

Text Mining for Discourse Markers

Karolin Boczoń ^{1,3} Jacques Jayez ^{1,2}

¹LORIA

²ENS de Lyon

³IDMC

13 June 2023



1 Challenges

2 Neural taggers

3 Grammars

4 Hybrid strategy

Ambiguity

The biggest challenge in recognizing discourse markers is their *ambiguity*.

bon, bien: adjective? adverb? noun? DM?

après: preposition? temporal adverb? concessive adverb?

- 1 Challenges
- 2 Neural taggers
- 3 Grammars
- 4 Hybrid strategy

Model

- CamemBERT [2] – large language model for French
finetuned on
- PERCEO [1] – high-quality corpus with morpho-syntactic labels

	<i>Bon</i>	,	tout	est	<i>bon</i>	.
camembert-perceo	INT	KON	PRO:ind	VER:pres	ADJ	FNO
spacy	ADJ	PUNCT	ADV	VERB	ADJ	PUNCT

Table 1: Example POS-tagging comparison

Unstable results

Après	il	y	a	un	problème	de	securite	.
ADV	PRO	VER		DET	NOM	PRP	NOM	FNO
Après	la	maison	<i>vous</i>	<i>trouvez</i>	un	champ	.	
PRP	DET	NOM	PRO	VER	DET	NOM		FNO
Après	la	maison	<i>on</i>	<i>trouve</i>	un	champ	.	
ADV	DET	NOM	PRO	VER	DET	NOM		FNO

Table 2: Adverbial and prepositional uses of *après*

Unwanted decomposition

afin de	ADV + PRP	⌚
afin que	ADV + KON	⌚
après tout	ADV	☺
au contraire	ADV	☺
d'abord	ADV	☺
du coup	ADV	☺

Table 3: Examples of composed discourse markers

But for *Paul a souffert du coup qu'il a reçu*, *du coup* remains ADV.

1 Challenges

2 Neural taggers

3 Grammars

4 Hybrid strategy

Unitex

It is possible to manually model the environment for the discourse marker use and other uses based on a native speaker's intuition.

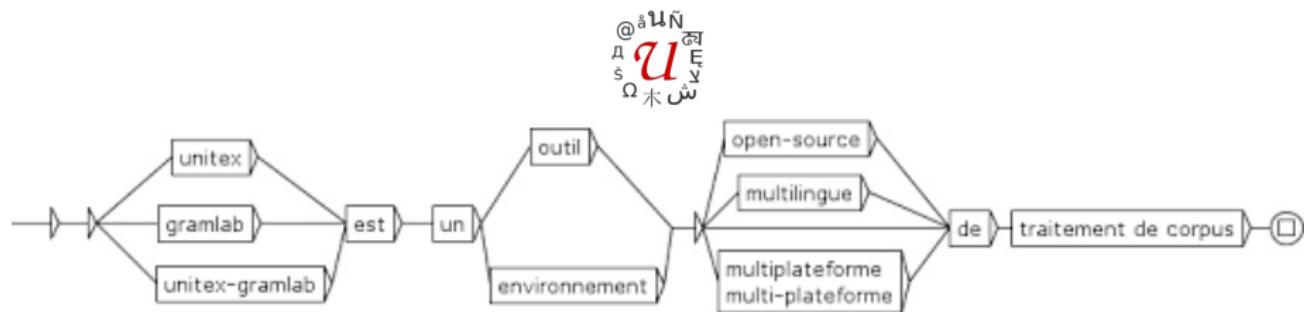


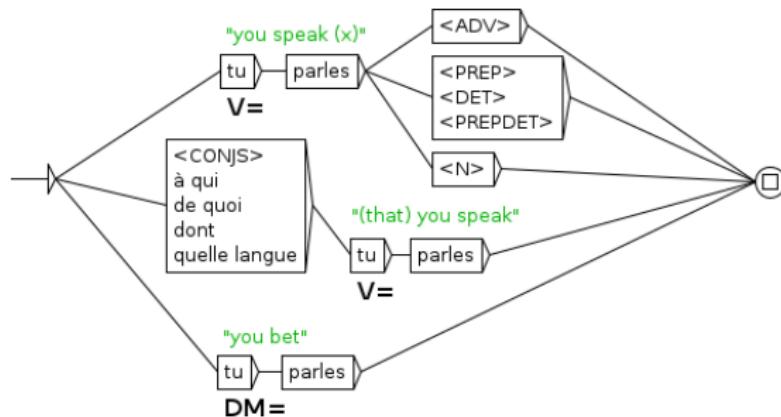
Figure 1: An example of a local grammar graph in Unitex

Example

Tu parles is a pronoun + verb sequence [=you speak] if it is

- followed by an adverb (*parler bien, vite, trop ...*)
- followed by a noun phrase (*parler [langue], à qqn, de qqch, avec qqn ...*)
- preceded by a conjunction or a relative pronoun (*que, quand, à qui, dont ...*)

In all the other cases it should be a discourse marker [=you bet].

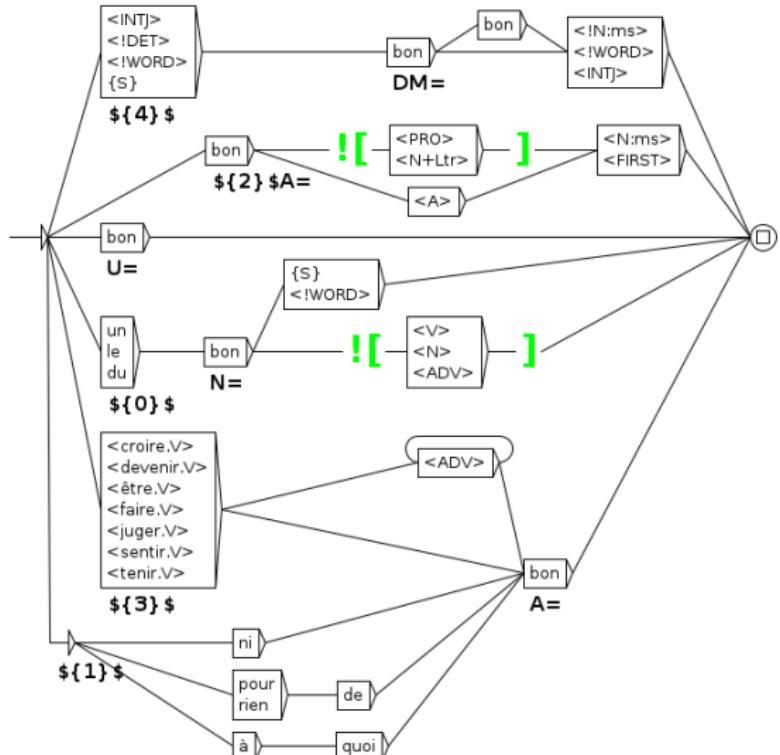


Example

left context	matched sequence	right context
tu te pointes devant un paysan, même si Ils sont restés sans voix : Enfin, Mais comment, me répond-on, Lorsque Regarde	V=tu parles bien _{ADV} V=tu parles pas _{ADV} V=tu parles chinois _N V=tu parles des _{DET} à qui à qui V=tu parles dont _{det} dont V=tu parles même si _{CONJS} V=tu parles quand _{CONJS} V=tu parles que _{CONJS} V=tu parles DM=tu parles DM=tu parles DM=Tu parles	français, faut surtout pas sérieusement ? ? règles de la radio, de ton souci , ont répliqué ses camarades . . bien français, faut surtout
On a vu l'accident, mais pas cette tête Quand tu te pointes devant un paysan, Tais-toi	V=tu parles V=tu parles V=tu parles	, que tu dises où sont tes amis. , qu'est-ce que tu en as de plus ? . . !
Tu nous a coûté 2 millions, il faudra Mais, quand même, on les voit ! _ Oh ! Ah, la fierté, le courage,	V=Tu parles d' _{DET} V=tu parles d' _{DET}	une histoire ! un avenir !
va pas en faire toute une histoire. _ Ils disent qu'on est l'avenir,		

Table 4: Results of parsing with the automaton

It is possible to manually model the environment for the discourse marker use and other uses based on a native speaker's intuition.



But it is also tedious,
time-consuming and not
very robust.

Figure 2: Local grammar for bon

- 1 Challenges
- 2 Neural taggers
- 3 Grammars
- 4 Hybrid strategy

A hybrid strategy

- Not always possible (results of taggers uninteresting).
- Script: tagging → ‘geometry’ of mistakes → finite automaton (FA).
- Typical mistakes
 - ① Ignorance of idioms.
 - ② Function/category confusion.
 - ③ Some black holes (DL black boxes).

A hybrid strategy, cont'd

- Example: *bon* as a DM in a Le Monde + CRFP mixed corpus.
- ≈ 61000 occurrences, sampled to 2000 and reviewed.
- ≈ 9% errors.
- Low precision (0.6), high recall (0.97, virtually no false negative).
- Uneven error distribution.

A hybrid strategy, cont'd

- In terms of types.

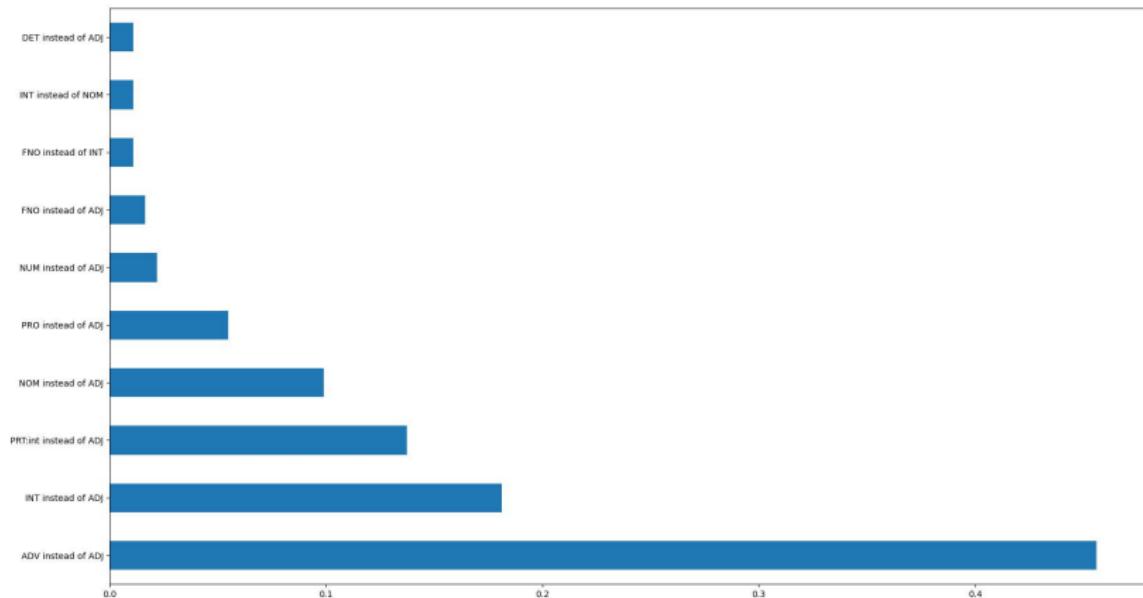


Figure 3: Error Type Distribution

A hybrid strategy, cont'd

- In terms of expressions.

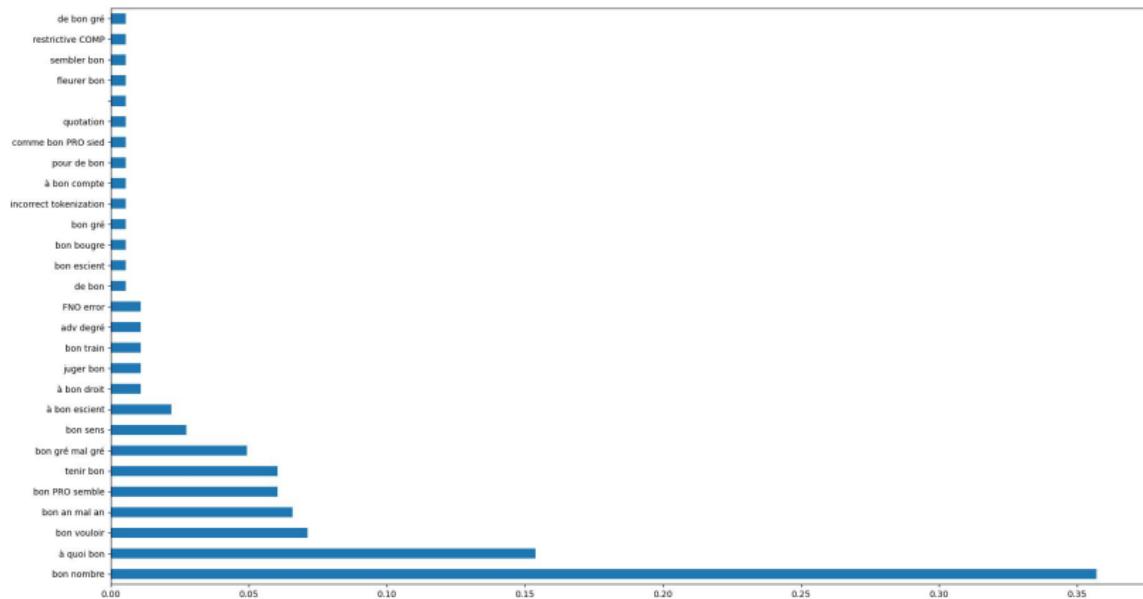


Figure 4: Error Distribution by Expression

A hybrid strategy, cont'd

- Tokenization may be a problem.

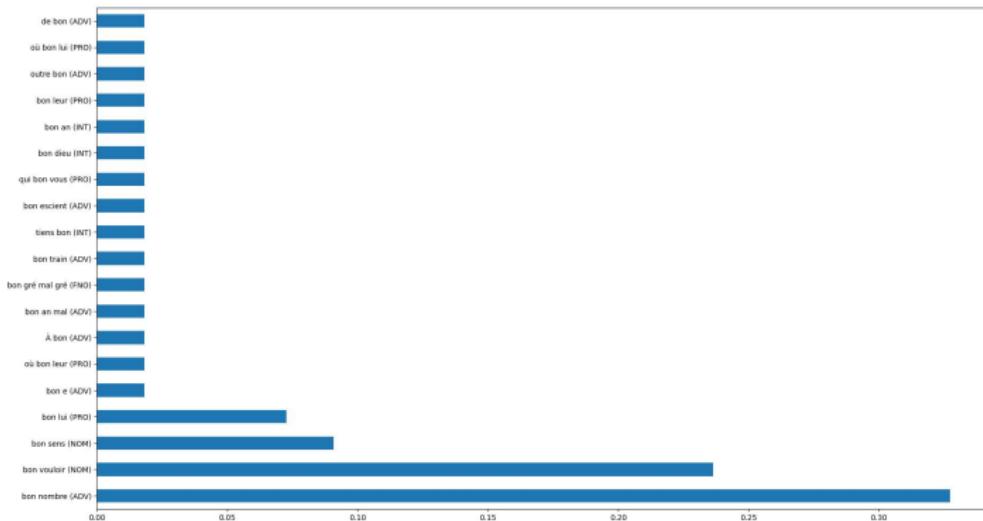


Figure 5: Particular tokenizations

A hybrid strategy, cont'd

- Using ‘repulsion’.
- Environments where *bon* **cannot** be a DM.
- Unitex subgraphs, for instance degree expressions.

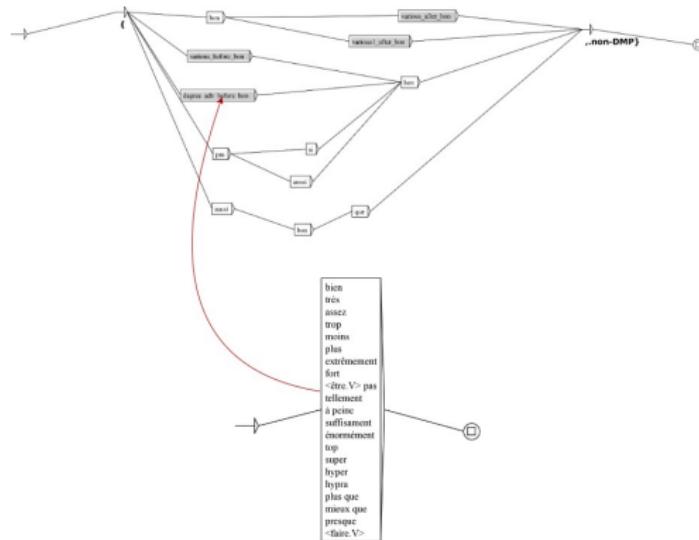


Figure 6: A subgraph in a cascade stage

A hybrid strategy, cont'd

- Final stage: keep the rest.

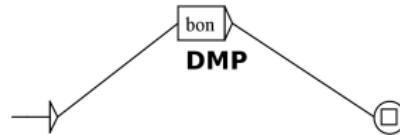


Figure 7: Last Stage (winning the bet)

A Hybrid strategy, cont'd

- Output for a subcorpus of ESLO (eslo.huma-num.fr/).
- 1800 results, out of 1,1511403 words, vocabulary = 17216.
- Eliminates all the previous problems.
- But not the problems with *après* or *tu parles...*

Final diagnostic

- Major problems when the category of an item depends on syntactic structure.
- *Attends le train* (V + NP) vs. *Attends le train est en retard* (V + S).
- Taggers practically useless.
- FA more flexible. Can start a cascade

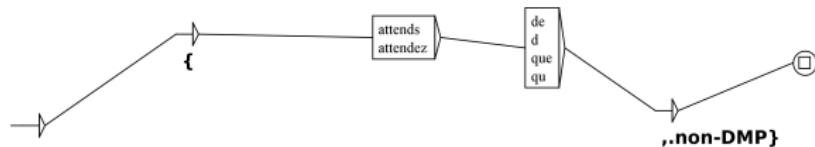


Figure 8: Initial stage

Final diagnostic

- Problem: constraints like.

If there is a finite verb on the right of *attends* and no intervening relative pronoun or complementizer, then we probably have a DM.

- Difficult (possible?) to express in Unitex.
- Easier with an external script.

Final diagnostic

- Conclusion.
 - ① Stick to Unitex for non-ambiguous items.
 - ② Apply the hybrid method when cascades are efficient.
 - ③ Use external scripts (possibly after a cascade).

References

-  ATILF, INIST & LIPN – “Perceo : un projet d'etiqueteur robuste pour l'ecrit et pour l'oral”, 2012, ORTOLANG (Open Resources and TOols for LANGuage) –www.ortolang.fr.
-  L. MARTIN, B. MULLER, P. J. ORTIZ SUÁREZ, Y. DUPONT, L. ROMARY, É. DE LA CLERGERIE, D. SEDDAH & B. SAGOT – “CamemBERT: a tasty French language model”, in *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics (Online)*, Association for Computational Linguistics, 2020, p. 7203–7219.