

Du TAL dans les écrits scolaires : PREMIÈRES APPROCHES

Claire Wolfarth, Claude Ponton, Catherine Brissaud

Lidilem

Université Grenoble Alpes



Le projet Scoledit

I – Corpus

II – Travaux TAL

III – Perspectives

Contexte

- Travaux en écriture moins fréquents qu'en lecture (Boscolo, 2008)
- Développement des recherches sur l'écriture
- Ces recherches s'appuient sur des corpus souvent restreints, spécifiques ou difficilement accessibles (Elalouf, 2005 ; Auriac-Slusarczyk, Gunnarsson, 2014)
- *Enquête « Lire – Écrire CP » de l'Institut Français de l'Éducation (2013-2015)*

Objectif

- Constitution d'un corpus scolaire longitudinal en poursuivant le recueil du corpus jusqu'en CM2 (2018)

Enjeux

- Pour la linguistique :
 - Description en synchronie et en diachronie des caractéristiques des écrits d'apprenants de 6 à 12 ans
 - Meilleure connaissance des phénomènes d'acquisition de l'écriture
- Pour la didactique:
 - Développement de séquences à partir de ces connaissances

Corpus : Structure

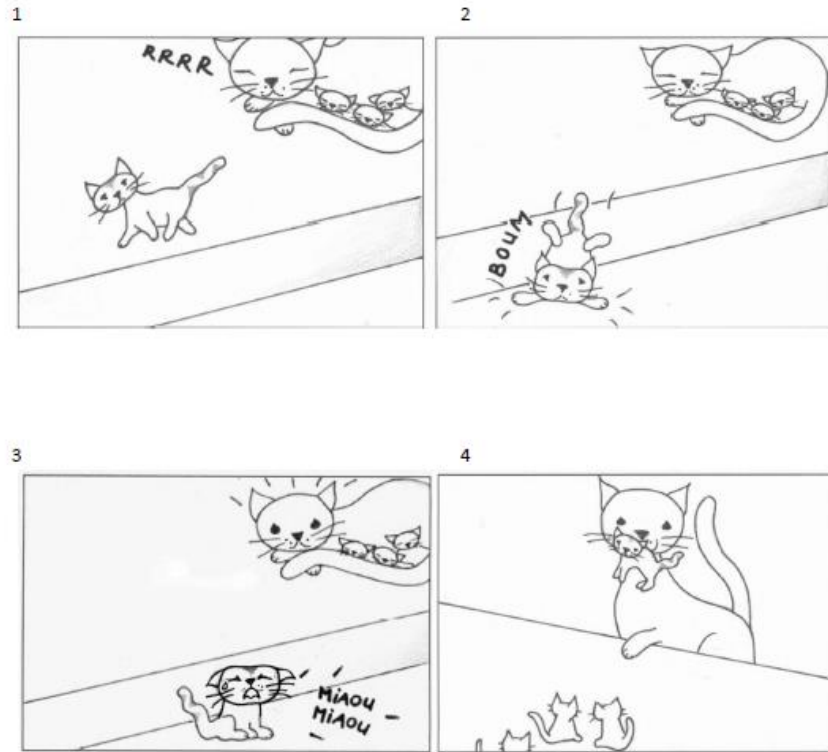
I – Corpus
 II – Travaux TAL
 III – Perspectives

CP	CE1	CE2	CM1	CM2
Dictées (Sept. 2013) 1151 dictées	Dictées (Juin 2015) 700-900 dictées	Productions de textes (Mai - Juin 2016) 900-1000 textes	Productions de textes (Mai - Juin 2017) 800-1100 textes attendus	Productions de textes (Mai - Juin 2018) 800-1100 textes attendus
Dictées (Juin 2014) 975 dictées				
Productions de textes (Juin 2014) 965 textes	Productions de textes (Juin 2015) 700-900 textes			

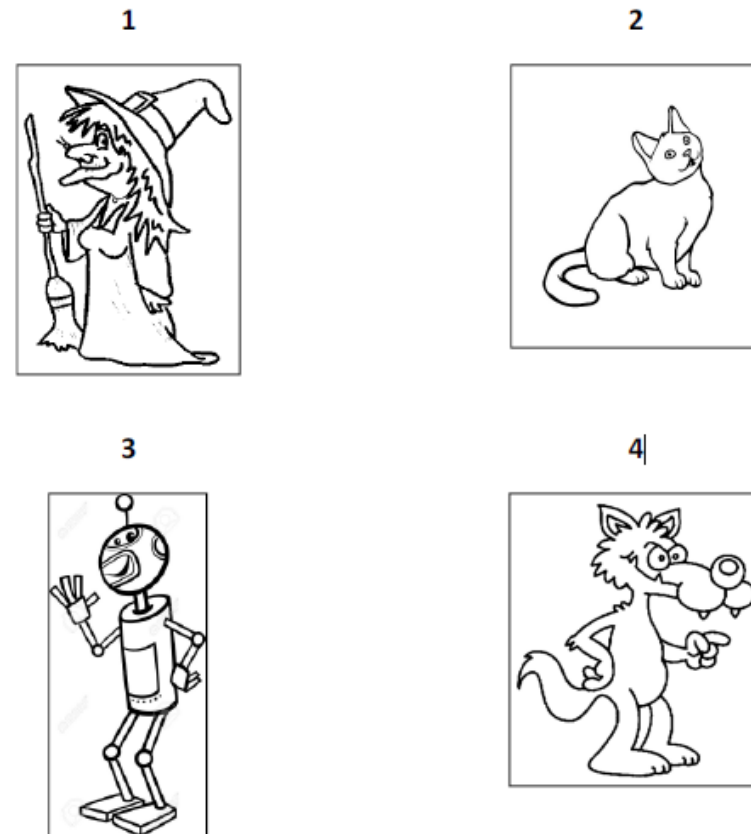
Corpus : Méthodologie

I - Corpus
II - Travaux TAL
III - Perspectives

Recueil CP (15 min)



Recueil CE1 (20 min) et CE2 (30 min)



Corpus : Données recueillies (CP)

(ex. élève 1558)

I - Corpus

II - Travaux TAL

III - Perspectives

Septembre 2013
(début CP)

AP
RA
É
TA



Juin 2014 (fin CP)

1. lapin
2. sa
3. éléphant
4. lors de l'attaque le sa
5. les la pinguouvent

le	chat	é	fati	gé	é	tombe	é	apré	
r	é	lle	sa	maman	é	cés	frère	é	sa
maman	le	ramèn	avèc	ses	frère				
é	an	site	il	dorme	avèc	ses	frère		

Corpus : Données recueillies (CE1)

(ex. élève 1558)

I - Corpus

II - Travaux TAL

III - Perspectives

Juin 2015 (fin CE1)

1/ patimp.

4/ récréations.

2/ patinon.

5/ charitable.

3/ capuchons.

6/ manifce

Les ~~Les~~ En été, les salade verte pousse
dans les jardins. Les jaunes canetons
picore le blé avec la poule noire.

Le loup se ~~promene~~ promene dans les bois
camp un chat ~~sur~~ sur des buisson et
lui di que faitu je me promene
pourquoi on peut se promene
ensemble et une bonne idet non pas
pourquoi pas suivre la riviere
il nia pas de chasseur sert sur
sinon il metir dessus et je
cour beaucoup mais que je
me cache dans un buisson
et il me chere et lire par
tout et je cour a nouveau

il faut faire attention il sont nombre
et met de pillege de partout et on
fait trais attention il met des
corde de par tout.

Corpus : Données recueillies (CE2)

(ex. élève 1558)

I - Corpus

II - Travaux TAL

III - Perspectives

Un chat blanc est noir pour la chasse
dans la forêt verte et marron. Il reste immobile
puis attaque sa proie vers un lièvre le chat
la ouvre et le ramène à son maître qui,
le récompense. Après le chat
joue avec les souris qui se cachent
derrière les arbres puis le chat fait une
neste quand il se réveille son maître fait
une balade dans la forêt. Le lendemain
il portera avec le bois de la forêt.

Juin 2016 (fin CE2)

La tante et des passants leur disent que fada
vous êtes il ya des loups - à bon souvenir
chez nous. Ben

Visualisation

Recherche par :

Forme

1558

Résultat de votre recherche :

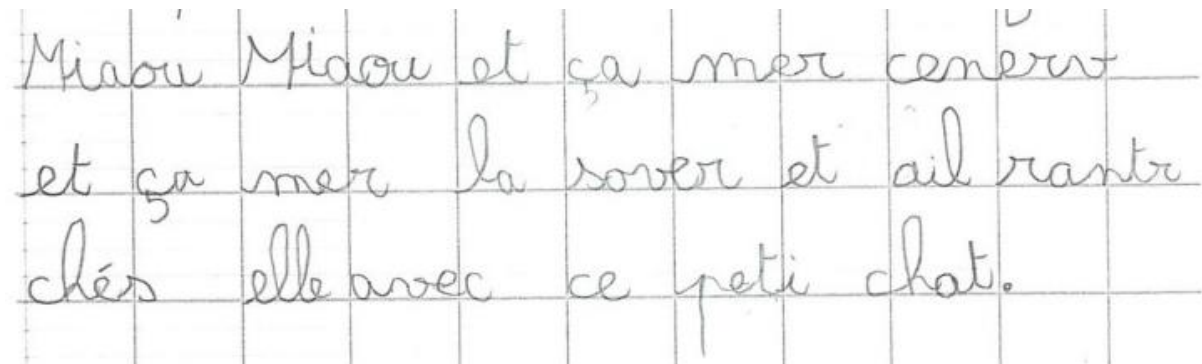
Il y a actuellement à l'étude : 1169 élèves.

Elève 1558

SCAN DICTÉE SEPTEMBRE	DICTÉE SEPTEMBRE	DICTÉE JUIN	SCAN DICTÉE JUIN
PRODUCTION	SCAN PRODUCTION		
<p>le chat éfati gé é tonbe éapré révéille sa maman é cés frère é sa maman le ramèn avec sec frère é an site il dorme avec séc frère</p>			

Spécificités du corpus : un défi pour le TAL

- Corpus manuscrit (transcription nécessaire)
- Corpus très « fautif »
- Longitudinal (évolutif)
- Taille du corpus
- Champ d'étude très récent en TAL



Miaou Miaou et ça mer cenéret
et ça mer la sover et ail rantr
chés elle avec ce peti chat.

Ex. élève 49 (CP)

TAL et corpus scolaires

- Hypothèses :
 - TAL comme aide à l'exploitation du corpus
 - TAL nourri par :
 - Le contexte de production
 - Les travaux dans différents domaines connexes
- Méthodologie :
 - Approche empirique
 - Appui sur les niveaux les plus éprouvés

TAL : État de l'art

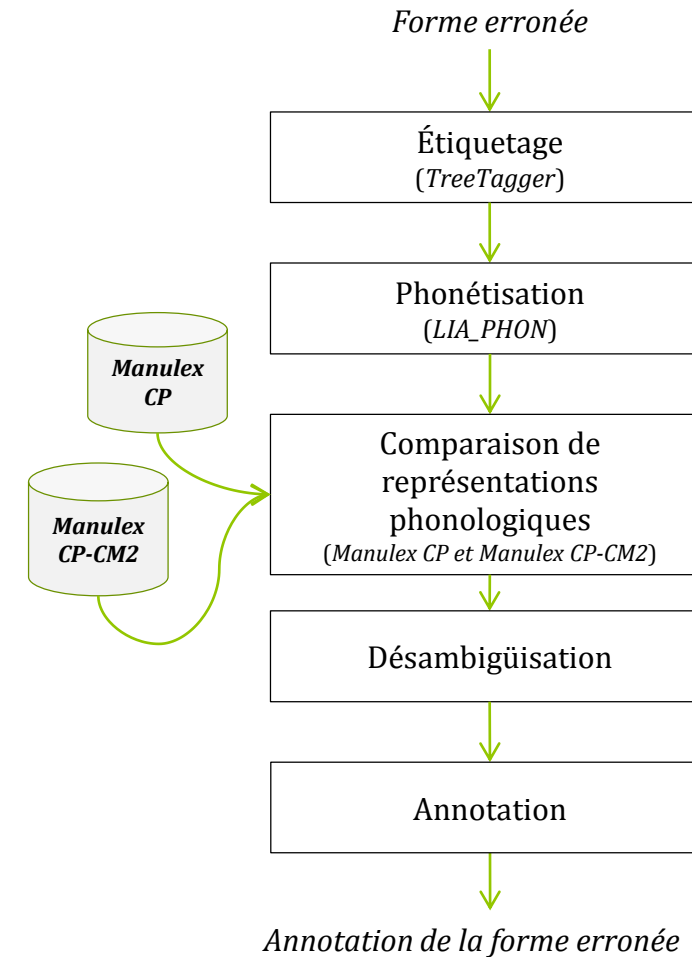
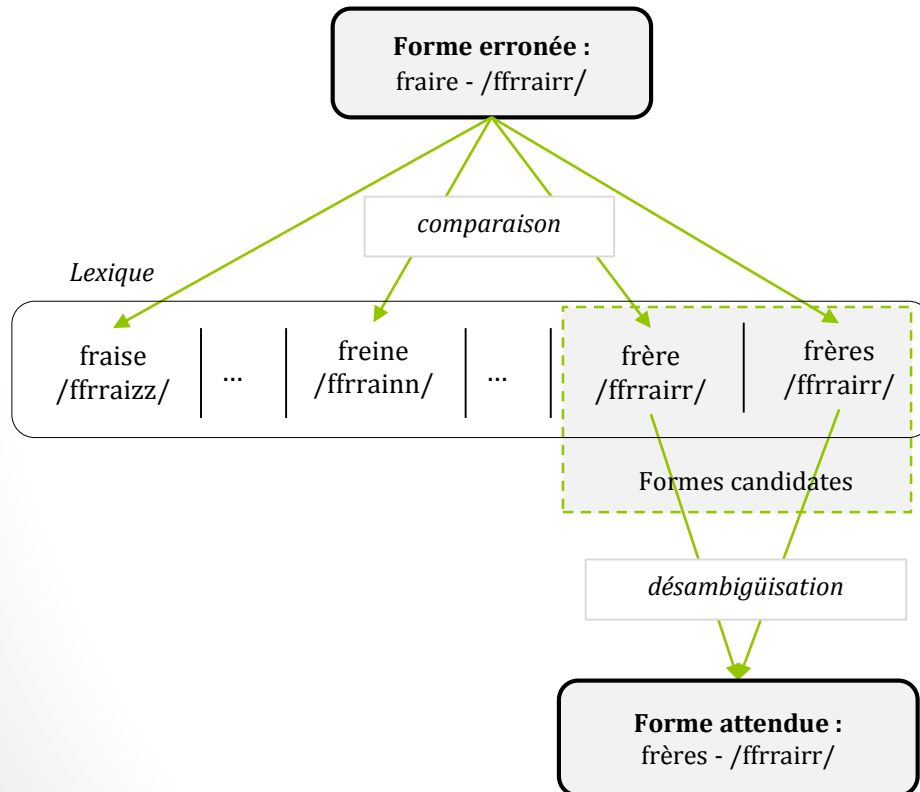
- Correction automatique de textes pour experts (Kukich, 1992)
- Normalisation de corpus peu normés
 - SMS, Tweets (Kobus *et al.*, 2008, Beaufort *et al.*, 2010, Baranes, 2012)
- TAL et apprentissage des langues
 - Freetext (Granger, Vandeventer, Hamel, 2001)
 - Exxelant (Antoniadis, Ponton, Zampa, 2010)

TAL : Particularités des productions de CP

- Écriture avec appui de l'oral
« tonb » (*tombe*)
- Segmentation en mots pas toujours maîtrisée
« Le petit chat sanva pan dan cesa maman dore. [...] » (*Le petit chat s'en va pendant que sa maman dort.*)
- Structures syntaxiques simples ou relativement simples
« le petit chas par et il sé fes male et mètñ il plere sa maman se fach » (*Le petit chat part et il s'est fait mal et maintenant il pleure. Sa maman se fâche.*)

TAL : Identification des formes normées par comparaison phonologique

Procédure d'identification de la forme attendue

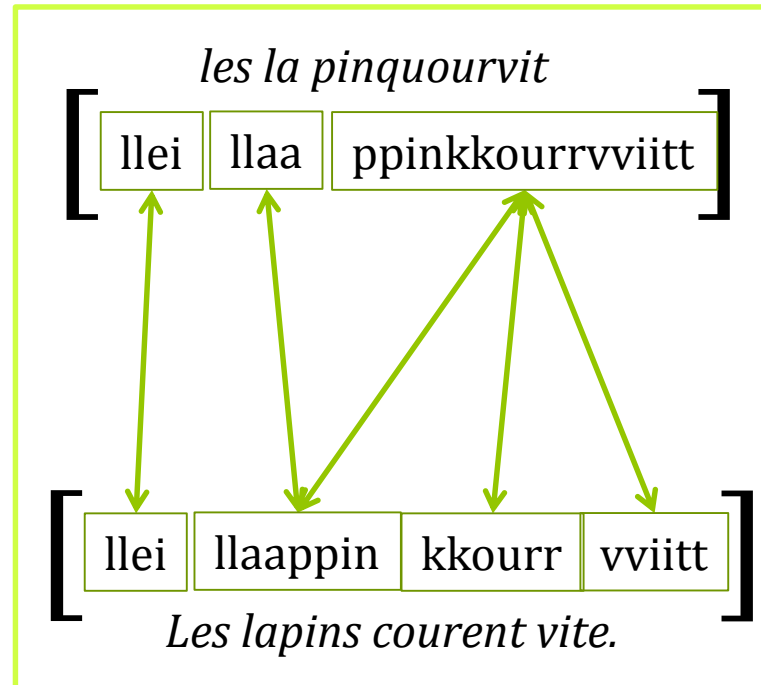


Processus d'annotation d'une forme erronée

TAL : Annotation par alignement avec un « corrigé »

Exemple

- « les la pinguourvit »
(*les lapins courent vite*)



Identification des formes attendues par alignement

Perspectives

- Poursuite du recueil en CM1 et en CM2
- Schéma d'annotation
- Alignement entre les productions et les « formes attendues »
- Réflexion sur le statut de la « forme attendue »
- Publication du corpus sur un site web à destination des enseignants et des chercheurs impliqués
 - Instaurer un dialogue entre les deux communautés
 - Améliorer la qualité des données
 - Développer des fonctionnalités adaptées
- Publication sur Ortolang

Merci pour votre attention

Références

- ANTONIADIS G., PONTON C., ZAMPA V. (2010). Exxelant et Mirto – Deux exemples d’environnement d’ALAO intégrant des outils TAL. *Multilinguisme et traitement des langues naturelles*. Montréal, Canada : PUQ.
- AURIAC-SLUSARCZYK E., GUNNARSON-LARGY C. (2014). *Écriture et réécritures chez les élèves: un seul corpus, divers genres discursifs et méthodologies d'analyse*. Academia.
- BARANES M. (2012). Vers la correction automatique de textes bruités: Architecture générale et détermination de la langue d'un mot inconnu. Actes de *RECITAL'2012 - Rencontre des Étudiants Chercheurs en Informatique pour le Traitement Automatique des Langues*, 95-108.
- BEAUFORT R., ROEKHAUT S., COUGNON L.-A. & FAIRON, C. (2010). Une approche hybride traduction/correction pour la normalisation des SMS, Actes de la 17e conférence sur le traitement automatique des langues naturelles (TALN'10), Montréal.
- BÉCHET F. (2001). LIA_PHON - Un système complet de phonétisation de textes, *Traitement Automatique des Langues (T.A.L.)* 42(1), 47-67.
- BOSCOLO, P. (2008). Writing in primary school. *Handbook of research on writing: History, society, school, individual, text*, 293-309.
- DENIS P., SAGOT B. (2012). Coupling an annotated corpus and a lexicon for state-of-the-art POS tagging. *Language Resources and Evaluation* 4(46), 721-736.
- ELALOUF M.-L. (2011). Constitution de corpus scolaires et universitaires : vers un changement d'échelle ? *Pratiques*, 149-150.
- ELALOUF M.-L. (dir) (2005). *Écrire entre 10 et 14 ans. Un corpus, des analyses, des repères pour la formation*, SCérén, CRDP de Versailles.
- GRANGER S., VANDEVENTER A., HAMEL M.-J. (2001). Analyse des corpus d'apprenants pour l'ELAO basée sur le TAL. *Traitement automatique des langues* 42(2), 609-621.
- KRAIF O., PONTON C. (2007). Du bruit, du silence et des ambiguïtés: que faire du TAL pour l'apprentissage des langues. Actes de *TALN* (Toulouse).
- KUKICH K. (1992). *Techniques for Automatically Correcting Words in Text*. *ACM Computing Surveys* 24 (4), 377–439.
- ORTÉGA É., LÉTÉ B. (2010). « eManulex: Electronic version of Manulex and Manulex-infra databases », <http://www.manulex.org>
- Wolfarth, C., Ponton, C., Totereau, C. (2016). Apports du TAL à la constitution et à l'exploitation d'un corpus scolaire. *Corpus*.