

ANR 2006 MDCA
PASSAGE

Produire des annotations syntaxiques à grande échelle

<http://atoll.inria.fr/passage>

Éric de la Clergerie

Eric.De_La_Clergerie@inria.fr



Grand Colloque STIC 2007
6 Novembre 2007

AGENCE NATIONALE DE LA RECHERCHE
ANR

Améliorer le traitement linguistique du français

GP 1		NV 2		GN 3		GA 4			
A	quoi	servent	les	ressources	linguistiques	?			
1	2	3	4	5	6	7			
CPL-V			MOD-N						
	SUJ-V								

GP 1	NV 2	GN 3		GA 4		
A	quoi	servent	les	ressources	linguistiques	?
1	2	3	4	5	6	7
CPL-V			MOD-N			
		SUJ-V				

Traitement syntaxique du langage :

- pas de consensus sur les formalismes
- diversité des phénomènes syntaxiques
- les mots gouvernent la syntaxe
⇒ lexiques syntaxiques riches et couvrants : manque pour le français
- les corpus pour saisir les usages des mots
⇒ Tendance : corpus **manuellement** annotés syntaxiquement ([TreeBank](#)) pour amorcer le processus.

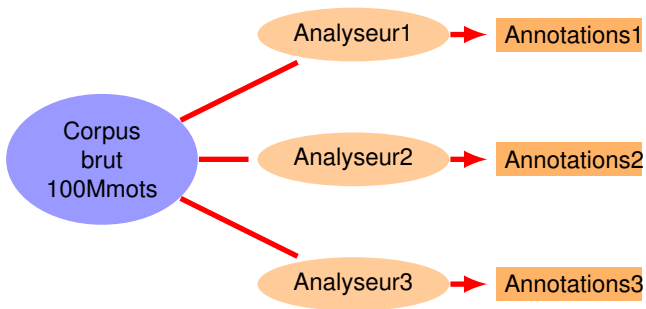
GP 1		NV 2		GN 3		GA 4		
A	quoi	servent	les	ressources	linguistiques	?		
1	2	3	4	5	6	7		
CPL-V			MOD-N					
SUJ-V								

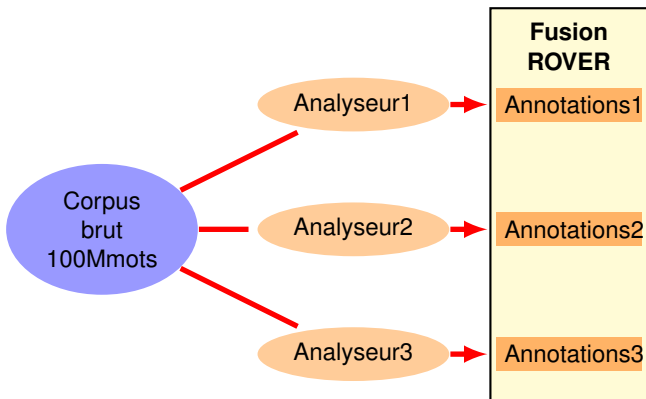
Traitement syntaxique du langage :

- pas de consensus sur les formalismes
- diversité des phénomènes syntaxiques
- les mots gouvernent la syntaxe
⇒ lexiques syntaxiques riches et couvrants : manque pour le français
- les corpus pour saisir les usages des mots
⇒ Tendance : corpus **manuellement** annotés syntaxiquement ([TreeBank](#)) pour amorcer le processus.

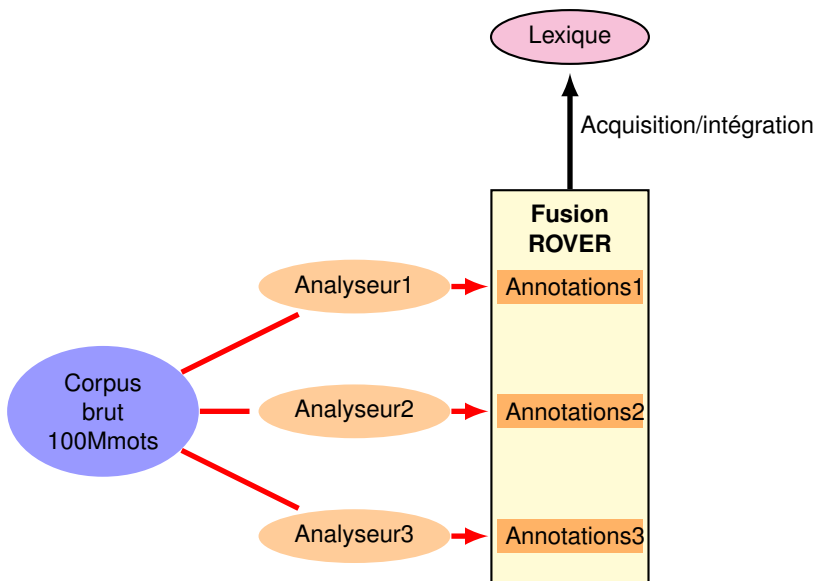
PASSAGE : proposition d'une approche plus intégrée

Corpus
brut
100Mmots

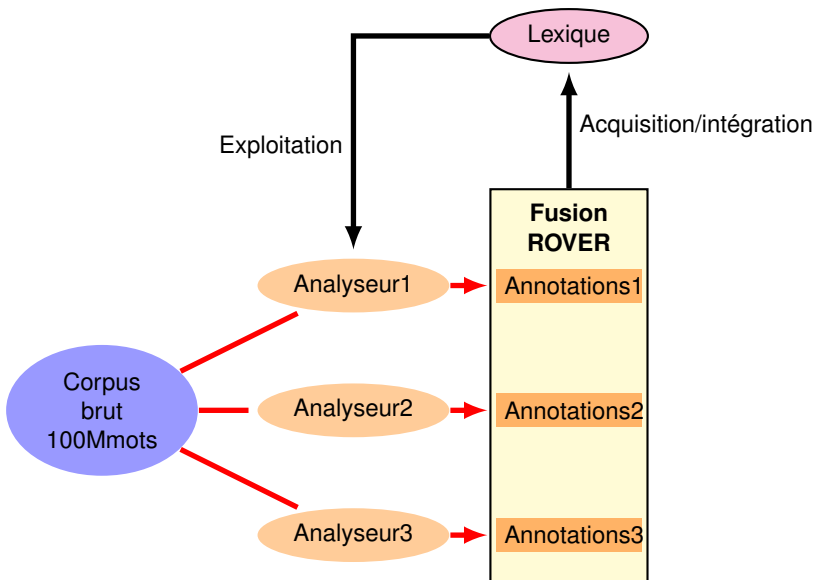




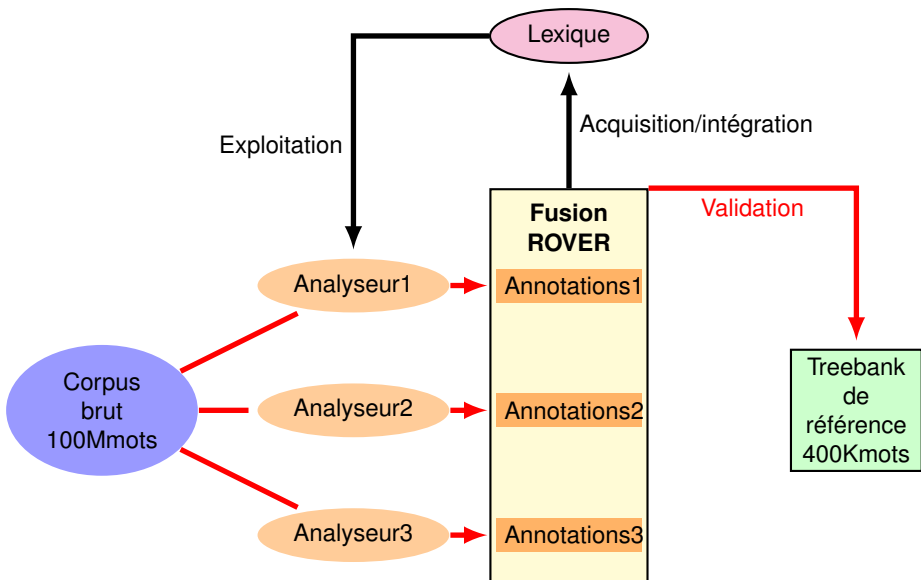
Établir un cercle vertueux entre outils et ressources



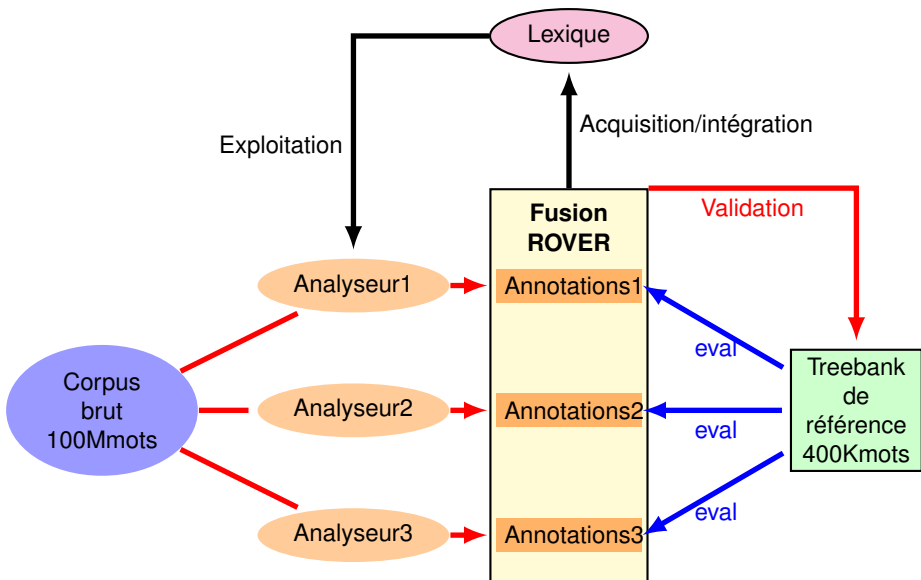
Établir un cercle vertueux entre outils et ressources



Établir un cercle vertueux entre outils et ressources



Établir un cercle vertueux entre outils et ressources



ALPAGE
INRIA Paris7



UNIVERSITÉ
PARIS
DIDEROT
PARIS 7

TALARIS/LORIA



LIC2M/CEA-LIST



LIR/LIMSI



ELDA



TAGMATICA

Tagmatica

ERSS



ALPAGE
INRIA Paris7



UNIVERSITÉ
PARIS
DIDEROT
PARIS 7

LPL



TALARIS/LORIA



LIC2M/CEA-LIST



SYNAPSE



LIR/LIMSI



XRCE

XEROX
Research Centre Europe

LIRMM



Faire coopérer une dizaine de parseurs

Une occasion unique, source de diversité (formalismes, technologies, ...)

Analyseur	Origine	Nature
FRMG	INRIA	TIG/TAG+DyALOG
SxLFG	INRIA	LFG+SYNTAX
LLP2	LORIA	TAG
LIMA	CEA-LIST	Système de règles
TAGPARSER	TAGMATICA	Induction + règles
SYNTAX	ERSS	Système de règles
GP1 & GP2	LPL	Grammaires de Propriétés
CORDIAL	SYNAPSE	
SYGMART	LIRMM	
XIP	XRCE	Cascade de règles

Exploiter de très grands corpus

Les corpus annotés syntaxiquement comme le **Penn Tree Bank** sont très importants pour le **Traitement Automatique des Langues [TAL]** mais **rare**s et **coûteux** à développer.

Exploiter de très grands corpus

Les corpus annotés syntaxiquement comme le **Penn Tree Bank** sont très importants pour le **Traitement Automatique des Langues [TAL]** mais **rare**s et **coûteux** à développer.

Par contre, il est maintenant possible d'accéder à de très grandes quantités de textes électroniques en français :

Corpus	Taille	Nature
Corpus EASy	1Mmots	multi-styles
Wikipedia Fr	~ 86Mmots	encyclopédique collaboratif libre
Wikisources	~ 80Mmots	littéraire libre
Monde Diplomatique	18Mmots	journalistique
FRANTEXT	20Mmots	littéraire libre
Europarl	28Mmots	débat Parlement européen
JRC-Acquis	39Mmots	juridique européen
Corpus Ester	1Mmots	oral transcrit
Total (actuel)	> 270 Mmots	

- Expertise de Technolanguage **EVALDA/EASy** sur l'évaluation des analyseurs du français :
 - ▶ **[Eval1]** Campagne Novembre 2007 (en cours)
 - ▶ **[Eval2]** Campagne fin 2009
- Format EASy :
chunks (GN, GA, NV, ...) + dépendances (Sujet-Verbe, Objet-Verbe, ...)
 - ▶ Existence guide d'annotation assez complet
 - ▶ Existence ~ 4K phrases annotées
 - ▶ Vers l'ajout des groupes récursifs (**constituance**)
 - ▶ Vers le respect des **standards** ISO TC37SC4.
- **ROVER** : Fusion des jeux d'annotations paramétrée par les évaluations + *feedback*
- Validation manuelle d'un sous-corpus

- Exploiter le ROVER pour acquérir des **connaissances lexicales** :
 - ▶ Information de valence (verbes, ...)
Sujet donne **Objet à-objet**
Sujet le lui donne
 - ▶ Classes d'alternation verbales (**Levin**)
Transfert de possession : vendre, céder, ...
 - ▶ Probabilités de désambiguation
catégories lexicales, constructions syntaxiques, restrictions de sélection, ...
 - ▶ Classes sémantiques (Hyp. distributionnelle **Harris**)
(**humain** ∨ **organisation**) annonce que Y
humain est nommé à la tête de **organisation**
 - ▶ Morphologie dérivationnelle avec transfert des informations syntaxiques
partir/départ ; X annoncer que Y //l'annonce de Y par X

- Les valider et intégrer dans un lexique

- Les exploiter dans certains analyseurs
⇒ Améliorer les analyseurs !

Besoin d'une infrastructure solide pour gérer des masses de données

- Utilisation de *fermes* de machines pour analyser de gros corpus par exemple GRID 5000 ?
- **EASYREF**, un prototype WEB de gestion collaborative d'annotations
visualisation, recherche, édition, *versionning*, comparaison, dépôt, fusion
 - ▶ Utilisé pour corriger les annotations EASy existantes
 - ▶ Utilisé pour annoter 500 nouvelles phrases
 - ▶ Démo à STIC 2007

► Rapports pour oral_delic_1:E96

Actions			Rld Bld	Corpus:Phrase:Rev	Auteur	Date	S C T
Del	Edit	Next	645	645 oral_delic_1 : E96 : r000	gil	2007-08-09 13:29:30	o
↶ Del	Edit	Next	979	645 oral_delic_1 : E96 : r000	nora	2007-10-15 12:30:13	x

Erreur sur relation: la relation ajout d'une relation COMP POST: COMP

New

► Erreurs potentielles détectées pour oral_delic_1:E96

► Corrections pour oral_delic_1:E96

Rev0001 created relation E96R4 with type 'COMP' and roles 'complementeur' => 'E96F6' 'verbe' => 'E96F10' -- nora -- lun

Annotations pour oral_delic_1:E96

S D ◀ ▶ Enoncé oral_delic_1 E96 -- Rev0001 -- Analyse complète FRMG											
			NV 1	GP 2			NV 3	GR 4	NV 5		GN 6
			NV 1	GP 2	GN 3		NV 4	GR 5	NV 6		GN 7
donc	c'	est	pour	ça	qu'	elle	a	pas	eu	son	C.A.P
1	2	3	4	5	6	7	8	9	10	11	12
COORD			MOD-N								
SUJ-V			COD-V								
ATB-SO			COMP								
CPL-V			SUJ-V			MOD-V			AUX-V		
			MOD-V			COD-V					

- Site Web : <http://atoll.inria.fr/passage>
- Liste de diffusion `passage@inria.fr`
- Zone wiki (**DOKUWIKI**)
 - ▶ Discussion et mise au point : corpus, protocoles, ...
 - ▶ Présentation et discussion : analyseurs
 - ▶ ...
- Zone INRIA GForge
 - ▶ Accès unifié aux logiciels et ressources
 - ▶ Versionning (SVN), gestionnaire de bugs, ...

Calendrier indicatif

2007				2008				2009	
	Corpus								Distrib
Format				Analyse Corpus					
		Eval1		Fusion Annotations					
				Validation Treebank					
				Acquisition					
							Intégration		
								Eval2	

- Rendre disponible des ressources linguistiques de qualité pour le français :
 - ▶ des corpus annotés, dont
 - ★ le **ROVER** (> 100 Mmots)
 - ★ le **TREEBANK** vérifié manuellement (400Kmots)
 - ▶ des **ressources lexicales**

Distribution : ELDA et CNRTL, en statut libre si possible

- Améliorer :
 - ▶ les analyseurs syntaxiques robustes efficaces du français
 - ▶ les méthodologies d'évaluation et de fusion des annotations
 - ▶ les méthodologies d'acquisition par bootstrap
- Permettre des applications par analyse syntaxique de corpus :
 - ▶ Acquisition de connaissances (extraction d'ontologies, ...)
 - ▶ Extraction d'informations (fouille de textes, veille, question-réponse, ...)

PASSAGE pour permettre d'amorcer d'autres projets :

- Corpus de Référence du Français (**FNC**), dans la lignée de **BNC** et **ANC**
- Expériences de transfert syntaxique entre langues à partir de corpus alignés comme **Europarl**, **JRC-Acquis** et **Wikipedia** ⇒ import de l'anglais, export vers d'autres langues européennes
- Acquisition de connaissances, à partir de **Wikipedia**

Discussions sur PASSAGE ⇒ Co-organisation ICGL'08 Workshop
«*Automated Syntactic Annotations for Interoperable Language Resources*»
Hong-Kong, Janvier 2008