

# SL0720X

## Étiquetage morpho-syntaxique

Franck Sajous (CLLE-ERSS) - [sajous@univ-tlse2.fr](mailto:sajous@univ-tlse2.fr)  
<http://w3.erss.univ-tlse2.fr/membre/fsajous/>

Notes :

- toutes les urls données dans ce document, ainsi que les textes à étiqueter sont présents sous forme de liens cliquables à l'adresse :  
<http://w3.erss.univ-tlse2.fr/membre/fsajous/SDL/SL0720X/>
- au cours des séances, vous utiliserez systématiquement l'éditeur de texte **TextPad** et le navigateur **Firefox**.

## 1 À la découverte de Treetagger

1. À la racine du disque C: doit se trouver un répertoire nommé SL0720. Si ce n'est pas le cas, créez-le : ouvrez le poste de travail, clic-droit sur C:, puis *ouvrir*. Clic-droit dans la fenêtre, puis *Nouveau/dossier*.
2. Ouvrez l'éditeur de texte Textpad (icône sur le bureau ou menu *Démarrer/Tous les programmes*).
3. Saisissez quelques phrases et enregistrez le fichier dans le répertoire C:\SL0720.
4. Ouvrez Firefox et saisissez l'url ci-dessous dans la barre d'adresse :  
<http://cental.fltr.ucl.ac.be/treetagger/>  
(ou cherchez *treetagger cental louvain* dans un moteur de recherche);
5. Cliquez sur le bouton *Parcourir* en face de *texte à étiqueter*, puis sélectionnez le fichier texte que vous avez créé dans la boîte de dialogue qui s'affiche et validez.
6. Cliquez sur envoyer (ou *submit query*).
7. Cliquez sur le lien *Télécharger le fichier étiqueté* (clic-droit pour enregistrer le fichier sur votre machine, clic-gauche pour l'afficher dans le navigateur).
8. Lorsque vous affichez le résultat dans votre navigateur, il est possible que les caractères comportant des diacritiques s'affichent mal. Dans ce cas, allez dans le menu *affichage* de votre navigateur, puis *Encodage des caractères/Occidental ISO-8859-1*.
9. Observez le résultat : de quelles informations dispose-t-on en plus du texte initial? Quelles sont, selon vous, les étapes préalables à l'étiquetage morpho-syntaxique?

## 2 Installation de l'interface

1. En salle machine, TreeTagger doit être installé dans le répertoire C:\TreeTagger. Si ce n'est pas le cas, installez-le en vous reportant à l'annexe D.

2. Affichez dans un navigateur la page suivante :  
<http://www.smo.uhi.ac.uk/~oduibhin/oideasra/interfaces/winttinterface.htm>
3. Dans l’item *Graphical interface* de la section *Links*, cliquez sur *Download the Windows interface to the tagger program*.
4. Enregistrez le fichier `.exe` dans le répertoire `C:\TreeTagger\bin`.
5. Créez un raccourci vers `wintreetagger.exe` depuis le bureau (cliquer-glisser avec le bouton droit de la souris depuis le fichier `wintreetagger.exe` vers le bureau, puis *créer un raccourci* . . .)

## 3 Étiquetage

### 3.1 Prise en main

1. Dans TextPad, créez un document texte et saisissez quelques phrases. Enregistrez le document.
2. Dans l’interface WinTreeTagger, sélectionnez la bonne langue, puis le fichier d’entrée (cliquez dans la zone de texte sous *Input File*).
3. Sous *lexical information*, cochez *none*.
4. Pour indiquer le fichier de sortie (résultat de l’étiquetage), cliquez dans la zone de texte sous *Output File*, puis :
  - sélectionnez un fichier existant (dans ce cas, le fichier sera écrasé) ;
  - ou entrez manuellement le nom d’un nouveau fichier (en le suffixant par `.txt`).
5. Cliquez sur *Run*.
6. Ouvrez le fichier résultat.

### 3.2 Casimir

Téléchargez le fichier `gloubi-boulga.txt`<sup>1</sup> et étiquetez-le.

Quels sont les formats particuliers pour les lemmes ?

Repérez les erreurs et classez-les en différents types et niveaux d’analyse. Repérez des incohérences dans l’étiquetage (des données similaires qui conduisent à des étiquetages différents).

Quels sont les différentes étiquettes pour les marques de ponctuations ? À quoi correspondent-elles ? Quels sont les cas particuliers dans la colonne des lemmes ?

Dans quels cas l’étiquetage des mots inconnus paraît correct/incorrect ? Cet étiquetage paraît-il cohérent (systématique) ?

### 3.3 Intervenir dans l’étiquetage

#### 3.3.1 Pré-étiquetage

Il est possible de fournir à TreeTagger un texte pour lequel certains éléments ont été pré-étiquetés : on impose ainsi que tel ou tel token ait, dans un contexte

<sup>1</sup>. Ce texte provient de l’article Wikipédia :  
<http://fr.wikipedia.org/wiki/Gloubi-boulga>

donné, tel ou tel lemme ou étiquette. Il faut pour cela produire, en entrée de TreeTagger, un fichier déjà segmenté (un token par ligne). Les lignes contenant les tokens que l'on veut pré-étiqueter doivent avoir le format suivant :

```
token    TAB    étiquette
```

ou

```
token    TAB    étiquette    TAB    lemme
```

(TAB représente une tabulation).

Reprendre la phrase « *Des soirées destinées aux adolescents combinant cosplay, rediffusion de dessins animés des années 1980 et fête étaient appelées Gloubi-boulga* » de l'exemple précédent, tokenisez-la à la main (*i.e.* mettez-la dans un fichier au format un token par ligne)<sup>2</sup>. Attention, dans le format "un token par ligne", les tokens ne doivent être ni précédés, ni suivis d'espace. Étiquetez-là avec l'option suivante : *Tokenization Options/none*. Répétez l'opération en faisant en sorte que *dessins animés* constitue un seul token. Observez le résultat. Ajoutez, dans le fichier d'entrée, l'étiquette et le lemme de *dessins animés* (chaque information étant séparée par une tabulation). Cochez *Input options/Manual tags have lemma*. Observez le résultat.

### 3.3.2 Lexique auxiliaire

La documentation de TreeTagger (cf. fichier `README.txt`, à la racine du répertoire `TreeTagger`) précise qu'il est possible d'utiliser un lexique auxiliaire dont elle précise le format. Un extrait de la documentation est reproduit ci-dessous :

Further optional command line arguments:

```
* -lex <f>: The file <f> contains additional lexicon entries to be used
by the tagger. The file format is identical to the format of the lexicon
argument of the training program (see below).
```

```
<f>: name of a file which contains the fullform lexicon. Each line
of the lexicon corresponds to one word form and contains the word form
and a sequence of tag-lemma pairs. Each tag is preceded by a tab character
and each lemma is preceded by a blank or tab character.
```

Example:

```
aback RB aback
abacuses NNS abacus
abandon VB abandon VBP abandon
abandoned JJ abandoned VBD abandon VBN abandon
abandoning VBG abandon
```

Considérez l'exemple précédent (version non tokenisée) et créez un fichier lexique qui contienne les informations relatives aux néologismes présents dans la phrase. Pour étiqueter la version non tokenisée, sélectionnez à nouveau *builtin* dans *Tokenization Options*. Pour utiliser le lexique que vous venez de créer, cochez, dans *Tagging Options*, la case *Auxiliary lexicon*, cliquez dans le champ

---

2. Il est possible également d'utiliser le tokenizer fourni avec TreeTagger : reportez-vous pour cela l'annexe E.

texte à droite et sélectionnez votre fichier lexique. Lancez l'étiquetage et observez.

### 3.3.3 Pré-étiquetage et lexique auxiliaire

La documentation concernant l'utilisation d'un lexique auxiliaire précise que l'on ne peut pas inclure dans ce lexique d'unité polylexicale telle que *dessin animé*. Pourquoi, selon vous ?

Reprenez le fichier tokenisé de la section 3.3.1 et étiquetez-là en utilisant conjointement le lexique créé à la section 3.3.2. Veillez à sélectionner dans l'interface les options adéquates de format d'entrée, de tokenisation et d'étiquetage.

## 3.4 Étiquetage en anglais

Téléchargez le fichier `holy_grail.txt`<sup>3</sup> et étiquetez-le.

Comparez le jeu d'étiquettes pour l'anglais à celui du français. Comparez également les types d'erreurs comises entre les deux langues.

## 4 Comparaison d'étiqueteurs

Téléchargez le fichier `dormeur_duval.txt` et étiquetez-le avec `TreeTagger`<sup>4</sup>.

1. Les deux premiers vers étiquetés par `Cordial` sont représentés figure 1. Quelles informations `Cordial` apporte-t-il en plus de `TreeTagger` ?
2. La version étiquetée (forme, lemme et partie du discours) par `Cordial` analyseur du *dormeur du val* est donnée dans l'annexe C. Quelle(s) différence(s) observez-vous ?
3. Relevez des erreurs d'étiquetage (ou résultats inattendus) pour chaque étiqueteur et comparez avec l'étiquetage de l'autre.

## 5 Étiquetage de textes balisés

Téléchargez le fichier « Extrait du corpus Air France » et décompressez-le. Visualisez-le dans `TextPad`, puis étiquez-le. Observez le résultat. Quel est le problème ?

Recommencez l'opération en cochant *Input Options/SGML tags present* et observez la différence.

---

3. Extrait issu du site *Unofficial Monty Python Home Page* :  
<http://www.muscomp.com/python.html>

4. Cet exemple est emprunté à Benoît Habert, *Instruments et ressources électroniques pour le français*. Ophrys (« L'Essentiel Français »), 2005

```

===== DEBUT DE PHRASE =====
1 C' ce PDS Pd-.n 1|1 S 1 est Principale
2 est être VINDP3S Vmip3s 2 V 1 est Principale
3 un un DETIMS Da-ms-i 4|4 B 1 est Principale
4 trou trou NCMS Ncms 4|4 B 1 est Principale
5 de de PREP Sp 6|4 B 1 est Principale
6 verdure verdure NCFS Ncfs 6|4 B 1 est Principale
7 où où PRI Ptr-.- 7|7 - 2 chante Relative
8 chante chanter VINDP3S Vmip3s 8 V 2 chante Relative
9 une un DETIFS Da-fs-i 10|10 T 2 chante Relative
10 rivière rivière NCFS Ncfs 10|10 T 2 chante Relative
===== FIN DE PHRASE =====

===== DEBUT DE PHRASE =====
1 Accrochant accrocher VPARPRES Vmpp-- 1 - 1 Indépendante suspendre/hang, suspend
2 follement follement ADV Rgp - 1 Indépendante
3 aux à le DETDP3G Da-p-d 4|4 - 1 Indépendante
4 herbes herbe NCFP Ncfp 4|4 - 1 Indépendante plante/grass
5 des de le DETDP3G Da-p-i 6|4 - 1 Indépendante
6 haillons haillon NCMP Ncmp 6|4 - 1 Indépendante
===== FIN DE PHRASE =====

```

51

FIGURE 1 – *Le dormeur du val* étiqueté par Cordial (deux premiers vers)

## A Jeu d'étiquettes pour le Français

ABR	abreviation
ADJ	adjective
ADV	adverb
DET:ART	article
DET:POS	possessive pronoun (ma, ta, ...)
INT	interjection
KON	conjunction
NAM	proper name
NOM	noun
NUM	numeral
PRO	pronoun
PRO:DEM	demonstrative pronoun
PRO:IND	indefinite pronoun
PRO:PER	personal pronoun
PRO:POS	possessive pronoun (mien, tien, ...)
PRO:REL	relative pronoun
PRP	preposition
PRP:det	preposition plus article (au, du, aux, des)
PUN	punctuation
PUN:cit	punctuation citation
SENT	sentence tag
SYM	symbol
VER:cond	verb conditional
VER:futu	verb futur
VER:impe	verb imperative
VER:impf	verb imperfect
VER:infi	verb infinitive
VER:pper	verb past participle
VER:ppre	verb present participle
VER:pres	verb present
VER:simp	verb simple past
VER:subi	verb subjunctive imperfect
VER:subp	verb subjunctive present

## B Jeu d'étiquettes pour l'Anglais : Penn Treebank Tagset

CC	Coordinating conjunction
CD	Cardinal number
DT	Determiner
EX	Existential there
FW	Foreign word
IN	Preposition or subordinating conjunction
JJ	Adjective
JJR	Adjective, comparative
JJS	Adjective, superlative
LS	List item marker
MD	Modal
NN	Noun, singular or mass
NNS	Noun, plural
NP	Proper noun, singular
NPS	Proper noun, plural
PDT	Predeterminer
POS	Possessive ending
PP	Personal pronoun
PP\$	Possessive pronoun
RB	Adverb
RBR	Adverb, comparative
RBS	Adverb, superlative
RP	Particle
SYM	Symbol
TO	to
UH	Interjection
VB	Verb, base form
VBD	Verb, past tense
VBG	Verb, gerund or present participle
VCN	Verb, past participle
VBP	Verb, non-3rd person singular present
VBZ	Verb, 3rd person singular present
WDT	Wh-determiner
WP	Wh-pronoun
WP\$	Possessive wh-pronoun
WRB	Wh-adverb

## C *Le dormeur du val* étiqueté par Cordial analyseur

==== DEBUT DE PHRASE ====	de	de	PREP	
C' ce PDS	rayons	rayon	NCMP	
est être VINDP3S	.	.	PCTFORTE	
un un DETIMS	==== FIN DE PHRASE =====			
trou trou NCMS	==== DEBUT DE PHRASE =====			
de de PREP	Un un DETIMS			
verdure verdure NCFS	soldat soldat NCMS			
où où PRI	jeune jeune ADJSIG			
chante chanter VINDP3S	, , PCTFAIB			
une un DETIFS	bouche bouche NCFS			
rivière rivière NCFS	ouverte ouvert ADJFS			
==== FIN DE PHRASE =====	, , PCTFAIB			
==== DEBUT DE PHRASE =====	tête nue tête nue			ADV
Accrochant accrocher VPARPRES	, , PCTFAIB			
follement follement ADV	==== FIN DE PHRASE =====			
aux à le DETDPIG	==== DEBUT DE PHRASE =====			
herbes herbe NCFP	Et et COO			
des de le DETDPIG	la le DETDFS			
haillons haillon NCMP	nuque nuque NCFS			
==== FIN DE PHRASE =====	baignant baigner VPARPRES			
==== DEBUT DE PHRASE =====	dans dans PREP			
D' de PREP	le le DETDMS			
argent argent ADJINV	frais frais ADJMIN			
; ; PCTFORTE	cresson cresson NCMS			
==== FIN DE PHRASE =====	bleu bleu ADJINV			
==== DEBUT DE PHRASE =====	, , PCTFAIB			
où où PRI	==== FIN DE PHRASE =====			
le le DETDMS	==== DEBUT DE PHRASE =====			
soleil soleil NCMS	Dort dormir VINDP3S			
, , PCTFAIB	; ; PCTFORTE			
de de PREP	==== FIN DE PHRASE =====			
la le DETDFS	==== DEBUT DE PHRASE =====			
montagne montagne NCFS	il il PPER3S			
fière fier ADJFS	est être VINDP3S			
, , PCTFAIB	étendu étendre VPARPMS			
==== FIN DE PHRASE =====	dans dans PREP			
==== DEBUT DE PHRASE =====	l' le DETDFS			
Luit luire VINDP3S	herbe herbe NCFS			
: : PCTFORTE	, , PCTFAIB			
==== FIN DE PHRASE =====	sous sous PREP			
==== DEBUT DE PHRASE =====	la le DETDFS			
c' ce PDS	nue nue NCFS			
est être VINDP3S	, , PCTFAIB			
un un DETIMS	==== FIN DE PHRASE =====			
petit petit ADJMS	==== DEBUT DE PHRASE =====			
val val NCMS	Pâle pâle ADJSIG			
qui qui PRI	dans dans PREP			
mousse mousser VINDP3S	son son DETPOSS			



lit	lit	NCMS	.	.	PCTFORTE	
vert	vert	ADJMS	=====	FIN DE PHRASE	=====	
où	où	PRI	=====	DEBUT DE PHRASE	=====	
la	le	DETDIFS	Les	le	DETDPIG	
lumière	lumière	NCFS	parfums	parfum	NCMP	
pleut	pleuvoir	VINDP3S	ne	ne	ADV	
.	.	PCTFORTE	font	faire	VINDP3P	
=====	FIN DE PHRASE	=====	pas	pas	ADV	
=====	DEBUT DE PHRASE	=====	frissonner	frissonner	VINF	
Les	le	DETDPIG	sa	son	DETPOSS	
pieds	pied	NCMP	narine	nar	NCFS	
dans	dans	PREP	;	;	PCTFORTE	
les	le	DETDPIG	=====	FIN DE PHRASE	=====	
glaïeuls	glaïeul	NCMP	=====	DEBUT DE PHRASE	=====	
,	,	PCTFAIB	Il	il	PPER3S	
il	il	PPER3S	dort	dormir	VINDP3S	
dort	dormir	VINDP3S	dans	dans	PREP	
.	.	PCTFORTE	le	le	DETDMS	
=====	FIN DE PHRASE	=====	soleil	soleil	NCMS	
=====	DEBUT DE PHRASE	=====	,	,	PCTFAIB	
Souriant	sourire	VPARPRES	la	le	DETDIFS	
comme	comme	SUB	main	main	NCFS	
=====	FIN DE PHRASE	=====	sur	sur	PREP	
=====	DEBUT DE PHRASE	=====	sa	son	DETPOSS	
Sourirait	sourire	VCONP3S	poitrine	poitrine	NCFS	
un	un	DETIMS	=====	FIN DE PHRASE	=====	
enfant	enfant	NCSIG	=====	DEBUT DE PHRASE	=====	
malade	malade	ADJSIG	Tranquille	Tranquille	NPMS	
,	,	PCTFAIB	.	.	PCTFORTE	
il	il	PPER3S	=====	FIN DE PHRASE	=====	
fait un somme	sourire	VINDP3S	=====	DEBUT DE PHRASE	=====	
:	:	PCTFORTE	Il	il	PPER3S	
=====	FIN DE PHRASE	=====	a	avoir	VINDP3S	
=====	DEBUT DE PHRASE	=====	deux	deux	ADJNUM	
Nature	Nature	NPMS	trous	trou	NCMP	
,	,	PCTFAIB	rouges	rouge	ADJPIG	
berce	bercer	VIMPP2S	au	à le	DETDMS	
-le	le	PPER3S	côté	côté	NCMS	
chaudemement	chaudemement	ADV	droit	droit	ADJMS	
:	:	PCTFORTE	.	.	PCTFORTE	
=====	FIN DE PHRASE	=====	=====	FIN DE PHRASE	=====	
=====	DEBUT DE PHRASE	=====	=====	DEBUT DE PHRASE	=====	
il	il	PPER3S	=====	FIN DE PHRASE	=====	
a froid avoir		VINDP3S				

## D Installation de TreeTagger

### 1. Téléchargement de TreeTagger

Ouvrez un navigateur web et saisissez dans la barre d'adresse l'URL :

<http://www.ims.uni-stuttgart.de/projekte/corplex/TreeTagger/>  
ou cherchez "treetagger" dans un moteur de recherche.

Trouvez le lien *Windows version* (entre les sections *Parameter files* et *tagsets*, ou Ctrl+F et *Windows version* dans la barre de recherche) et effectuez un clic-droit, puis sélectionnez *enregistrer la cible du lien sous* dans le menu contextuel. Enregistrez le fichier `tree-tagger-windows-3.1.zip` dans le répertoire C:\.

### 2. Extraction de l'archive TreeTagger

Effectuez un clic-droit sur le fichier `.zip` et sélectionnez dans le menu contextuel *7-zip/extract here*.

Le répertoire C:\TreeTagger contient à présent un sous-répertoire TreeTagger qui contient lui-même les sous-répertoires `bin`, `cmd` et `lib`.

### 3. Téléchargement des fichiers de paramètres

Pour utiliser TreeTagger avec une langue donnée, il faut lui fournir un fichier de paramètres correspondant. Dans la section *Parameter files for PC* de la page web de TreeTagger, effectuez un clic-droit sur les liens *English parameter file* et *French parameter file (Latin1)* et enregistrez les fichiers dans C:\TreeTagger\lib.

### 4. Extraction des fichiers de paramètres

Le répertoire C:\TreeTagger\lib doit contenir les fichiers suffixés `.bin.gz`. Décompressez chacun de ces fichiers par un clic-droit, puis *7-zip/extract here*. Une fois décompressés, les fichiers portent l'extension `.par`. Renommez-les (clic-droit, puis *Renommer* ou sélection, puis F2) respectivement en `french.par` et `english.par` s'ils portent des noms différents.

## E Utiliser séparément le segmenteur de TreeTagger

TreeTagger fournit un programme qu'il utilise pour *tokeniser* les textes avant de les étiqueter. La procédure pour segmenter un texte est décrite ci-dessous. On suppose que TreeTagger est installé dans le répertoire C:\TreeTagger, que le texte à segmenter est C:\SL0720\texte.txt et que l'on veut produire le résultat de la segmentation dans le fichier C:\SL0720\texte.tok.

1. Ouvrir l'invite de commande : sous Windows, bouton *Démarrer/Exécuter*, puis saisir `cmd` et validez avec la touche *Entrée*.
2. Saisissez la commande suivante :  
`perl c:\TreeTagger\cmd\tokenize.pl -f c:\720\texte.txt > c:\720\texte.tok`
3. Pour l'anglais, remplacer l'option `-f` par `-e`