier # à ouvrir # et
\$ plus fraternel # plus volontaire # il aura les couleurs # que nous lui donnerons

Perspectives

26

- *Critères perceptifs utilisés lors de l'annotation manuelle des constituants prosodiques ?*

Références

Aston, G.: 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.: « Lexico-grammaire et textométrie : identification et visualisation de schémas lexico-grammaticaux caractéristiques dans deux corpus juridiques comparables en français. » *Corpus 17* (2017).

Fleury, S., Zimina, M.: Trameur: A Framework for Annotated Text Corpora Exploration. *COLING 2014 the 25th International Conference on Computational Linguistics: System Demonstrations*, August 2016, Dublin, Ireland, pp. 57-61.

Lacheret, A., Kahane, S., Beliao, J., Dister, A., Gerdes, K., Goldman, J-P., Obin, N., Pietrandrea, P., Tchobanov, A.: 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.: 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.): Phraséologie et genres de discours. Patrons, motifs, routines. *LIDIL 53* (2016).