

Aligner production et normalisation : une première approche pour l'étude d'écrits scolaires

Claire Wolfarth

Lidilem

Université Grenoble Alpes

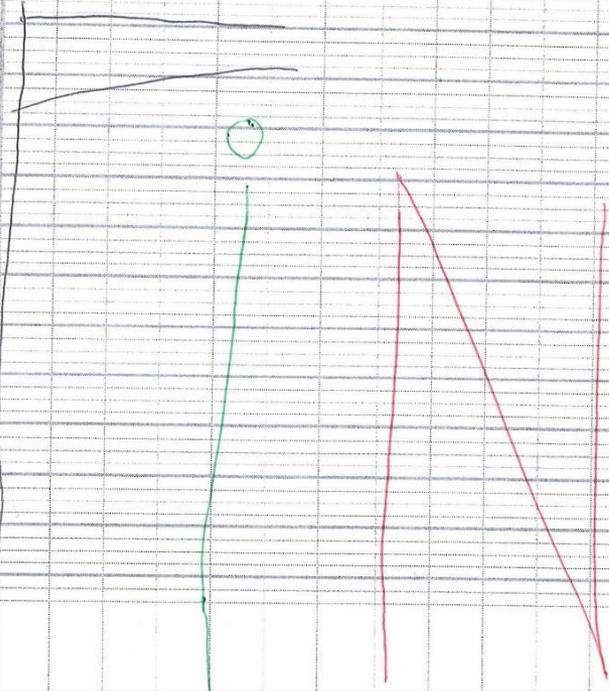


Quelques exemples



Il était une fois un loup qui avait un petit chape-
ron rouge ^{qui} sortis de sa maison le petit chape-
ron rouge aller voir sa grand-mère qui était malade et
lui apporter du lait et un gâteau le petit chape-
ron rouge rencontre un loup le loup dit quel est ce que
tu fais la le petit chape-ron rouge dit j'apporte du
lait et un gâteau a ma ~~maman~~ grand-mère qui
est malade le loup dit je passe par ce chemin
et toi l'autre le petit chape-ron rouge dit d'accord
le petit chape-ron rouge ses fait avoir le loup et passe
par le chemin le plus court et le petit chape-ron
rouge a pris le plus long le loup avança il dit
toc toc toc la grand-mère dit quel est la le
loup dit ses moi le petit chape-ron rouge qui
t'apporte du lait et un gâteau la grand-mère dit
gentre donc mon enfant le loup entra et mangent
la grand-mère il se déguisent en grand-mère
le petit chape-ron rouge arriva et dit toc toc toc
le loup dit qui est la le petit chape-ron rouge dit
ses moi le petit chape-ron rouge qui t'apporte du

lait et un gâteau le loup dit entre mon enfant le
petit chape-ron rouge dit bien mangen sa va
le quérir le petit chape-ron rouge dit que vous
avez de grande dent que vous avez de grandes
oreilles que vous avez de grandes bouche et appe
le petit chape-ron rouge sais fait mangaien par le
loup.



Ma copine est la mère on aisé un chat
dans la forêt abandonné dans la forêt et
de loup côté un soir d'halloween.



Un chat march sur la route.

Le chat tombe sur le route.

Le chat pleure.

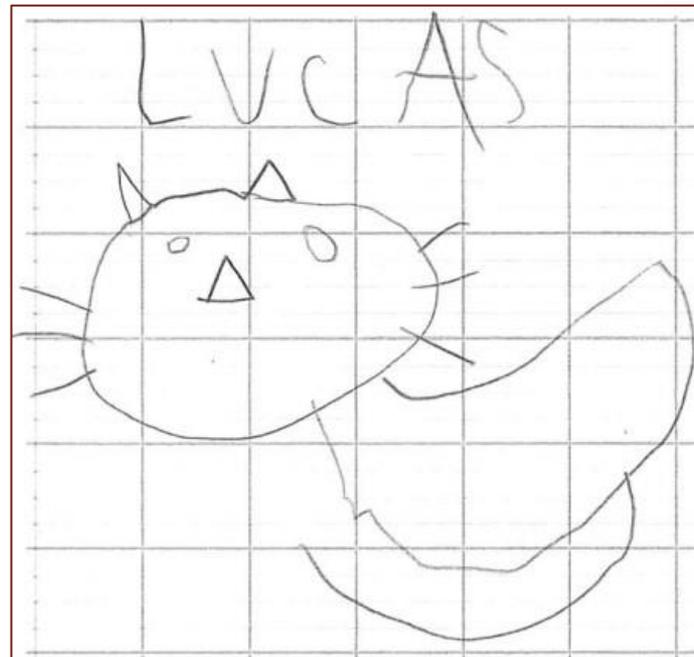
La maman vient à le plan par la
bouche.

l	e	b	a	e	r	t	p	a	t	i		
g	l	i	c	e	→					«glise»		
p	e	l	e	r	→					«pleure»		
e	r	t	c	e	r	t	f	e	m	a	l	e
e	r	t	f	e	m	a	l	e	r	i	e	

↳ «s'est fait mal»

↳ «s'est fait mal au pied»

47



LAP
RA
ÉL
TERA

1/ petit patates

4/ cadre reiro

2/ patison

5/ reioaoniharitlen

3/ cadre

6/ manise

les salade reioa reem - les jardin.

les jeune caton ~~pic~~ picore le blé avec la poule reois.

1/ patin

4/ recreation

2/ patron

5/ charitable

3/ capuchon

6/ manifique

En été, les salade reioa pousse dans les jardins.
Les jeune caton ~~pic~~ picore le blé avec la poule reois.

IAH
KAHU
MNAA
HHUUA

Corpus scolaires : état des lieux

- Développement des recherches sur l'écriture
- Emergence de projets de constitution de grands corpus
 - Elalouf et al., 2005 ;
 - Gunnarsson-Largy et Auriac-Slusarczyk, 2013 ;
 - Garcia-Debanc, 2014 ;
 - David et Doquet, 2016;
 - *Scoledit*
- Annotation manuelle des corpus scolaires



I – Corpus *Scoledit*

II – Principe d'alignement

III – Résultats et perspectives

I – Corpus scolaire : *Scoledit*



Scoledit

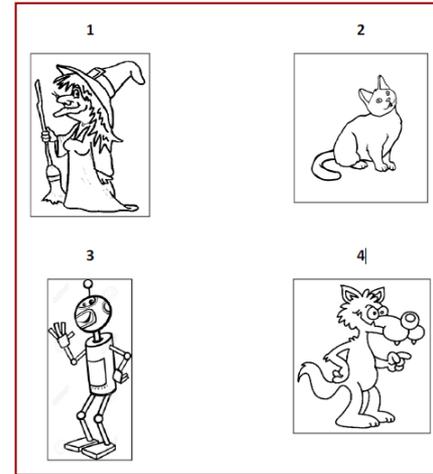
I – Corpus Scoledit

II – Principe d’alignement

III – Résultats et perspectives

- Corpus longitudinal de textes scolaires et de dictées
 - 5 ans (CP – CM2)
 - 5 académiques, 40 écoles
 - 800 – 1200 élèves par an
 - Texte + dictée

<http://otus.u-grenoble3.fr/scoledit/>



Enjeux

- Pour la linguistique :
 - Description en synchronie et en diachronie des caractéristiques des écrits d’apprenants de 6 à 11 ans
 - Meilleure connaissance des phénomènes d’acquisition de l’écriture (évolution de l’acquisition de la syntaxe, de l’orthographe, de la ponctuation, ...)
- Pour la didactique:
 - Développement de séquences et de dispositifs à partir de ces connaissances
- Pour le traitement automatique des langues (TAL) :
 - Type de textes encore peu étudiés



Méthodologie

Constats :

- Erreurs nombreuses
- Erreurs peu régulières

Hypothèse :

- Comparer la production des apprenants avec une norme de cette production devrait permettre de faciliter son annotation
 - Développement d'un aligneur tokens par tokens des productions transcrites et des productions normées

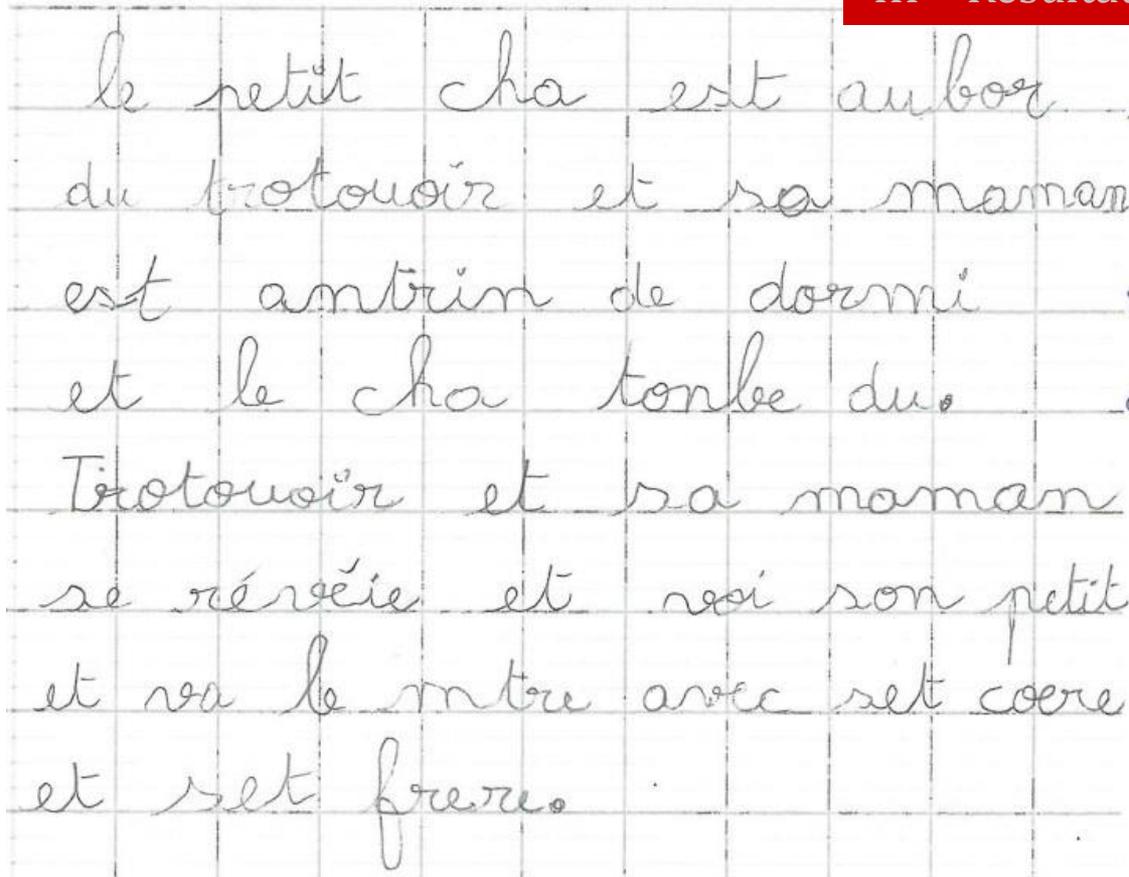


Spécificités du corpus

I – Corpus Scoledit

II – Principe d'alignement

III – Résultats et perspectives



le petit cha est au bord
du trottoir et sa maman
est en train de dormir
et le cha tombe du
Trottoir et sa maman
se réveille et voit son petit
et va le mettre avec ses sœurs
et ses frères.

Production de l'élève 96 en fin de CP

Version orthographiée selon la norme : « Le petit chat est au bord du trottoir et sa maman est en train de dormir et le chat tombe du trottoir et sa maman se réveille et voit son petit et va le mettre avec ses sœurs et ses frères. »

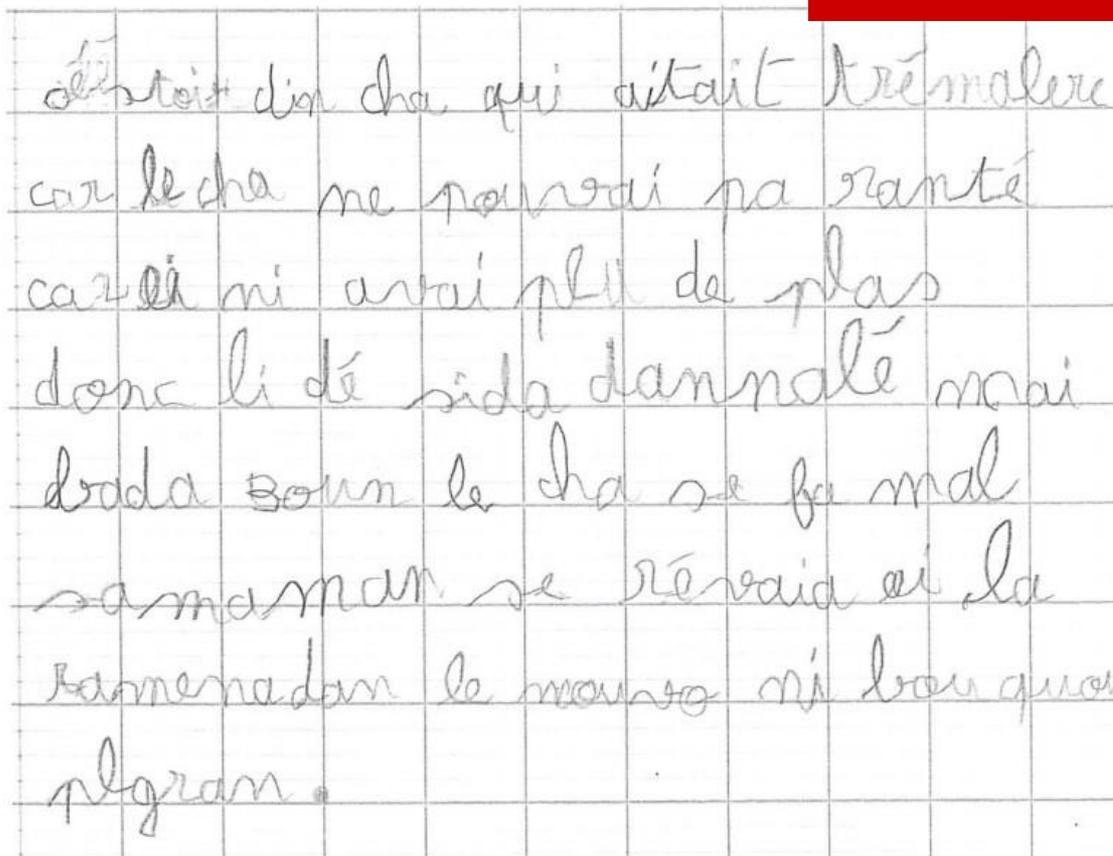


Spécificités du corpus

I – Corpus Scoledit

II – Principe d'alignement

III – Résultats et perspectives



étoit d'un chat qui était très malheure
car le chat ne pouvait pas rentrer
car il n'y avait plus de place
donc il décida d'aller mais
badaboum le chat se fait mal
sa maman se réveilla et la
ramena dans le nouveau nid beaucoup
plus grand.

Production de l'élève 2972 en fin de CP

Version orthographiée selon la norme : « C'est l'histoire d'un chat qui était très malheureux car le chat ne pouvait pas rentrer car il n'y avait plus de place donc il décida de s'en aller mais badaboum le chat se fait mal <segmentation/> sa maman se réveilla et le ramena dans le nouveau nid beaucoup plus grand. »

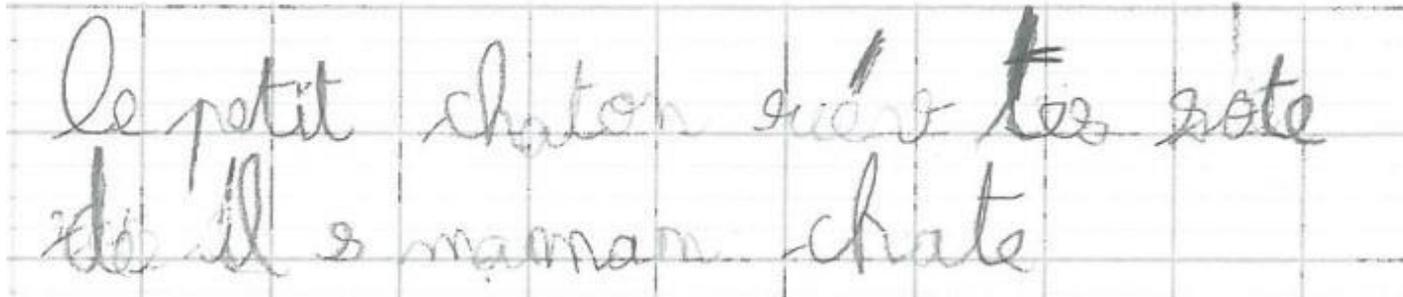


Spécificités du corpus

I – Corpus *Scoledit*

II – Principe d'alignement

III – Résultats et perspectives



Production de l'élève 2453 en fin de CP

Version orthographiée selon la norme : « *le petit chaton se réveille et saute de [...] maman chatte.* »



Composition du corpus

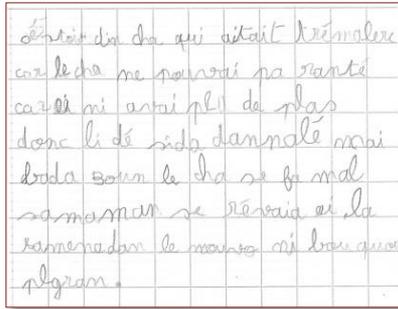
I – Corpus *Scoledit*

II – Principe d'alignement

III – Résultats et perspectives

Pour chaque production :

- Scan



- Transcription

cé <revision/> stoir din cha qui aïtât tré malere / car le cha ne pouvai pa ranté / car li ni avai plu de plas / donc le dé sida d'annalé mai / bada Boum le cha se fa mal / sa maman se révaia <letMF>a</letMF>i la / ramena dan le nouvo ni bou quou / plgran.

- Normalisation

c'est l'histoire d'un chat qui était très malheureux car le chat ne pouvait pas rentrer car il n'y avait plus de place donc il décida de s'en aller mais badaboum le chat se fait mal sa maman se réveilla et le ramena dans le nouveau nid beaucoup plus grand.



II – Principe d'alignement



Alignement : enjeux

I – Corpus scolaires : spécificités

II – Principe d'alignement

III – Résultats et perspectives

Requête : *Comment les élèves en fin de CP réalisent les verbes à l'imparfait ?*

Transcription :

[...] qui aitait très malheure car le cha ne pouvai pas ranté car li ni avai plu [...]

ALIGNEMENT

Normalisation :

[...] qui était très malheureux car le chat ne pouvait pas rentrer car il n' y avait plus [...]

Étiquettes

morphosyntaxiques :

...	PRO	VER	ADV	ADJ	KON	DET	NOM	ADV	VER	ADV	VER	KON	PRO	ADV	PRO	VER	ADV	...
REL		<u>impf</u>				ART			<u>impf</u>		<u>inf</u>		PER		PER	<u>impf</u>		

Résultat :

- *était*, orthographié « aitait »
- *pouvait*, orthographié « pouvai »
- *avait*, orthographié « avai »
- ...



Etat de l'art

I – Corpus scolaires : spécificités

II – Principe d'alignement

III – Résultats et perspectives

- Dictées
 - Santiago-Oriola (1998)
 - Beaufort et Roekhaut (2011)
- Traduction
 - Desmet, Hérogue, 2005



Principes d'alignement

I – Corpus scolaires : spécificités

II – Principe d'alignement

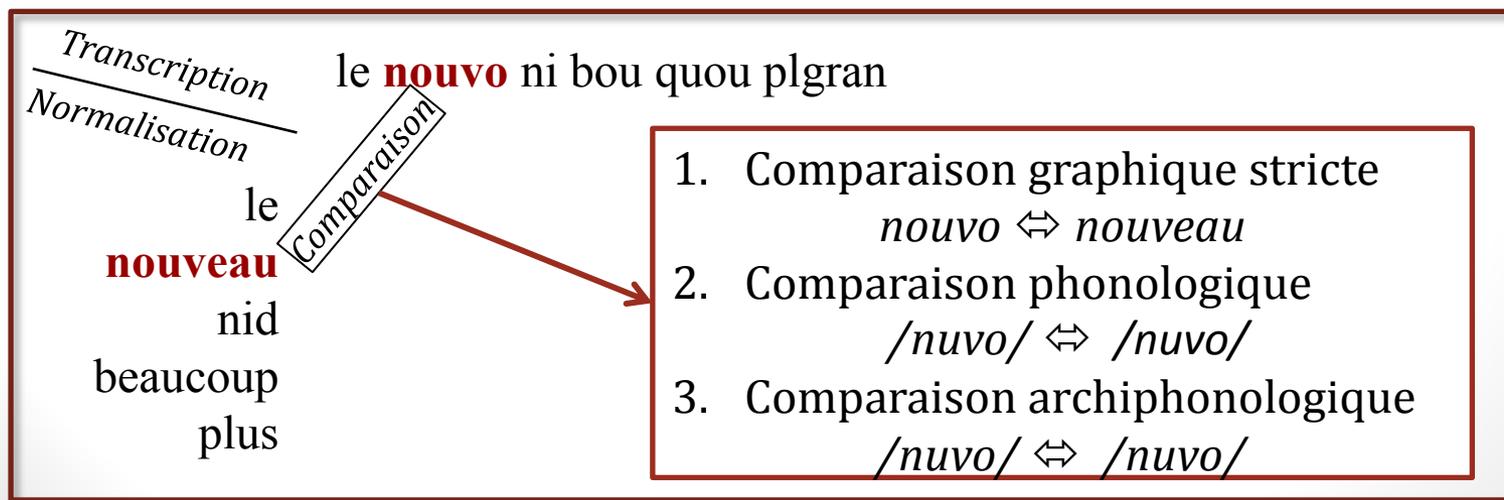
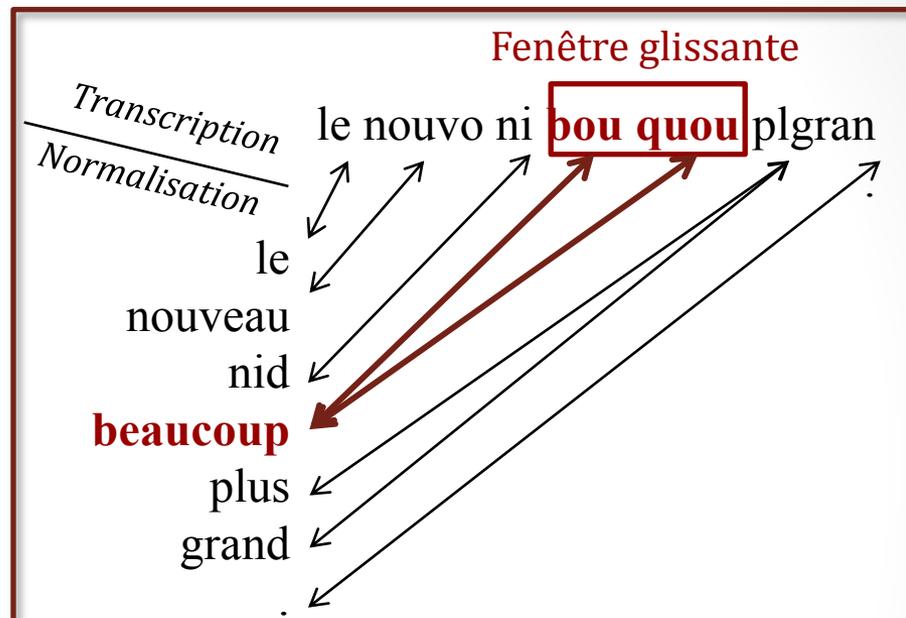
III – Résultats et perspectives

Unité de comparaison

- Alignement de tokens
- Fenêtre glissante

Modes de comparaison

- Indices phonologiques
- Indices phonologiques avec relâchement de contraintes



Exemple

I – Corpus scolaires : spécificités

II – Principe d'alignement

III – Résultats et perspectives

Transcription

Normalisation

é
il
miole
sa
maman
li
di
plere
pa
petit
chaton
et
elle
latrape
.



et
il
miaule
sa
maman
lui
dit
ne
pleure
pas
petit
chaton
et
elle
l'
attrape
.



Exemple

I – Corpus scolaires : spécificités

II – Principe d'alignement

III – Résultats et perspectives

Transcription

Normalisation

é
il
miole
sa
maman
li
di
plere
pa
petit
chaton
et
elle
latrape
.



et
il
miaule
sa
maman
lui
dit
ne
pleure
pas
petit
chaton
et
elle
l'
attrape
.

① Comparaison graphique stricte



Exemple

I – Corpus scolaires : spécificités

II – Principe d'alignement

III – Résultats et perspectives

Transcription

é
il
miole
sa
maman
li
di
plere
pa
petit
chaton
et
elle
latrape
.

Normalisation

et
il
miaule
sa
maman
lui
dit
ne
pleure
pas
petit
chaton
et
elle
l'
attrape
.

2

- 1 Comparaison graphique stricte
- 2 Comparaison phonologique

é et
/e/ ↔ /e/

Exemple

I – Corpus scolaires : spécificités

II – Principe d'alignement

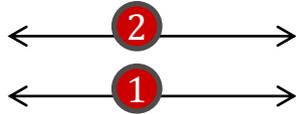
III – Résultats et perspectives

Transcription

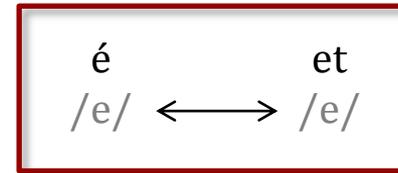
é
il
miole
sa
maman
li
di
plere
pa
petit
chaton
et
elle
latrape
.

Normalisation

et
il
miaule
sa
maman
lui
dit
ne
pleure
pas
petit
chaton
et
elle
l'
attrape
.



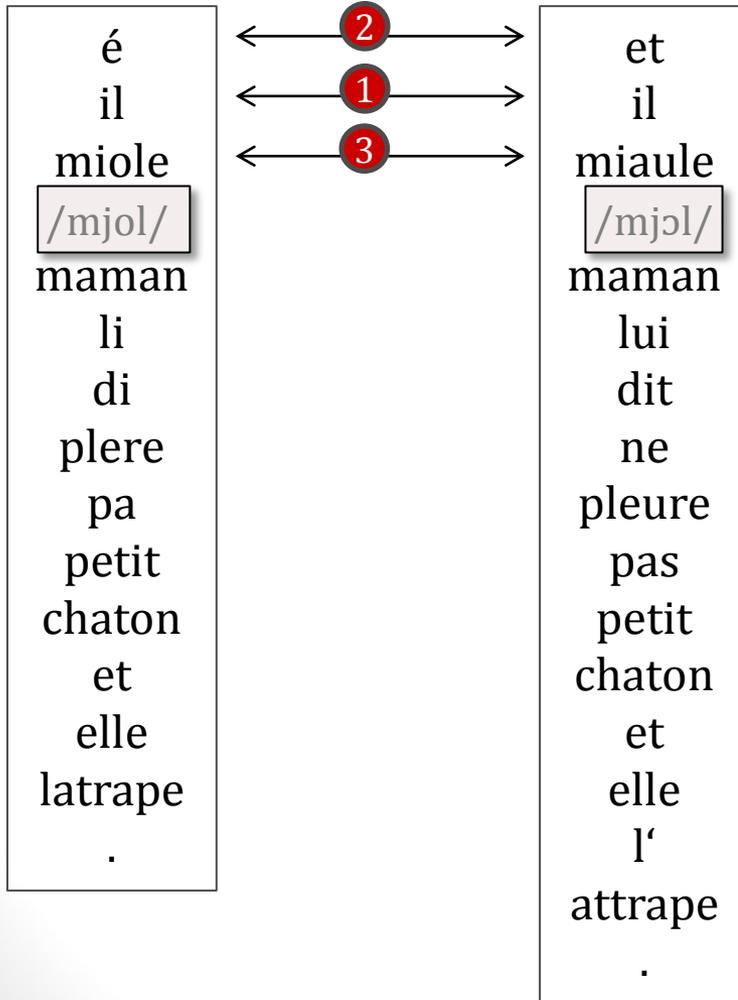
- 1 Comparaison graphique stricte
- 2 Comparaison phonologique



Exemple

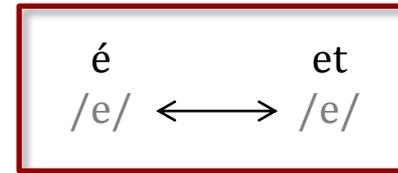
Transcription

Normalisation



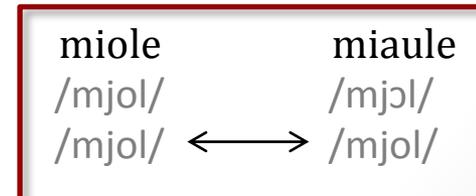
1 Comparaison graphique stricte

2 Comparaison phonologique



3 Comparaison phonologique avec relâchement de contraintes

e	ə	œ	ɔ	∅ final
ɛ	∅	ẽ	o	schwa



Exemple

Transcription

Normalisation

é	← 2 →	et
il	← 1 →	il
miole	← 3 →	miaule
sa	← 1 →	sa
maman	← 1 →	maman
li		lui
di		dit
plere		ne
pa		pleure
petit		pas
chaton		petit
et		chaton
elle		et
latrape		elle
.		l'
		attrape
		.

① Comparaison graphique stricte

② Comparaison phonologique

é	↔	et
/e/		/e/

③ Comparaison phonologique avec relâchement de contraintes

e	ə	ẽ	ɔ	∅ final
ɛ	∅	ẽ	o	schwa

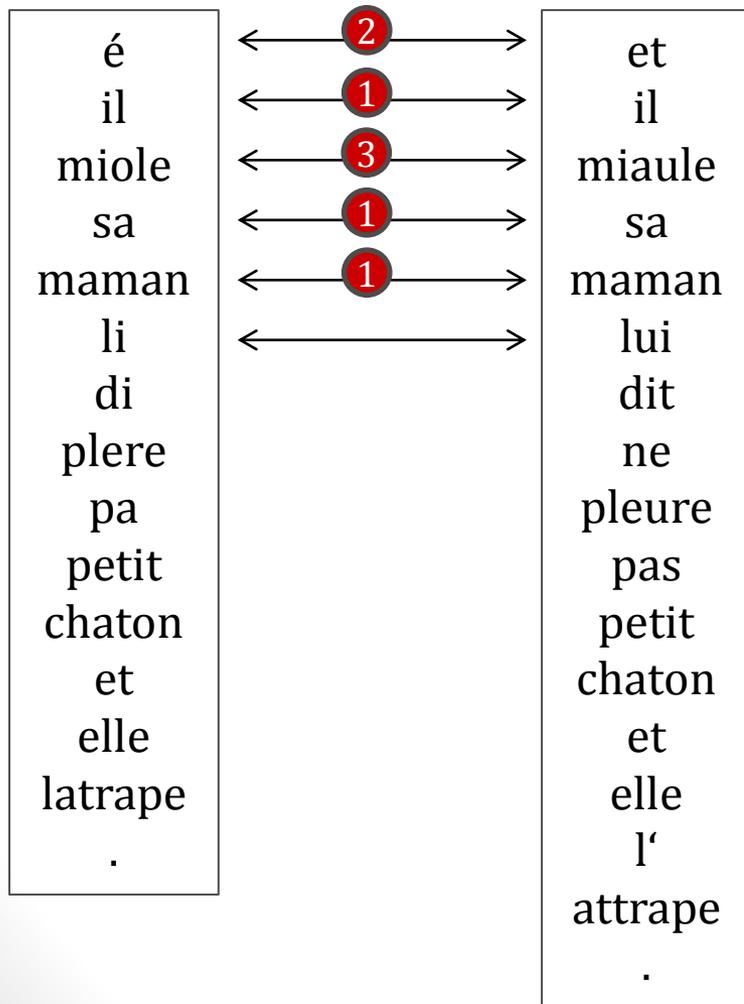
miole	↔	miaule
/mjol/		/mjol/
/mjol/		/mjol/



Exemple

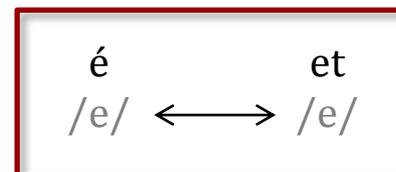
Transcription

Normalisation



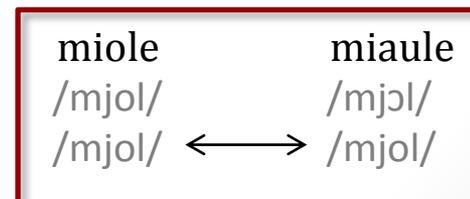
① Comparaison graphique stricte

② Comparaison phonologique



③ Comparaison phonologique avec relâchement de contraintes

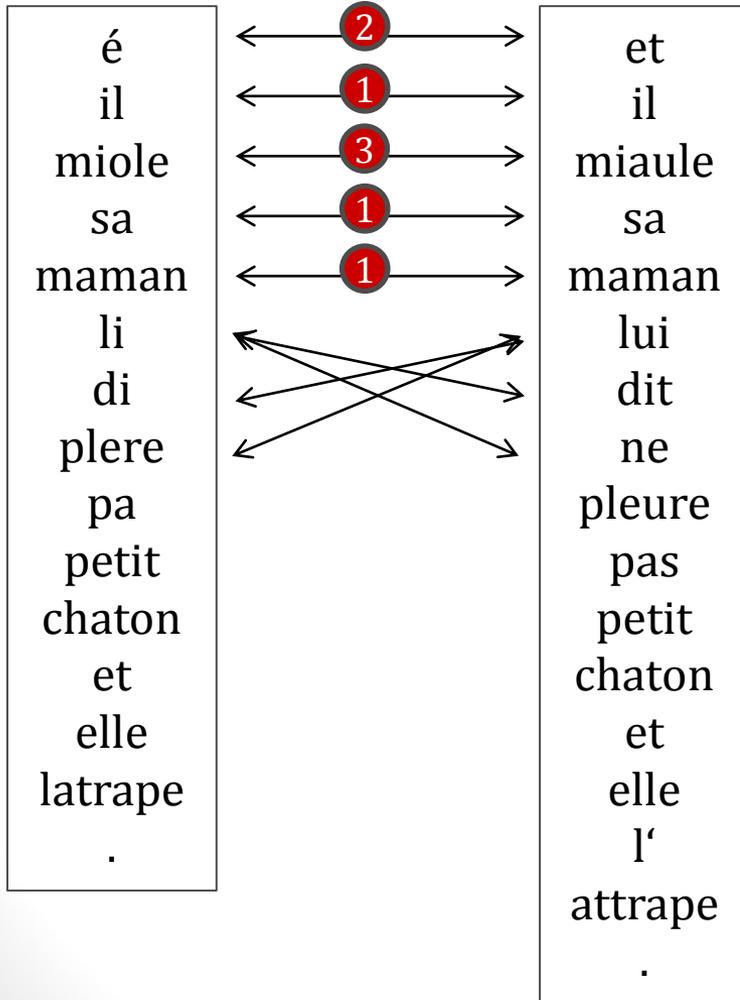
e	ə	œ	ɔ	∅ final
ɛ	∅	ẽ	o	schwa



Exemple

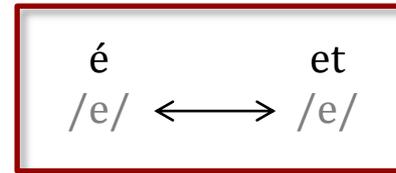
Transcription

Normalisation



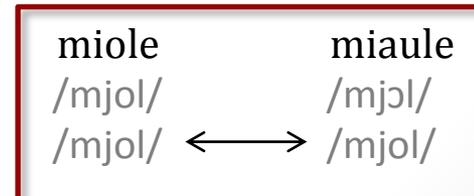
1 Comparaison graphique stricte

2 Comparaison phonologique



3 Comparaison phonologique avec relâchement de contraintes

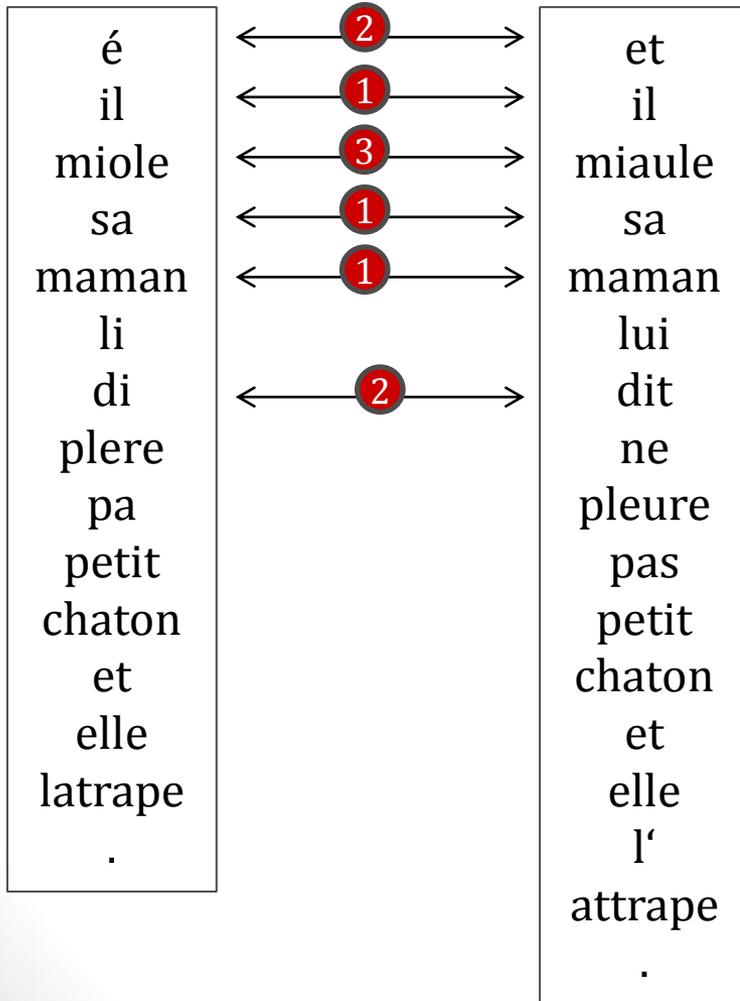
e	ə	œ	ɔ	∅ final
ɛ	∅	ẽ	o	schwa



Exemple

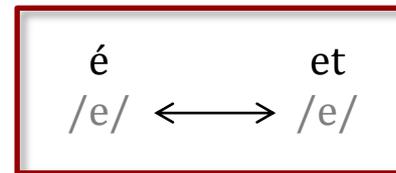
Transcription

Normalisation



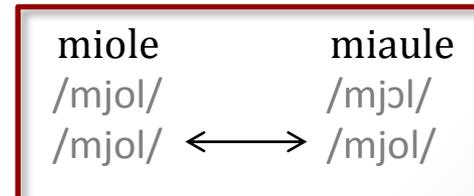
① Comparaison graphique stricte

② Comparaison phonologique



③ Comparaison phonologique avec relâchement de contraintes

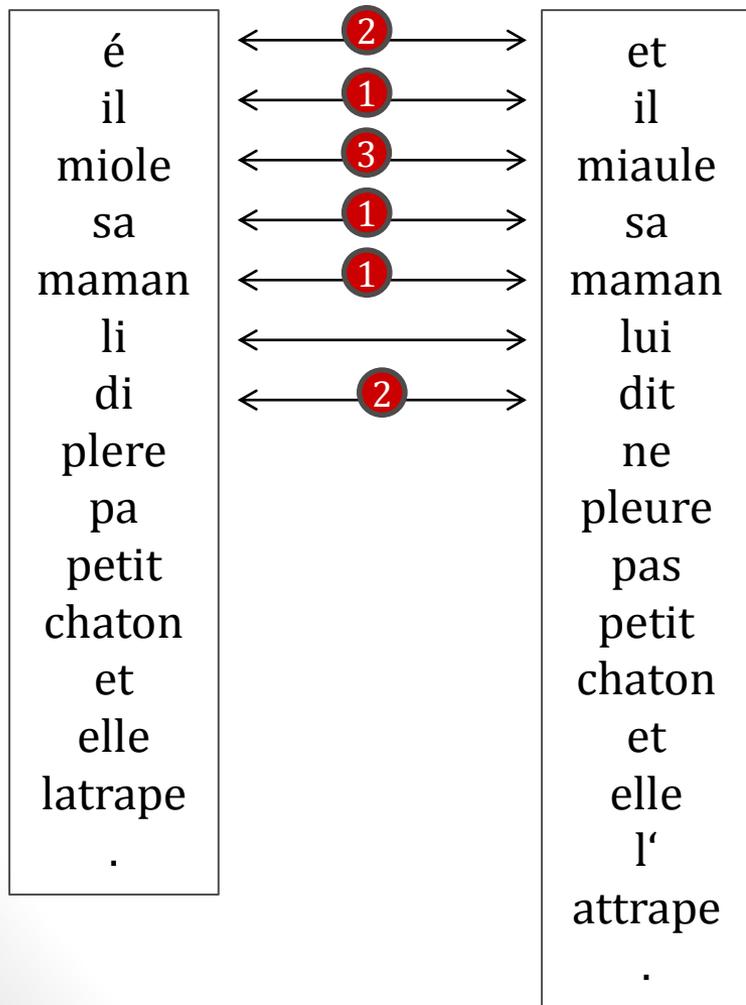
e	ə	ẽ	ɔ	∅ final
ɛ	∅	ẽ	o	schwa



Exemple

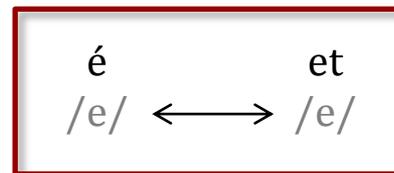
Transcription

Normalisation



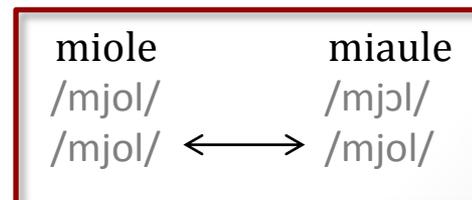
1 Comparaison graphique stricte

2 Comparaison phonologique



3 Comparaison phonologique avec relâchement de contraintes

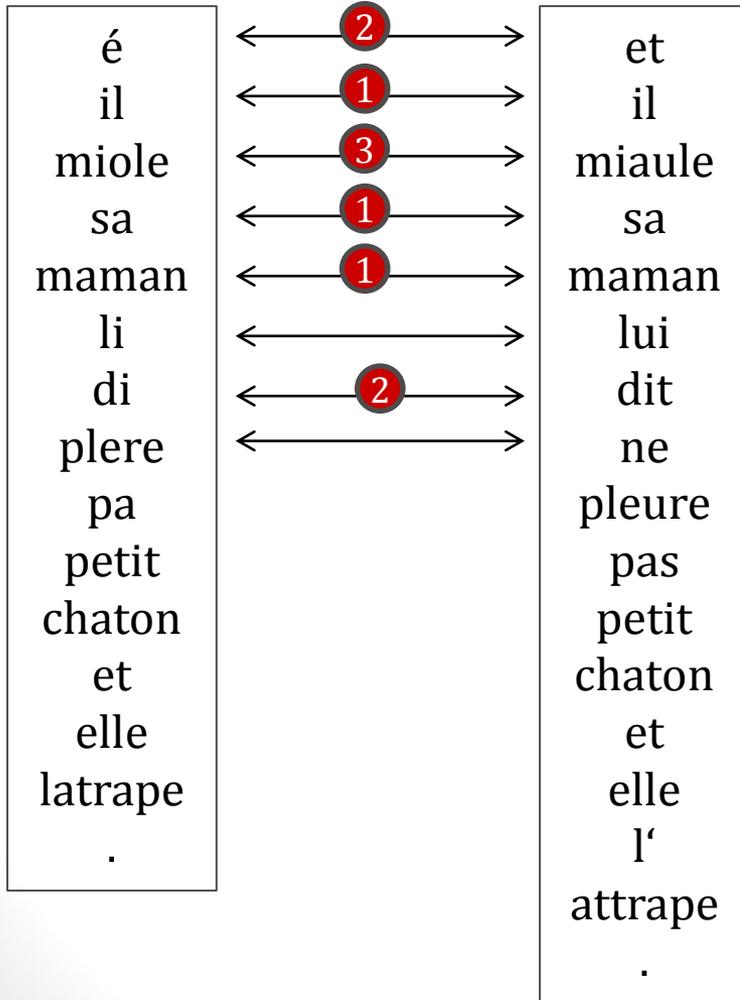
e	ə	œ	ɔ	∅ final
ɛ	∅	ẽ	o	schwa



Exemple

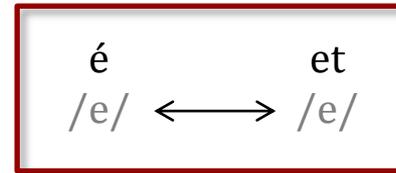
Transcription

Normalisation



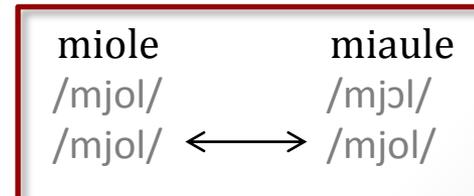
① Comparaison graphique stricte

② Comparaison phonologique



③ Comparaison phonologique avec relâchement de contraintes

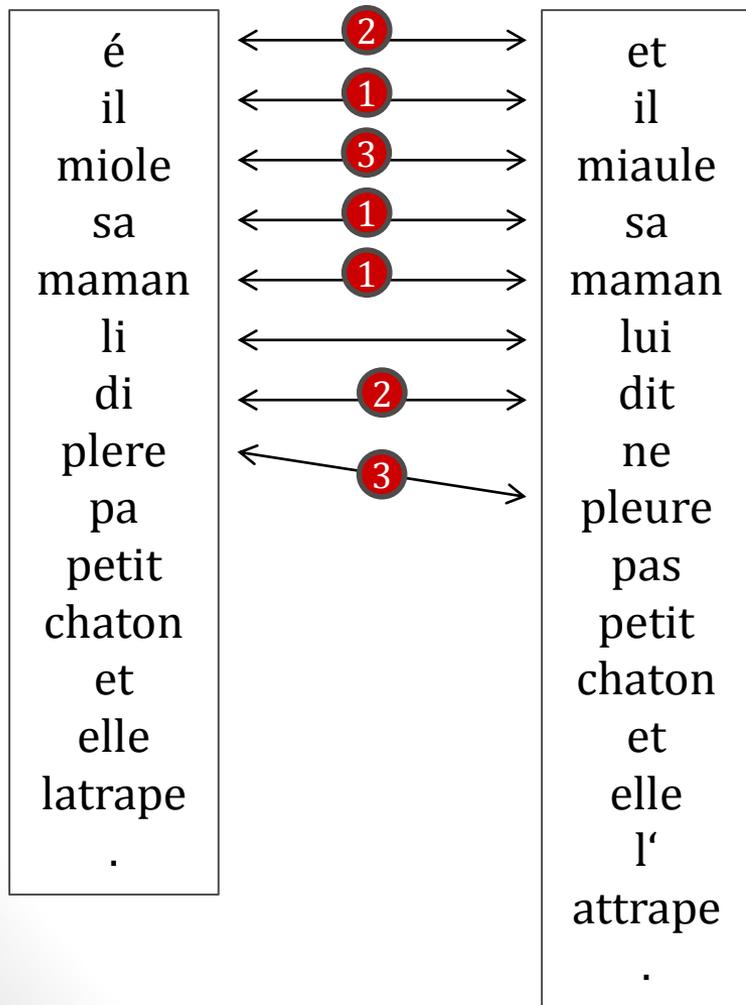
e	ə	œ	ɔ	∅ final
ɛ	∅	ẽ	o	schwa



Exemple

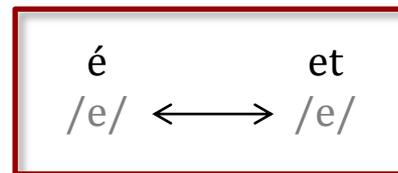
Transcription

Normalisation



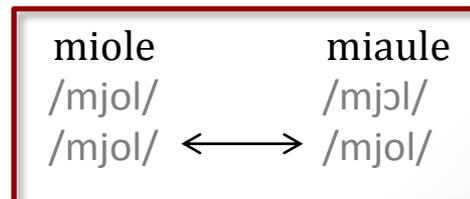
① Comparaison graphique stricte

② Comparaison phonologique



③ Comparaison phonologique avec relâchement de contraintes

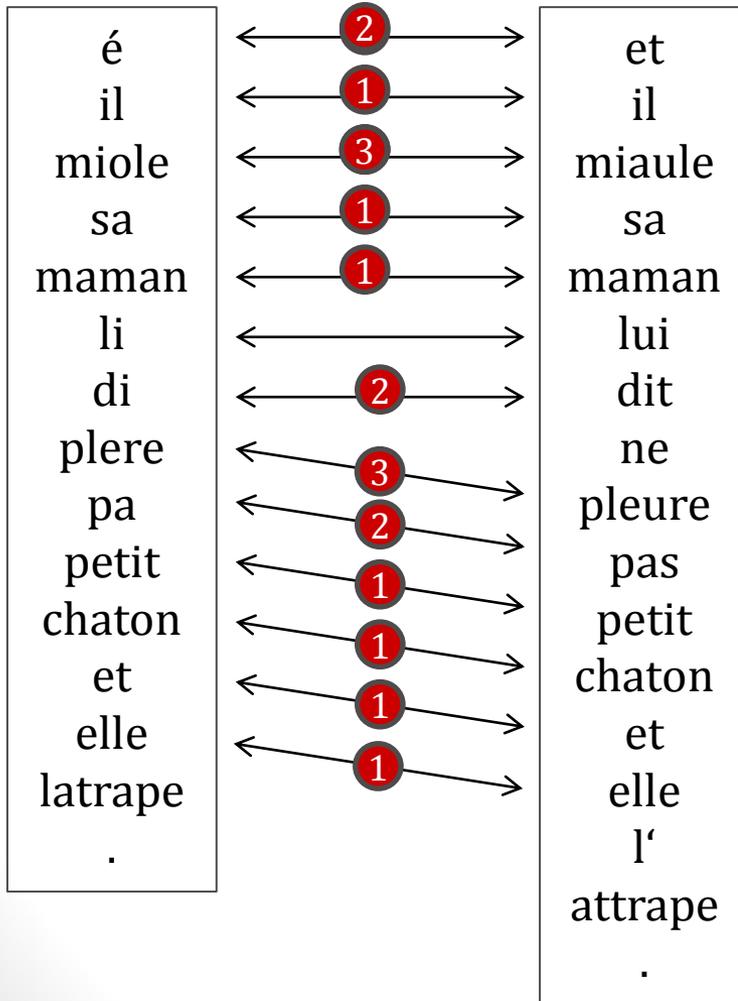
e	ə	œ	ɔ	∅ final
ɛ	∅	ẽ	o	schwa



Exemple

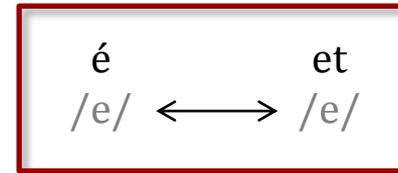
Transcription

Normalisation



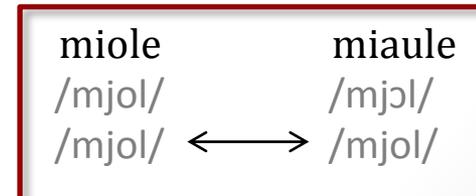
1 Comparaison graphique stricte

2 Comparaison phonologique



3 Comparaison phonologique avec relâchement de contraintes

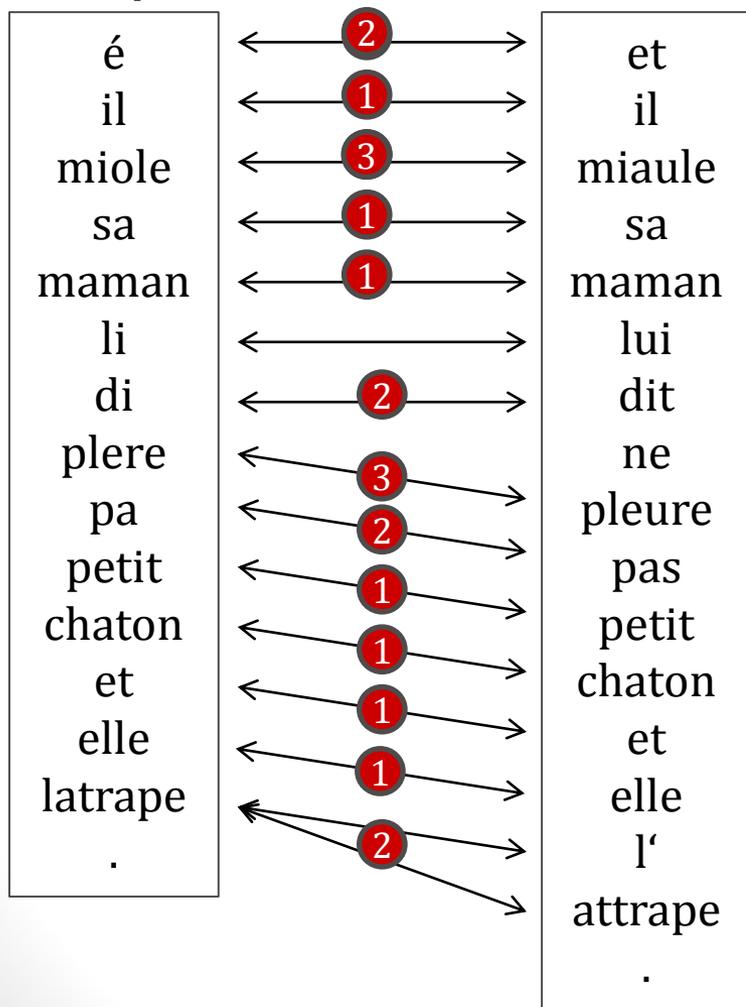
e	ə	œ	ɔ	∅ final
ɛ	∅	ẽ	o	schwa



Exemple

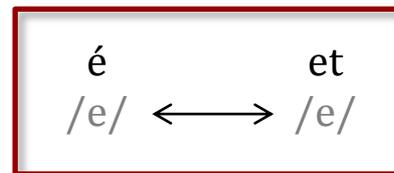
Transcription

Normalisation



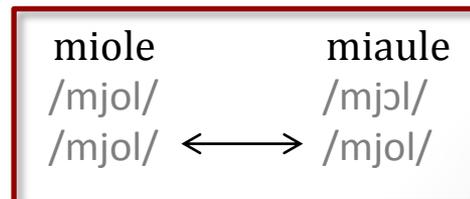
① Comparaison graphique stricte

② Comparaison phonologique



③ Comparaison phonologique avec relâchement de contraintes

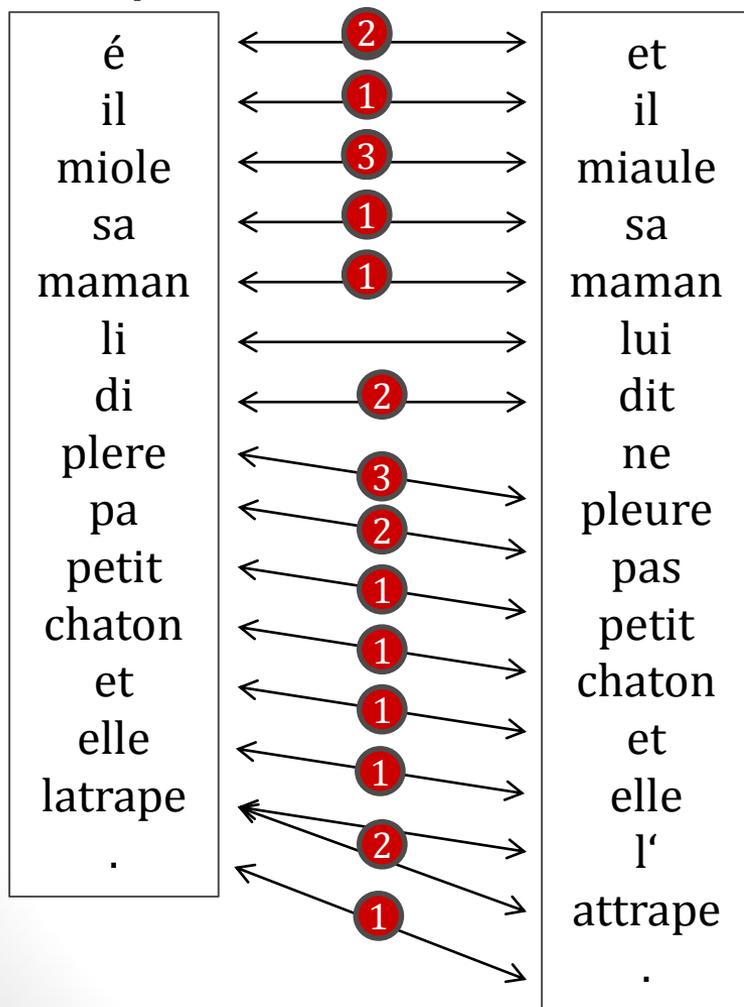
e	ə	œ	ɔ	∅ final
ɛ	∅	ẽ	o	schwa



Exemple

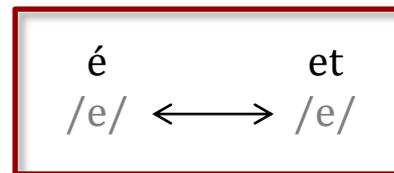
Transcription

Normalisation



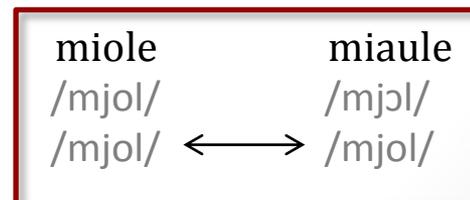
① Comparaison graphique stricte

② Comparaison phonologique



③ Comparaison phonologique avec relâchement de contraintes

e	ə	œ	ɔ	∅ final
ɛ	∅	ẽ	o	schwa



III – Résultats et perspectives



Résultats

I – Corpus scolaires : spécificités
II – Principe d'alignement
III – Résultats et perspectives

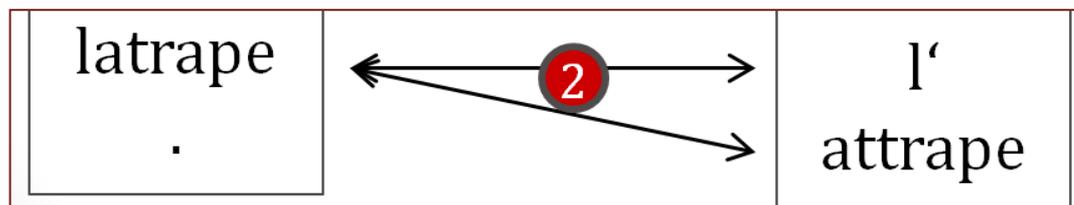
Mesure	Alignement par comparaison graphique stricte	Alignement par comparaison phonologique	Alignement par comparaison phonologique avec relâchement de contraintes
Précision	83,3 %	86,3 %	86,7 %
Rappel	59,9 %	72,6 %	74 %
F-mesure	69,7 %	78,9 %	80 %



Résultats

I – Corpus scolaires : spécificités
II – Principe d'alignement
III – Résultats et perspectives

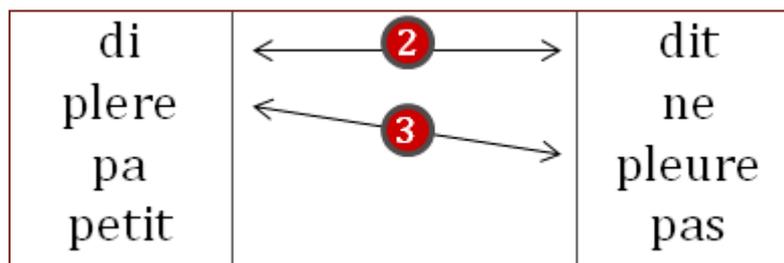
Taille de la fenêtre	Aucune concaténation	Concaténation n.(n+1)	Concaténation n.(n+1).(n+2)	Concaténation n.(n+1).(n+2).(n+3)
Précision	87,82%	88,22%	88,17%	88,17%
Rappel	71,74%	75,48%	75,60%	75,60%
F-mesure	78,97%	81,35%	81,41%	81,41%



Résultats

I – Corpus scolaires : spécificités
II – Principe d'alignement
III – Résultats et perspectives

Distance	Comparaison à n+1	Comparaison à n+2	Comparaison à n+3
Précision	86,89%	88,17%	86,93%
Rappel	66,35%	75,60%	74,21%
F-mesure	75,24%	81,41%	80,07%



Formes non reconnues

I – Corpus scolaires : spécificités

II – Principe d'alignement

III – Résultats et perspectives

de	de
dormire	dormir
a#lor	alors
le	le
petit	petit
des#den	descend
dans	dans
lescalie	l'#escalier
boum	boum
.	.
il#miel	Il#miaule#<dialogue>
miaou	miaou
	</dialogue>
la	la
maman	maman
chat	chat
se	se
revei#re#mei	réveille#remet
son	son
petit	petit
é	et
il	il
sen#dors	s'#endort
.	.



Formes non reconnues

I – Corpus scolaires : spécificités

II – Principe d'alignement

III – Résultats et perspectives

il est une fois un chat ci
marché à un autre chat ci
arrivé de gonfler le chat ton
le le chat se le chat ramasse
les chat ils vivent en commun
tréméalisme pour toujours.

le	le
cha	chat
sé#la#cha#rama	<incomprehensible/>
les	les
cha	chats
éilvunmcomon#tréméalism#p	et#ils#vivent#en#commun#tr
ourtoujour#.	ès#<incomprehensible/>
	pour
	toujours#.#



Visualisation Corpus CP

Afficher la CE1

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 >

Rechercher un élève :

N° d'élève

ok

Il y a actuellement à l'étude : 1400 élèves.

Élève 47

CP

CE1

SCAN DICTÉE SEPTEMBRE

LP lapin
RAPH rat
LPIH éléphant
LLULULIH
"ton jeu avec le rat"

DICTÉE SEPTEMBRE

LP
RAPH
LPIH
LLULULIH

DICTÉE JUIN

lapin
rat
à l'éfant
ton jeu avec le rat
les lapin cour vite

SCAN DICTÉE JUIN

lapin
rat
à l'éfant
ton jeu avec le rat
les lapin cour vite

PRODUCTION

le chat et tonba dent
les zoccalié et il
pleur parce que il
sas fo tra tra male
et sa maman le
protg et il le
loch.

SCAN PRODUCTION

le chat et tonba dent
les zoccalié et il
pleur parce que il
sas fo tra tra male
et sa maman le
protg et il le
loch.

Élève 49

CP

CE1



Lexique 476 mots distincts toutes classes confondues.

Afficher les formes produites par les élèves

Afficher les formes rectifiées (normées)

Formes du dictionnaire (vocabulaire)

Vocabulaire CP

476 mots distincts sur 12686 répertoriés à ce jour (3,75%).

Mots	Nb d'occurrences	Mots	Nb d'occurrences
le	1565	chat	1052
il	892	et	822
petit	483	maman	482
son	454	un	438
être	398	se	393
tomber	393	pleurer	230
faire	174	avoir	172
miaou	156	réveiller	149
de	147	qui	139
marche	120	dormir	119
chaton	118	mal	116

Visualisation des mots dictés (Vue administrateur)

CP : 1169 élèves

éléphant

lapin

rat

Septembre

Juin

Septembre

Juin

Septembre

Juin

789 versions différentes

327 versions différentes

601 versions différentes

114 versions différentes

431 versions différentes

77 versions différentes

Septembre		Juin		Septembre		Juin		Septembre		Juin	
Production	Nombre										
E	23	éléphant	151	A	56	lapin	676	RA	279	ra	338
L	20	éléphant	114	LAP	52	la pin	43	ra	80	rat	258
dessin	19	éléfan	112	LAPIN	29	lapain	17	A	55	ras	144
ELF	13	éléfen	32	LA	28	la pun	12	R	37	rae	59
ÉLÉFAN	12	éléfant	15	dessin	24	lapun	12	Ra	20	rats	19
ÉLF	10	éléfane	14	L	22	Lapin	10	RAE	19	rad	7
EL	9	éléphants	12	LAPUN	22	lapine	9	AR	17	le ra	6
A	7	éléphen	11	LAPA	21	lapni	7	dessin	15	rar	5
ELEFAN	7	éléfon	10	LP	18	la pain	6	RARA	8	Ra	4
ELFA	6	élephan	7	LAPN	16	lapins	6	RAT	8	RA	3
é	6	éléphan	7	lapun	15	lapai	4	rae	8	le ras	3
éléfan	6	éléfent	7	LAPE	11	la	3	RAA	7	raes	3
LE	5	éléphane	7	la	11	la pine	3	a	7	un ra	3
T	5	elefan	6	LPA	9	la pins	3	RAL	6	un rat	3



Bibliographie

- BEAUFORT R., ROEKHAUT S. (2011). AUTOMATION OF DICTATION EXERCISES. *COMPUTATIONAL LINGUISTICS IN THE NETHERLANDS JOURNAL*, (1), 1-20.
- BÉCHET F. (2001). LIA_PHON - Un système complet de phonétisation de textes, *Traitement Automatique des Langues (T.A.L.)* 42(1), 47-67.
- BORÉ C., ELALOUF M.-L. (2017). Deux étapes dans la construction de corpus scolaires : problèmes récurrents et perspectives nouvelles. In Doquet C., David J., Fleury S., *Spécificités et contraintes des grands corpus de textes scolaires : problèmes de transcription, d'annotation et de traitement*, Corpus, 16 | déc. 2016, 185-214.
- CATACH N. (1980, 3e édition, 1995). L'orthographe française : traité théorique et pratique avec des travaux d'application et leurs corrigés (avec la collaboration de Gruz C. et Duprez D.). Paris : Nathan.
- DAVID J., DOQUET C. (2016). LES ÉCRITS D'ÉLÈVES: UN CORPUS DE RÉFÉRENCE POUR LE FRANÇAIS CONTEMPORAIN. Actes de *SHS WEB OF CONFERENCES* (VOL. 27, p. 11001). EDP SCIENCES.
- DESMET P., HÉROGUEL A. (2005). LES ENJEUX DE LA CRÉATION D'UN ENVIRONNEMENT D'APPRENTISSAGE ÉLECTRONIQUE AXÉ SUR LA COMPRÉHENSION ORALE À L'AIDE DU SYSTÈME AUTEUR IDIOMA-TIC. Actes d *ALSIC. APPRENTISSAGE DES LANGUES ET SYSTÈMES D'INFORMATION ET DE COMMUNICATION*, 8(1), 281-303.
- ELALOUF M.-L. (dir) (2005). *Écrire entre 10 et 14 ans. Un corpus, des analyses, des repères pour la formation*, SCérén, CRDP de Versailles.
- GARCIA-DEBANC C., BONNEMAISON K. (2014). LA GESTION DE LA COHÉSION TEXTUELLE PAR DES ÉLÈVES DE 11-12 ANS: RÉUSSITES ET DIFFICULTÉS. Actes de *SHS WEB OF CONFERENCES* (VOL. 8, PP. 961-976). EDP SCIENCES.
- SANTIAGO ORIOLA C. (1998). Système vocal interactif pour l'apprentissage des langues. La synthèse de la parole au service de la dictée. Thèse. Toulouse III.
- WOLFARTH C., PONTON C., TOTEREAU C. (2017). Apports du TAL à la constitution et à l'exploitation d'un corpus scolaire au travers du développement d'un outil d'annotation orthographique. *Corpus* 16, janv. 2017, 185-214.



Merci de votre attention



Aligner production et normalisation : une première approche pour l'étude d'écrits scolaires

Claire Wolfarth

Lidilem

Université Grenoble Alpes