

# MarsaTag, a tagger for French written texts and speech transcriptions

Stéphane Rauzy, Grégoire de Montcheuil & Philippe Blache

Laboratoire Parole et Langage  
CNRS & Aix Marseille Université, France

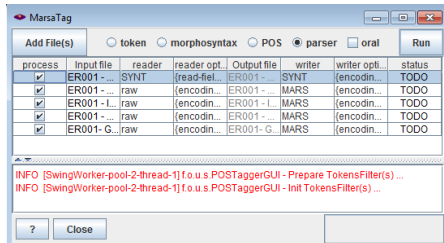
`stephane.rauzy@lpl-aix.fr`

In input, French productions :

- Written texts
- Speech transcriptions

In output, the associated syntactic information :

- Part-Of-Speech tagging, e.g. **Det**, **Noun**, ...
- Groups and syntactic tree structure, e.g. **NP**, **VP**, ...
- Functional relations, e.g. **SUB**, **OBJ**, ...



## MarsaTag

([hdl:11041/sldr000841](https://hdl.handle.net/11041/sldr000841))

The processing chain :

- Pre-lexical treatment, the tokenizer
- Lexical access
- Stochastic tagger
- Stochastic chunker and deep parser

Based on Patterns model approach (HMM stochastic model), the grammar is learned/extracted from an annotated corpus. The performance of the tools depends on :

- The size and the coverage of the training corpus
- The quality of the Gold Standard annotations

# The tokenizer

Split the sentence in a list of tokens and identify tokens (not in the main lexicon) requiring a special treatment :

- Numbers, dates, ...
- Punctuation marks
- Acronyms, abbreviations, ...
- Proper names

*Le budget de l'Armée sera augmenté le 23 août, a indiqué aujourd'hui M. Vasseur à 14H45.*

Le |<sub>L</sub> budget |<sub>L</sub> de |<sub>L</sub> l' |<sub>L</sub> Armée |<sub>L</sub> sera |<sub>L</sub> augmenté |<sub>L</sub> le |<sub>L</sub> 23 |<sub>N</sub> août |<sub>L</sub> , |<sub>P</sub>  
a |<sub>L</sub> indiqué |<sub>L</sub> aujourd'hui |<sub>L</sub> M. |<sub>A</sub> Vasseur |<sub>N</sub> à |<sub>L</sub> 14H45 |<sub>H</sub> . |<sub>P</sub>

A main lexicon allows to associate to each form of the sentence its corresponding lexical distribution of tags, e.g.

form	lemma	sampa	tag	frequency
est	être	E	Aux	1671
est	être	E	Verb	21395
est	est	Est	Noun	422

At this stage, an ambiguous list of tags is proposed for the tokens of the sentence, e.g.

Sentence :        La    valise    est    dans    le    coffre    .  
Propositions :    Det    Noun    Noun    Noun    Det    Noun    Pct  
                         Noun                    Verb    Prep    Pro    Verb  
                         Pro                            Aux

# Stochastic grammars, probabilistic models

Probabilistic models (e.g. Patterns model) allows to associate a probability to any sequence of categories (e.g. tags). The model describes the regularities of the stochastic grammar.

Training stage :

- Learn the parameters of the model on an annotated corpus, e.g. 853 occurrences of **Noun Verb Det Adj**, followed by :

pattern context	category	occurrences	probability
	Pct	12	0.015
Noun Verb Det Adj	Coord	34	0.045
	Noun	807	0.94

Tagging and parsing process :

- Apply the model to the ambiguous data input, the best solution is the one maximizing the probability over the sequence.

# Part-Of-Speech tagging

Form	Solution	Score	A-score	Propositions
<i>La</i>	Da-fs--d-	-6.6362705	A	Pp3fsj- Da-fs--d- Ncm---
<i>définition</i>	Ncfs--	-6.4358025	A	Ncfs--
<i>connait</i>	Vmip3s--	0.74682426	C	Vmip3s--
<i>des</i>	Spd+Da-mp--id	-4.5408936	A	Da-mp--i- Sp-+Da-fp--dd Spd+Da-fp--id Spd+Da-mp--dd Spd+Da-mp--id
<i>nuances</i>	Ncfp--	-5.697094	A	Ncfp-- Vmip2s-- Vmsp2s--
<i>importantes</i>	Afpfp-	-1.0718307	A	Afpfp-
<i>selon</i>	Sp-	-3.0093784	A	Sp-
<i>ces</i>	Dd-mp----	-0.7771969	A	Dd-fp---- Dd-mp----
<i>différents</i>	Afpmp-	-0.046440125	A	Afpmp- Ai-mp- Di-mp----
<i>domaines</i>	Ncmp--	-5.917076	A	Ncmp--
.	Wd	-1.1673737	A	Wd

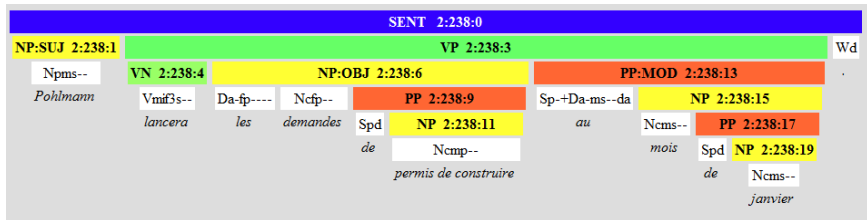
- MarsaLex lexicon : 595.000 entries with frequency computed from a 140 Megawords corpus (mainly newspapers).
- MarsaTag tagger : Trained on the 700,000 manually tag corrected LPL-Grace corpus. Tags set of 51 categories, score of 0.975 (F-Measure) on written texts.

<b>GN</b>	<b>NV</b>	<b>GN</b>	<b>GA</b>	<b>GP</b>	Wd
Da-fs--d- Ncfs-- La définition	Vmip3s-- connaît	Spd+Da-mp--id Ncfs-- des nuances	Afpfp- importantes	Sp- Dd-mp---- Afpmp- Nemp-- selon ces différents domaines	.

Objective : Insert frontiers and label chunks to form syntactic constituents with flat structure.

- MarsaTag chunker : Trained on the 100.000 Easy gold standard. Easy grammar of seven constituents (**GN**, **GP**, **NV**, ...). Score of 0.93 (F-Measure) on written texts.





Objective : Form syntactic constituents and tree structure, and indicate functional relations between the constituents.

- MarsaTag parser : Trained on the 100.000 words LPL-FT corpus, a grammar of 15 constituents (NP, VP, VN, ...) and 9 relations (SUB, OBJ, COORD, ...). Evaluation in progress.

The tools have been adapted for the treatment of speech transcriptions by investigating the CID corpus (Bertrand et al. 2008)

- CID - Corpus of Interactive Data - 8 hours of spontaneous dialogues in French
- Enriched Orthographic Transcription aligned on signal
- 115,000 tokens with manually corrected POS-tags
- Disfluencies annotation
- Various others multi-level annotations available...

# Speech transcription treatment - Filtering

Remove from the transcription the phenomena which are not found in written french :

- hesitation, pause and filled pause (e.g. *heu*, #, +, ...)
- laughter (e.g. @)
- word truncations (e.g. *remp-*)

TOE : # *alors moi j'y étais allée déjà je comprends rien à ce qu'ils font*  
+ *mais euh j'y étais allée pour remp- pour remplir la salle #*

Filtered TOE : *alors moi j'y étais allée déjà je comprends rien à ce qu'ils font mais j'y étais allée pour pour remplir la salle*

# Speech transcription treatment - Punctuation marks

Punctuation marks are not available in the transcription. The tagger allows to insert these marks based on written french model :

- **Wd** : Strong punctuations (e.g. ., !) delimiters of sentences
- **Wm** : Soft punctuations (e.g. ,) delimiters of smaller syntactic units

Example of pattern with punctuation marks insertion :

pattern context	insertion	category	proba
	-	<b>SubPro</b>	0.01
<b>Verb Det Noun</b>	<b>Wd</b>	<b>SubPro</b>	0.34
	<b>Wm</b>	<b>SubPro</b>	0.75

Punctuation marks are inserted if they increase the probability of the sequence treated.

# Speech transcription treatment - Lexicon adaptation

- Some forms are used with different and supplementary purposes than in written productions (e.g. spoken discourse markers). These forms are few (less than 40) but frequent (10% of the CID corpus). Their lexical tags distribution have been modified :

form	tag	CID frequency	probability
bon	Interjection	634	0.95
bon	Adverb	29	0.04
bon	Adjective	5	0.01
...	...	...	...

- Forms proper to speech production have been added to the lexicon (word reductions, onomatopoeia, ...)

# Speech transcription - Evaluation

1	f	@	d	e	k	i	n	u	z	o~	s	R	v	i	m	a~	Z	#	s	phon (19566)
2	en fait	euh	dès qu'		ils	nous	ont		servi		manger		#	s'					token (7363)	
3	en fait		dès qu'		ils	nous	ont		servi		manger			s'					token (719,7422)	
4		Wm																Wm	S-punctual (1243)	
5	I		Cs	P	Pp1-pj-	Vai	Vmps-sm-		Vmn----		Px3fp-								S-morpho (7422)	
6	interjectio		conjunction	pr	pronoun	aux	verb		verb		pronou								S-category (7422)	

323-107208 Visible part 2.373047 seconds 325-480316

MarsaTag speech tagger : an F-measure of 0.947 evaluated on the manually corrected CID corpus.

Chunking looks robust but deep parsing does not work properly :

- Disfluencies disturb the formation of tree structures ?
- Phrase Structure Grammar is not appropriate for modeling the syntax of spontaneous speech in interaction ?

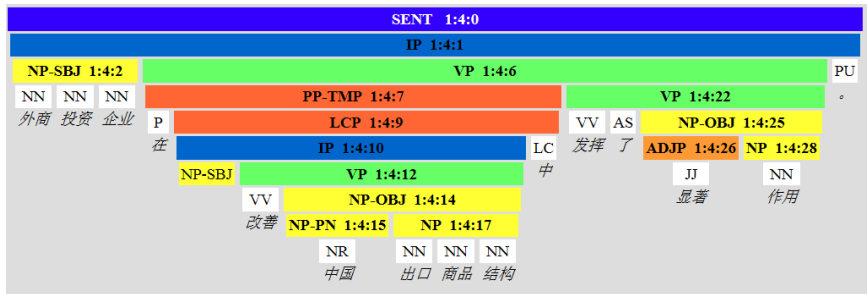
MarsaTag is efficient for extracting the syntactic information from French productions :

- Written texts
- Speech transcriptions (except for deep parsing)

Perspectives :

- Add to the processing chain a model taking into account diffluencies phenomena in speech productions.
- Apply our corpus-based approach to other languages with existing resources :
  - e.g. The Chinese Penn Treebank (Xue et al. 2005) for written texts, PolyU Corpus of Spoken Chinese (Yap et al. 2012) for speech transcriptions.

# Chinese Penn Treebank



Chinese Mandarin Penn Treebank (Xue et al. 2005) :

- Number of sentences : 51,443
- Number of tokens : 1,127,832



# Chinese Penn Treebank

Chinese Mandarin Penn Treebank : lexicon, starting with character 也

Form	Frequency	Category	Examples
也	5244	AD	<a href="#">1003:4:8</a> <a href="#">1003:5:8</a> <a href="#">1005:8:30</a> <a href="#">1006:5:25</a> <a href="#">1007:4:28</a> <a href="#">1008:2:11</a> <a href="#">1010:7:21</a> <a href="#">1013:5:33</a> <a href="#">1014:6:17</a> <a href="#">1016:2:63</a> <a href="#">1016:2:10</a>
也	77	CC	<a href="#">185:7:31</a> <a href="#">220:9:32</a> <a href="#">2249:6:52</a> <a href="#">239:32:86</a> <a href="#">242:5:59</a> <a href="#">245:8:35</a> <a href="#">250:20:34</a> <a href="#">258:9:54</a> <a href="#">282:17:57</a> <a href="#">291:3:55</a> <a href="#">359:11:78</a> <a href="#">460</a>
也	1	NN	<a href="#">2124:0:14</a>
也	30	SP	<a href="#">2232:325:31</a> <a href="#">2249:6:102</a> <a href="#">2271:20:8</a> <a href="#">2271:26:8</a> <a href="#">2271:28:8</a> <a href="#">2271:32:8</a> <a href="#">2271:34:8</a> <a href="#">2271:50:8</a> <a href="#">2271:58:8</a> <a href="#">2271:64:8</a> <a href="#">227</a>
也可以说	1	AD	<a href="#">2374:112:23</a>
也好	34	AD	<a href="#">2150:238:20</a> <a href="#">2150:238:30</a> <a href="#">2150:238:40</a> <a href="#">2155:333:19</a> <a href="#">2157:232:52</a> <a href="#">2157:232:61</a> <a href="#">2157:253:35</a> <a href="#">2157:253:62</a> <a href="#">2225:270</a>
也好	3	SP	<a href="#">1133:8:26</a> <a href="#">1133:8:43</a> <a href="#">2262:15:7</a>
也就是	8	AD	<a href="#">2155:222:8</a> <a href="#">2156:167:29</a> <a href="#">2157:177:11</a> <a href="#">2171:1:40</a> <a href="#">2292:11:49</a> <a href="#">2307:18:6</a> <a href="#">833:20:46</a> <a href="#">885:46:4</a>
也就是说	29	AD	<a href="#">1557:6:41</a> <a href="#">1592:2:89</a> <a href="#">1932:25:35</a> <a href="#">2156:68:31</a> <a href="#">2156:245:3</a> <a href="#">2156:280:71</a> <a href="#">2156:342:4</a> <a href="#">2156:360:67</a> <a href="#">2157:35:4</a> <a href="#">2157:14</a>
也是	8	AD	<a href="#">2223:9:20</a> <a href="#">357:5:37</a> <a href="#">620:3:27</a> <a href="#">762:13:145</a> <a href="#">770:2:90</a> <a href="#">770:80:116</a> <a href="#">854:69:23</a> <a href="#">864:25:47</a>
也是	3	CC	<a href="#">739:15:43</a> <a href="#">797:56:43</a> <a href="#">809:99:39</a>
也有说	1	AD	<a href="#">2208:2:3</a>
也罢	1	SP	<a href="#">2262:15:14</a>
也罢	1	VV	<a href="#">1352:1:73</a>
也许	107	AD	<a href="#">1108:4:42</a> <a href="#">1122:2:36</a> <a href="#">1302:6:12</a> <a href="#">1336:6:41</a> <a href="#">1364:12:5</a> <a href="#">1381:3:63</a> <a href="#">1457:7:8</a> <a href="#">1785:10:7</a> <a href="#">1803:9:6</a> <a href="#">1829:0:25</a> <a href="#">1921:10:</a>
也通社	1	NR	<a href="#">617:6:18</a>
也门	54	NR	<a href="#">1091:0:35</a> <a href="#">1241:1:20</a> <a href="#">1241:2:26</a> <a href="#">1528:3:20</a> <a href="#">1528:4:7</a> <a href="#">1533:0:44</a> <a href="#">1533:1:18</a> <a href="#">1578:1:30</a> <a href="#">1578:2:52</a> <a href="#">1753:33:24</a> <a href="#">1753:3</a>
也门人	4	NN	<a href="#">1528:3:52</a> <a href="#">1533:1:52</a> <a href="#">982:3:52</a> <a href="#">982:4:44</a>

Lexicon extracted from the Chinese Penn Treebank :

- Number of entries : 71, 629