

# Étiquetage morpho-syntaxique

Benoît Habert

LIMSI – CNRS & université Paris X–Nanterre

`http ://www.limsi.fr/Individu/habert - habert [ @ ]  
limsi.fr`

Cours PluriTAL 12/01/2006

# Plan

- **Exemples d'étiquetage morpho-syntaxique**
- **Forme des étiquetages**
- **Ajuster traitement et données**

# Étiqueter : dimensions

[Ide & Romary 03]

- « Répartitoire » (*annotation scheme*) : ensemble des catégories distinguées ; définition en extension/intension
- Choix d'étiquettes (*encoding scheme*) : représentation de surface des catégories ('Nom commun masculin singulier'  $\equiv$  NCMS)
- Mode d'annotation (*data architecture*) : réalisation « pratique »

# Le Dormeur du Val

1 C'est un trou de verdure où chante une rivière  
2 Accrochant follement aux herbes des haillons  
3 D'argent ; où le soleil, de la montagne fière,  
4 Luit : c'est un petit val qui mousse de rayons.

5 Un soldat jeune, bouche ouverte, tête nue,  
6 Et la nuque baignant dans le frais cresson bleu,  
7 Dort ; il est étendu dans l'herbe, sous la nue,  
8 Pâle dans son lit vert où la lumière pleut.

9 Les pieds dans les glaïeuls, il dort. Souriant comme  
10 Sourirait un enfant malade, il fait un somme :  
11 Nature, berce-le chaudement : il a froid.

12 Les parfums ne font pas frissonner sa narine ;  
13 Il dort dans le soleil, la main sur sa poitrine  
14 Tranquille. Il a deux trous rouges au côté droit.

# Étiquetage avec *Cordial* (1/2)

m	§	P	d	lemme	n	Typegram	Code	Codegram	S
o		h	i		b		Gram		y
t		r	a						n
		a	l		a				t
		s	o		m				a
		e	g		b				g
			u		i				m
			e		g.				e

==== DEBUT DE PHRASE ====

1	1	Le	1	le	A2	DETDMS	0xA000	D a - m s - d	2
1	1	Dormeur	2	dormeur	A2	NCMS	0xA000	N c m s	2
1	1	du	3	du		DETDMS	0xA000	D a - m s - d	4
1	1	Val	4	val		NCMS	0xA000	N c m s	4

==== FIN DE PHRASE ====

==== DEBUT DE PHRASE ====

\r\r

# Étiquetage avec *Cordial* (2/2)

P	...	N	P	Type	Sens
h	...	u	i	prop.	du
r	...	m	v		mot
a	...		o		
s	...	p	t		
e	...	r			
		o			
		p			
		.			

==== DEBUT DE PHRASE ====

Le	...	1		Indépendante	
Dormeur	...	1		Indépendante	personne
du	...	1		Indépendante	
Val	...	1		Indépendante	

==== FIN DE PHRASE ====

==== DEBUT DE PHRASE ====

\r\r

# Étiquetage avec le *TreeTagger*

Le	DET :ART	le
Dormeur	NOM	dormeur
du	PRP :det	du
Val	NOM	val
C'	PRO :DEM	ce
est	VER :pres	être
un	DET :ART	un
trou	NOM	trou
de	PREP	de
verdure	NOM	verdure
où	PRO :REL	où
chante	VER :pres	chanter
une	DET :ART	un
rivière	NOM	rivière

# Étiquetage avec *Unitex/Nooj/Intex*

{Pâle,pâle.A+z1 :ms :fs}  
{dans,dans.PREP+z1}  
{son,son.N+z1 :ms}  
{son,son.DET+z1 :ms}  
{lit,lire.V+z1 :P3s}  
{lit,lit.N+z1 :ms}  
{vert,vert.N+z1 :ms}  
{vert,vert.A+z1 :ms}  
{où,où.PRO+z1}  
{la,le.PRO+z1 :3fs}  
{la,le.DET+z1 :fs}  
{la,la.N+z1 :ms :mp}  
{lumière,lumière.N+z1 :fs}  
{pleut,pleuvoir.V+z1 :P3s}

# Morphologie avec *Flemm/Dérif* (1/2)

⋮

Dormeur NOM :Nc-s– dormeur [ [ dormir VER] eur NOM] (dormeur/NOM, dormir/VER) "  
(Agent habituel - Auteur exceptionnel - Instrument) de dormir"

⋮

chante VER(pres) :Vmip1s–1 chanter [ chanter VER] (chanter/VER) "chanter"

chante VER(pres) :Vmip3s–1 chanter [ chanter VER] (chanter/VER) "chanter"

chante VER(pres) :Vmmp2s–1 chanter [ chanter VER] (chanter/VER) "chanter"

chante VER(pres) :Vmisp1s–1 chanter [ chanter VER] (chanter/VER) "chanter"

chante VER(pres) :Vmisp3s–1 chanter [ chanter VER] (chanter/VER) "chanter"

⋮

mousse NOM :Nc-s– mousse [ mousse NOM] (mousse/NOM) "mousse"

⋮

tête NOM :Nc-s– tête [ tête NOM] (tête/NOM) "tête"

nue NOM :Ncfs– nue [ [nue ADJ] NOM] (nue/NOM, nue/ADJ) "Entité dont la propriété vue  
comme saillante est d'être nue"

⋮

# Morphologie avec *Flemm/Dérif* (2/2)

⋮

⋮

Dort VER(pres) :Vmip3s-3 dormir [ dormir VER] (dormir/VER) "dormir"

⋮

étendu VER(pper) :Vmpps-sm- étendre [ é [tendre VER ] VER] (étendre/VER, tendre/VER)  
"tendre jusqu'au bout"

⋮

nue NOM :Ncfs- nue [ [nue ADJ] NOM] (nue/NOM, nue/ADJ) "Entité dont la propriété vue  
comme saillante est d'être nue"

⋮

# Étiquetages...

*Cordial* (30 étiquettes réalisées sur 309)

ADJFS (3 o.), ADJINV (3 o.), ADJMIN (1 o.), ADJMS (4 o.), ADJNUM (1 o.), ADJPIG (1 o.), ADJSIG (4 o.), ADV (4 o.), COO (1 o.), DETDFS (6 o.), DETDMS (6 o.), DETDPIG (5 o.), DETIFS (1 o.), DETIMS (5 o.), DETPOSS (3 o.), NCFP (1 o.), NCFS (13 o.), NCMP (5 o.), NCMS (10 o.), NCSIG (1 o.), NHMIN (1 o.), PCTFAIB (12 o.), PCTFORTE (12 o.), PDS (2 o.), PREP (11 o.), PRI (4 o.), SUB (1 o.), VIN (1 o.), VPARPMS (1 o.), VPARPRES (3 o.)

*TreeTagger* (19 étiquettes réalisées sur 59)

ADJ (11 o.), ADV (5 o.), DET :ART (19 o.), DET :POS (3 o.), KON (1 o.), NOM (33 o.), PRO :DEM (2 o.), PRO :PER (6 o.), PRO :REL (4 o.), PRP (11 o.), PRP :det (3 o.), PUN (18 o.), SENT (5 o.), VER :aux :pres (1 o.), VER :cond (1 o.), VER :infi (1 o.), VER :pper (1 o.), VER :ppre (3 o.), VER :pres (12 o.)

# Comparaison d'étiquetages (1/2)

v.	<i>Cordial</i>				<i>TreeTagger</i>			
	<i>mot</i>	<i>lemme</i>	<i>cat.</i>	<i>POS</i>	<i>mot</i>	<i>lemme</i>	<i>cat.</i>	<i>POS</i>
		⋮				⋮		
1	chante	chanter	VIND- P3S	V	chante	chanter	VER pres	: V
1	une	un	DETIFS	D	une	un	DET ART	: D
1	rivière	rivier	NCFS	N	rivière	<b>rivière</b>	NOM	N
		⋮				⋮		
3	argent	argent	ADJINV	A	argent	argent	<b>NOM</b>	N
		⋮				⋮		
4	mousse	<b>mousser</b>	<b>VIND- P3S</b>	V	mousse	mousse	NOM	N
		⋮				⋮		
5	tête	tête	NCFS	N	tête	tête	NOM	N
5	nue	<b>nu</b>	<b>ADJFS</b>	A	nue	nue	NOM	N
		⋮				⋮		

# Comparaison d'étiquetages (2/2)

v.	<i>Cordial</i>				<i>TreeTagger</i>			
	<i>mot</i>	<i>lemme</i>	<i>cat.</i>	<i>POS</i>	<i>mot</i>	<i>lemme</i>	<i>cat.</i>	<i>POS</i>
		⋮				⋮		
9	comme	<b>comme</b>	<b>SUB</b>	C	comme	comme	ADV	R
		⋮				⋮		
11	Nature	nature	ADJINV	A	Nature	nature	<b>NOM</b>	N
11	berce	bercer	<b>VIMP- P2S</b>	V	berce	bercer	VER	: V
		⋮				⋮	pres	

# Plan

- Exemples d'étiquetage morpho-syntaxique
- **Forme des étiquetages**
- Ajuster traitement et données

# Conventions implicites vs. explicites

## TreeTagger

Le → DET :ART → le¶  
Dormeur → NOM → dormeur¶  
du → PRP :det → du¶  
Val → NOM → val¶

Les sorties de TreeTagger obéissent à un *format délimité*, du type des tables de SGBD (csv : *comma separated values*) : une entité ≡ une ligne (**n** ou **\r\n**) ; une facette ≡ une colonne (**t**)

	<i>forme</i>	<i>catégorie</i>	<i>lemme</i>
<i>mot</i> <sub>1</sub>	Le	DET : ART	le
<i>mot</i> <sub>2</sub>	Dormeur	NOM	dormeur
<i>mot</i> <sub>3</sub>	du	PRP :det	du
<i>mot</i> <sub>4</sub>	Val	NOM	val

# Notation idiosyncrasique

m	§	P	d	lemme	n	Typegram	Code	Codegram	S
o		h	i		b		Gram		y
t		r	a						n
		a	l		a				t
		s	o		m				a
		e	g		b				g
			u		i				m
			e		g.				e

==== DEBUT DE PHRASE ====

1	1	Le	1	le	A2	DETDMS	0xA000	D a - m s - d	2
1	1	Dormeur	2	dormeur	A2	NCMS	0xA000	N c m s	2
1	1	du	3	du		DETDMS	0xA000	D a - m s - d	4
1	1	Val	4	val		NCMS	0xA000	N c m s	4

==== FIN DE PHRASE ====

==== DEBUT DE PHRASE ====

\r\r

# Transformation *ad hoc*

<i>m</i>	<i>§</i>	<i>P</i>	<i>d</i>	<i>lemme</i>	<i>n</i>	<i>Typegram</i>	<i>Code</i>	<i>Codegram</i>	<i>S</i>
<i>o</i>		<i>h</i>	<i>i</i>		<i>b</i>		<i>Gram</i>		<i>y</i>
<i>t</i>		<i>r</i>	<i>a</i>						<i>n</i>
		<i>a</i>	<i>l</i>		<i>a</i>				<i>t</i>
		<i>s</i>	<i>o</i>		<i>m</i>				<i>a</i>
		<i>e</i>	<i>g</i>		<i>b</i>				<i>g</i>
			<i>u</i>		<i>i</i>				<i>m</i>
			<i>e</i>		<i>g.</i>				<i>e</i>

==== DEBUT DE PHRASE ====

<i>1</i>	<i>1</i>	<b>Le</b>	<i>1</i>	<b>le</b>	<i>A2</i>	<b>DETDMS</b>	<i>0xA000</i>	<b>D a - m s - d</b>	<i>2</i>
<i>1</i>	<i>1</i>	<b>Dormeur</b>	<i>2</i>	<b>dormeur</b>	<i>A2</i>	<b>NCMS</b>	<i>0xA000</i>	<b>N c m s</b>	<i>2</i>
<i>1</i>	<i>1</i>	<b>du</b>	<i>3</i>	<b>du</b>		<b>DETDMS</b>	<i>0xA000</i>	<b>D a - m s - d</b>	<i>4</i>
<i>1</i>	<i>1</i>	<b>Val</b>	<i>4</i>	<b>val</b>		<b>NCMS</b>	<i>0xA000</i>	<b>N c m s</b>	<i>4</i>

==== FIN DE PHRASE ====

==== DEBUT DE PHRASE ====

|r|r

# Format vs catégories

- TreeTagger

<w lemma="le" type="DET :ART">Le</w>

<w lemma="dormeur" type="NOM">dormeur</w>

<w lemma="du" type="PRP :det">du</w>

<w lemma="val" type="NOM">val</w>

- Cordial (jeu<sub>1</sub>)

<w lemma="le" type="DETDMS">Le</w>

<w lemma="dormeur" type="NCMS">dormeur</w>

<w lemma="du" type="DETDMS">du</w>

<w lemma="val" type="NCMS">val</w>

- Cordial (jeu<sub>2</sub>)

<w lemma="le" type="Da-ms-d">Le</w>

<w lemma="dormeur" type="Ncms">dormeur</w>

<w lemma="du" type="Da-ms-d">du</w>

<w lemma="val" type="Ncms">val</w>

# Choix dans les notations

- Notation idiosyncrasique vs. homologuée (*de facto* ou norme)
- Notation implicite vs. explicite
- Notation validable vs. sans validation explicite
- Notation traitant un niveau / plusieurs niveaux d'annotation
- Notation à plat | enchâssante | enchevêtrée

# Plan

- Exemples d'étiquetage morpho-syntaxique
- Forme des étiquetages
- **Ajuster traitement et données**

# Ajuster traitement et données (1/2)

- Régler un instrument / adapter les données
  - Exemple : *Dormeur* en phrases (vs vers)

Le dormeur du val

C' est un trou de verdure où chante une rivière accrochant follement aux herbes des haillons d' argent ; où le soleil , de la montagne fière , luit : c' est un petit val qui mousse de rayons . . . . Les parfums ne font pas frissonner sa narine ; il dort dans le soleil , la main sur sa poitrine tranquille . Il a deux trous rouges au côté droit .

# Ajuster traitement et données (2/2)

*vers*

...

rivière	rivier	NCFS
===== FIN DE PHRASE =====		
===== DEBUT DE PHRASE =====		
Accrochant	accrocher	VPARPRES
follement	follement	ADV
aux	au	DETDPIG
herbes	herbe	NCFP
des	de	DETDPIG
haillons	haillon	NHMIN
===== FIN DE PHRASE =====		
===== DEBUT DE PHRASE =====		
D'	de	PREP
argent	argent	<b>ADJINV</b>
;	;	PCTFORTE
===== FIN DE PHRASE =====		

*phrase*

...

rivière	rivier	NCFS
===== FIN DE PHRASE =====		
===== DEBUT DE PHRASE =====		
accrochant	accrochant	VPARPRES
follement	follement	ADV
aux	au	DETDPIG
herbes	herbe	NCFP
des	de	DETDPIG
haillons	haillon	NHMIN
===== FIN DE PHRASE =====		
===== DEBUT DE PHRASE =====		
d'	de	PREP
argent	argent	<b>NCMS</b>
;	;	PCTFORTE
===== FIN DE PHRASE =====		