

# Which Method to Develop a Natural Language Processing Tool to automatically analyse First Language Learner Corpora?

Claire Wolfarth, Claude Ponton, Catherine Brissaud

Lidilem

Université Grenoble Alpes



IAH  
XAHU  
MNA  
HHUA

Dictations – begin of 1st year

LAP  
RA  
ÉL  
TERA

1. lapin
2. rat
3. éléphant
4. comme jouer avec le rat
5. les lapin cour vite

Dictations – end of 1st year

Il était une fois un loup qui avait un petit chapeau rouge. Il sortait de sa maison le petit chapeau rouge aller voir sa grand-mère qui était malade et lui apporter du lait et un gâteau. Le petit chapeau rouge rencontre un loup. Le loup dit: "C'est ce que tu fais?" Le petit chapeau rouge dit: "J'apporte du lait et un gâteau à ma maman, grand-mère qui est malade." Le loup dit: "Je passe par ce chemin et toi l'autre." Le petit chapeau rouge dit: "D'accord." Le petit chapeau rouge se fait avoir le loup et passe par le chemin le plus court. Le petit chapeau rouge a pris le plus long. Le loup arriva. Il dit: "Toc toc toc." La grand-mère dit: "Qui est là?" Le loup dit: "C'est moi, le petit chapeau rouge qui t'apporte du lait et un gâteau." La grand-mère dit: "Entre donc, mon enfant." Le loup entra et mangent la grand-mère. Il se déguisa en grand-mère. Le petit chapeau rouge arriva et dit: "Toc toc toc." Le loup dit: "Qui est là?" Le petit chapeau rouge dit: "C'est moi, le petit chapeau rouge qui t'apporte du

Narrative stories – end of 1st year

Un chat marcher sur la route.  
Le chat tombe sur le route.  
Le chat pleure.  
La maman vient à le pleurer par la bouche.

lait et un gâteau. Le loup dit: "Entre mon enfant." Le petit chapeau rouge dit: "Bien manger de la viande." Le petit chapeau rouge dit: "C'est ce que vous avez de grande dent que vous avez de grandes oreilles que vous avez de grandes bouche et que le petit chapeau rouge s'est fait manger par le loup."

Narrative stories – end of 3rd year

Narrative stories – end of 1st year

le garçon la femme le animal le geste

# Inventory of school corpora

- 2<sup>nd</sup> language learner corpora
- Emergence of large school corpus projects
  - Longitudinal
    - Juel, 1988 (English)
    - Lancaster Corpus of Children's Project Writing, 1998 (English)
    - Corpus ÉMA, Boré & Elalouf, 2017 (French)
  - Non longitudinal
    - Elalouf, 2005 (French)
    - Project Ecriscol, David et Doquet, 2016 (French)
- Length, accessibility, processes



# I – *Scoledit*: a French longitudinal school Corpus

*Project funded by French research national Agency (E-CALM project)*



# *Scoedit* corpus

## Issues

- In linguistics
  - Synchronous and diachronic description of the characteristics of learners' writings from 6 to 11 years old
  - Better knowledge of writing acquisition phenomena (evolution of syntax acquisition, spelling, punctuation, etc.)
- In automatic language processing (NLP)
  - Type of texts still little studied
- In didactics
  - Development of sequences and devices based on this knowledge



# Scoledit corpus

- A longitudinal corpus of narrative writings and dictations
  - 5 years (primary school)
  - 40 schools: 800 – 1200 pupils per year
  - narrative writings + dictations

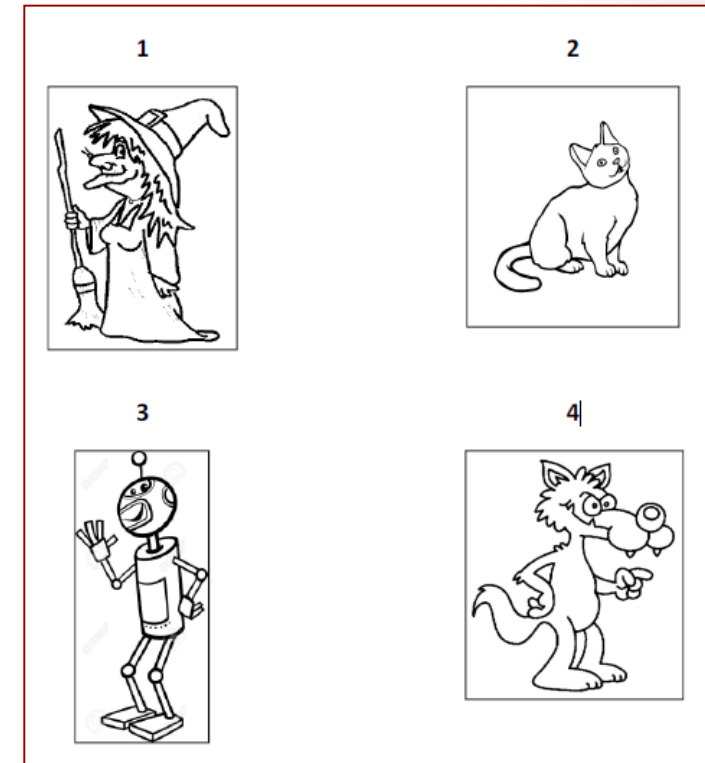
<http://scoledit.org/scoledit/>

Grade	Number of writings	Number of words
1 <sup>st</sup> year (2013-2014)	977	16 000
2 <sup>nd</sup> year (2014-2015)	876	67 000
3 <sup>rd</sup> year (2015-2016)	1190	153 500
4 <sup>th</sup> year (2016-2017)	1102	180 000
5 <sup>th</sup> year (2017-2018)	1000	/
<b>Total</b>	<b>5244</b>	<b>416 300</b>

I - *Scoledit* Project

II - Alignment/Comparison approach

III - Perspectives



# Children's writing specificities

I – *Scoedit* Project

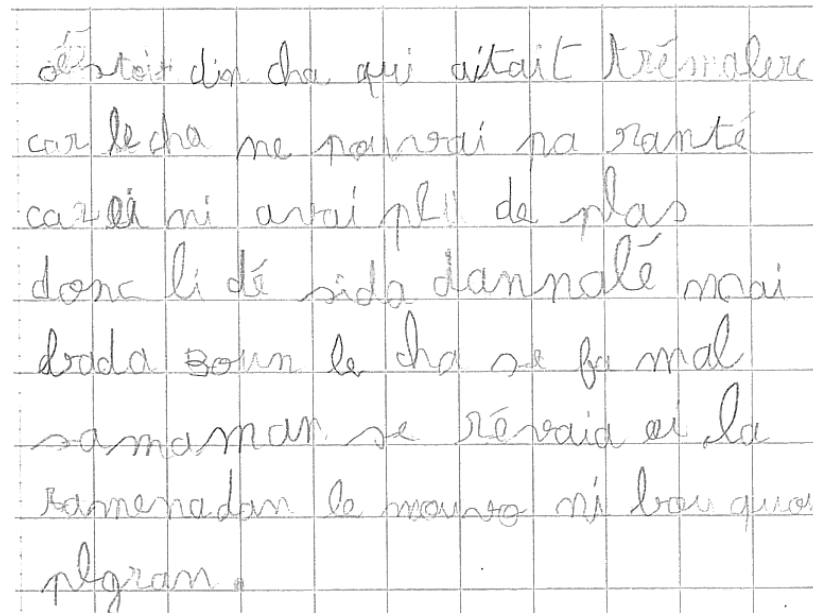
II – Alignment/Comparison approach

III – Perspectives

Children 297, at the end of the first year

« *sé stoir din cha qui aïtait tré malere car le cha ne pouvai pa ranté car li ni avai plu de plas donc li dé sida dannalé mai bada Boum le cha se fa mal sa maman se révaia ai la ramena dan le nouvo ni bou quou plgran.* »

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.



*sé stoir din cha qui aïtait tré malere  
car le cha ne pouvai pa ranté  
car li ni avai plu de plas  
donc li dé sida dannalé mai  
bada boum le cha se fa mal  
sa maman se révaia ai la  
ramena dan le nouvo ni bou quou  
plgran.*



# Children's writing specificities

I – *Scoedit Project*

II – Alignment/Comparison approach

III – Perspectives

Children 297, at the end of the third year

« Il était une fois une petite chatte qui *saplai* Boubou, on l'avait *apelée* Boubou par *cequ* elle adorait faire des blages. Mais un jour, cette petite *farsez* a *u* l'idée de *dir* une blage *mechate* aux policier. Mais *plus quel* etait toute petite elle ne *savai* pas *quel* *pouvé* aller en prison. *mai* quand elle *vu* les *abi* du monsieur *tombée* elle *vu* que *seté* un *robeau* qui lui *di* « C est *carnavale* aujourd'hui » *vien* je vais te *dégiser*. Et il *ver* la fête tout la journée, et *depuis* *se* jour le chat et le *robo* ne se *quiter* plus "FIN" »

Il était une fois une petite chatte qui s'appelait Boubou, on l'avait appelée Boubou parce qu'elle adorait faire des blagues. Mais un jour, cette petite farceuse a eu l'idée de dire une blague méchante aux policiers. Mais puisqu'elle était toute petite elle ne savait pas qu'elle pouvait aller en prison. Mais quand elle vit les habits du monsieur tomber elle vit que c'était un robot qui lui dit "C'est carnaval aujourd'hui" viens je vais te déguiser. Et il fait la fête toute la journée, et depuis ce jour le chat et le robot ne se quittèrent plus "FIN"





## II – Analysis of school corpora : Alignment/Comparison approach



# Processing chain

I – Scoledit Project

II – Alignment/Comparison approach

III – Perspectives

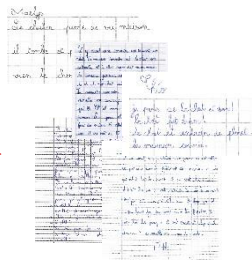


Data collection (protocol)



**Digitalisation**

Scanning,  
Storage, backup,  
loss of data

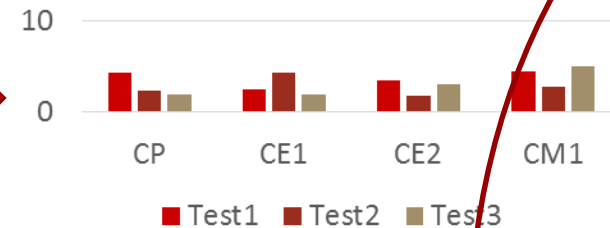


**Transcription**

Choices (loss of data),  
A description language,  
Interpretation, errors

La sorcière cherche  
sons chat / dans  
toute sa maison mes  
/ trouve pas sons  
chat. /# Alors elle  
se mes a pleurer /  
a <revision/>  
<revision/> pleurer  
qu'elle / tonbe par  
terre. Tout a coup  
/ Elle entant la  
porte souvrir /  
Elle dit « c'est  
mon chat. » / Sons  
chat sote dans  
c'est / <revision/>  
bras est il fut  
eurreu / jusca la  
fin des  
<revision/>tans.

Analyses



Alignment /  
comparison  
modules



```
<BE-ILV-008-005>
Christian Bex en
direct du
<F><MAJ><NOP>
#Vendée$ vendée
</NOP></MAJ></F>
globe challenge, une
aventure périlleuse
<G><CLA><POR> #qui$
qu'
</POR></CLA></G>est
<G><CLA><ADE> #une$
la </ADE></CLA></G>
course à la voile.
```

Extrait corpus Frida  
(Granger, 2001)

**Annotation**

Data enrichment,  
Some errors



# Annotation methods

- Classical approach

```

- <item type="forme" pos="19">
  - <f>
    s(remplacement)(b)(del)é(/del)(/b)(a)e(/a)(/remplacement)rai
  </f>
  <c>CAT</c>
  <l>serai</l>
</item>
- <item type="delim" pos="20">
  <f> </f>
  <c>DELIM</c>
  <l> </l>
</item>
- <item type="forme" pos="21">
  - <f>
    (remplacement)(b)(del)inventore(/del)(/b)(a)inventerr(/a)
  </f>
  <c>CAT</c>
  <l>inventeur</l>
</item>

```

```

<BE-ILV-008-005>
Christian Bex en direct du <F><MAJ><NOP> #Vendée$
vendée </NOP></MAJ></F> globe challenge, une aventure
périlleuse <G><CLA><POR> #qui$ qu' </POR></CLA></G>est
<G><CLA><ADE> #une$ la </ADE></CLA></G> course à la
voile.

```

```

Je pense qu'aller à l'école en vélo et bien et en effet f
ça ne <polue>_<pollue> pas. Le vélo <et>_<est> très bon pour la santé. §
En premier lieu ça permet de mieux <conaitre>_<connaître> la f
nature. De plus on peut sauter par dessus les <obstacles>_<obstacles> §
Mais quand il n'y a pas de <trottoir>_<trottoir> ou de piste f
<cyclabes>_<cyclable> le vélo peut être dangereux. §
Finalement, je pense donc que le vélo est une f
bonne idée. §

```

# Annotation methods

- **Classical approach**

```

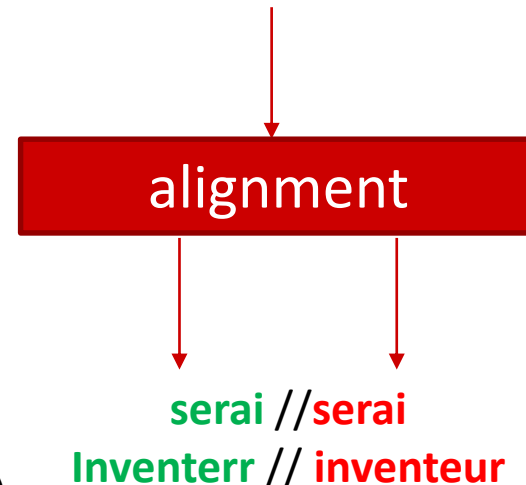
<item type="forme" pos="19">
  <f>
    s(remplacement)(b)(del)é(/del)(/b)(a)e(/a)(/remplacement)rai
  </f>
  <c>CAT</c>
  <l>serai</l>
</item>
<item type="delim" pos="20">
  <f> </f>
  <c>DELIM</c>
  <l> </l>
</item>
<item type="forme" pos="21">
  <f>
    (remplacement)(b)(del)inventore(/del)(/b)(a)inventerr(/a)(/remplacement)
  </f>
  <c>CAT</c>
  <l>inventeur</l>
</item>

```

- **Scoledit approach**

Transcription: serai inventerr

Normalisation: serai inventeur



# Scoledit approach: normalisation

Children 651 transcription: " *sé stoir din cha qui aitait tré malere ...* "

Children 651 normalisation: " C'est l'histoire d'un chat qui était très malheureux... "

- Advantages
  - No heavy formalism
  - Simpler than annotation
  - Reusability of procedures: e.g. dictation, rewriting exercises...
  - More freedom in the exploitation?
- Disadvantages
  - Automatic alignment => increased risk of error

## Related works

- Dictations
  - Santiago-Oriola (1998)
  - Beaufort et Roekhaut (2011)
- Translation
  - Desmet, Hérogue, 2005

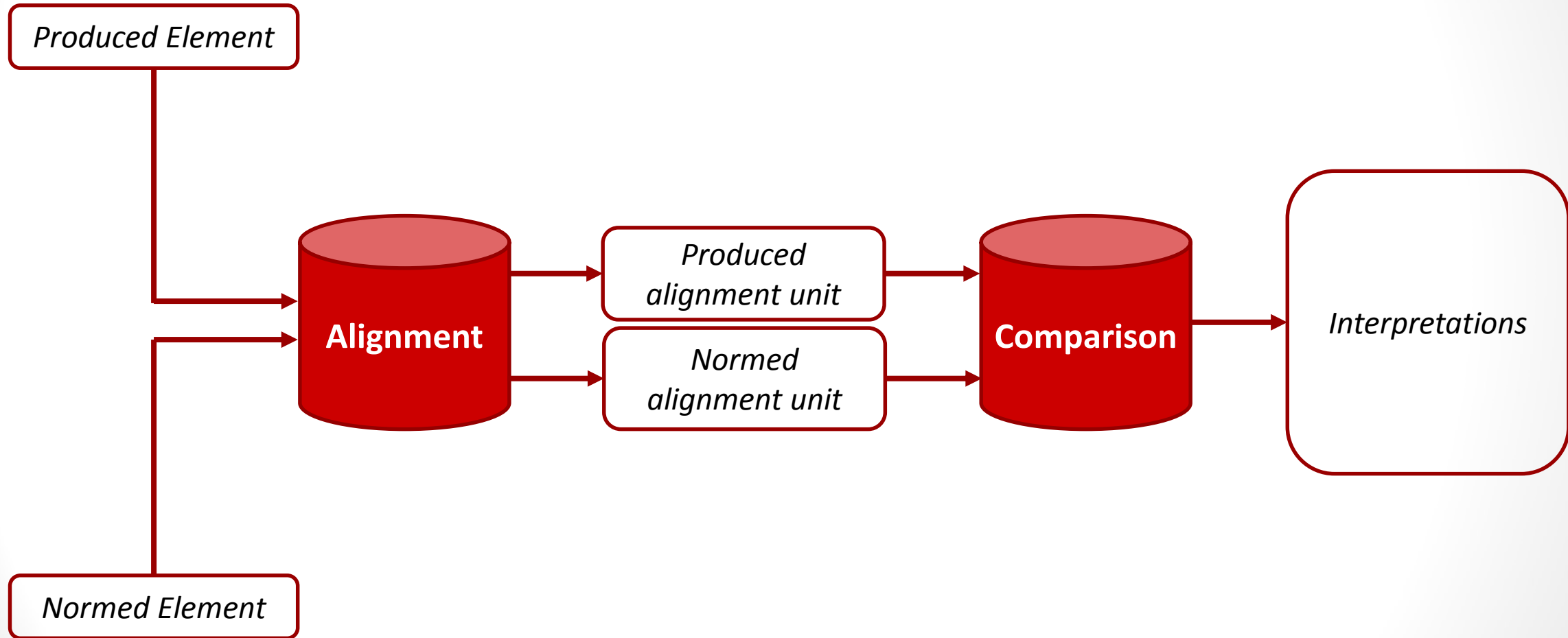


# Alignment/Comparison approach

I – Scoledit Project

II – Alignment/Comparison approach

III – Perspectives

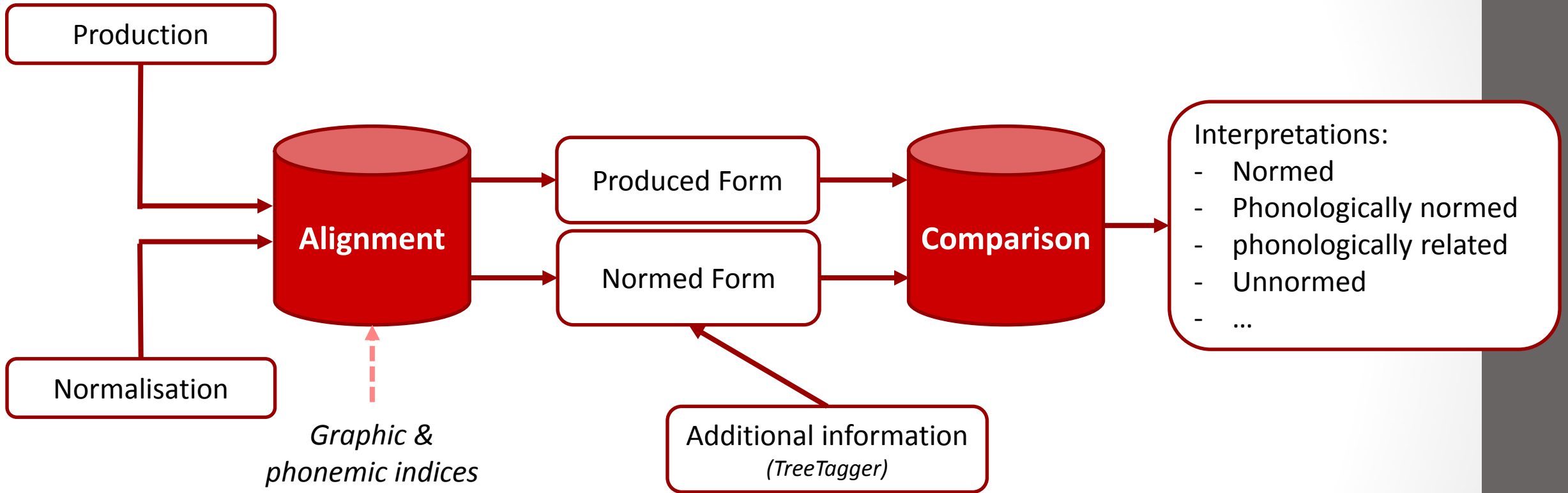


# Alignment/Comparison approach

I – Scoedit Project

II – Alignment/Comparison approach

III – Perspectives



Transcription : <i>din cha</i> → /dɛ̃ʃa/	➔	d	d'	PRP	de	phonologically normed	Under-segmented
Normalisation : <i>d'un chat</i> → /dœ̃ʃa/		in	un	DET:ART	un	phonologically related	Under-segmented
		cha	chat	NOM	chat	phonologically normed	Well-segmented

# Alignment/Comparison approach

I – *Scoledit* Project

II – Alignment/Comparison approach

III – Perspectives

**Vocabulaire CP**

477 mots distincts sur 12852 répertoriés à ce jour (3,71%).

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
	166	miaou	156
réveiller	149	de	147
qui	139	marche	120
dormir	119	chaton	118
mal	116	foi fois	96
bébé	91	sur	91

**Extract from Scoledit's website**





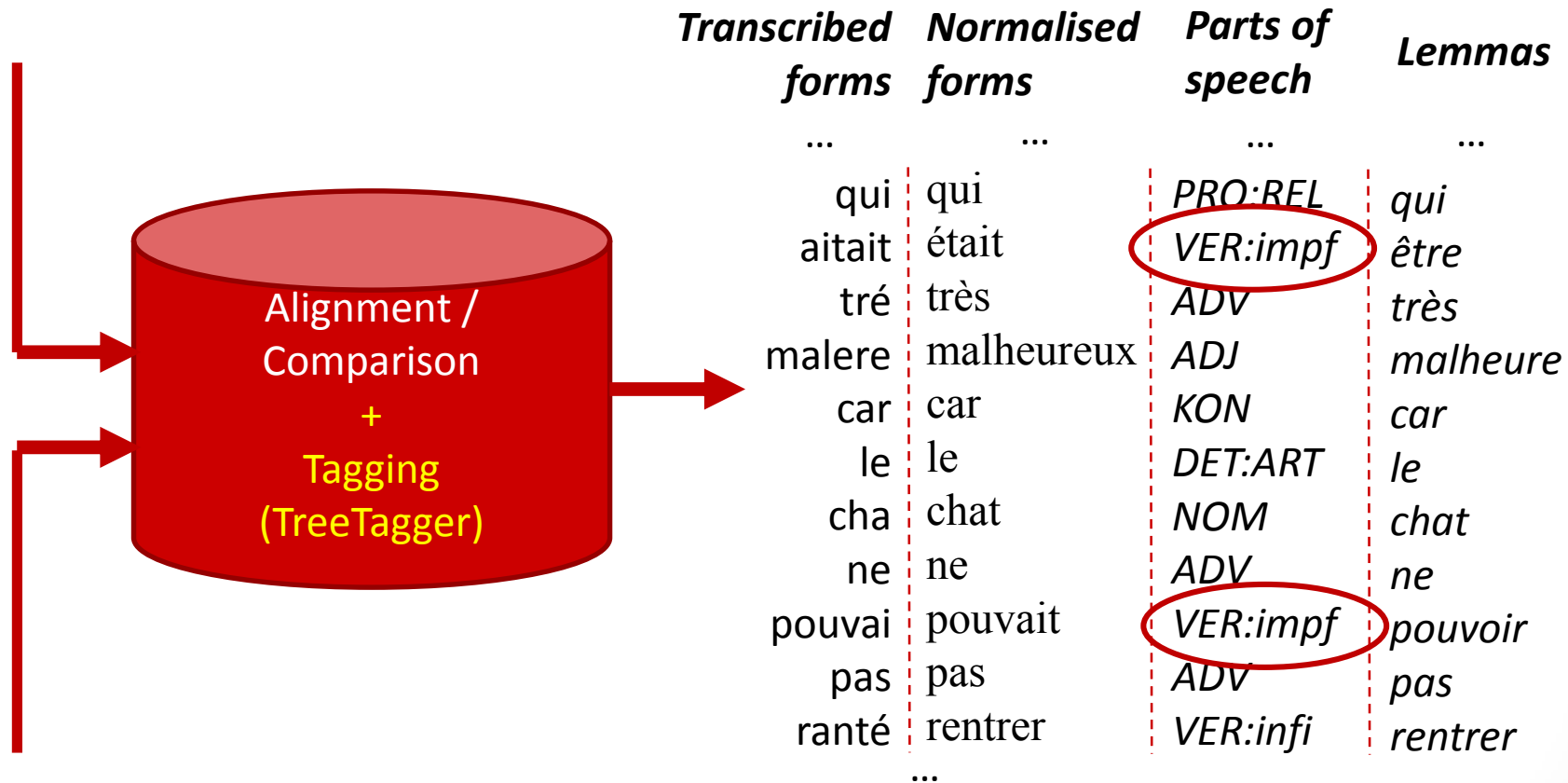
# Alignment/Comparison approach

I – Scoedit Project

II – Alignment/Comparison approach

III – Perspectives

Children 651 transcription: " sé stoir din cha qui aïtait tré malere car le cha ne pouvai pa ranté ... "



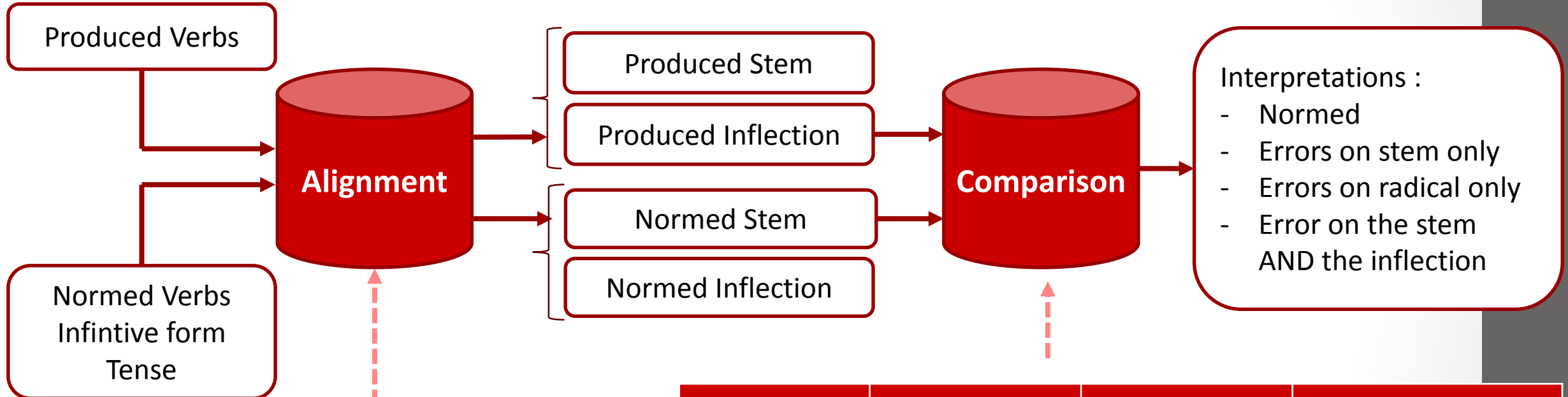
Children 651 normalisation: " C'est l'histoire d'un chat qui était très malheureux car le chat ne pouvait pas rentrer... "

# Alignment/Comparison approach

I – Scoledit Project

II – Alignment/Comparison approach

III – Perspectives



*Linguistic theory  
(Martinet, 1979 ; Meleuc et  
Fauchart, 1999; Blanche-  
Benveniste, 2002 ; Pellat,  
2009)*

Normed	Error on the stem only	Error on the inflection only	Error on the stem AND the inflection
di + t // di + t	fai+re // fé+re	av+ait // av+ai	miaul+e // miol + ∅

# Alignment/Comparison approach

I – *Scoedit* Project

II – Alignment/Comparison approach

III – Perspectives

statut	forme	baseForme	desiForme	prod	baseProd	desiProd
normé	suis	suiv	ant	suis	suiv	ant
phonologi	fais	fai	s	fait	fai	t
non normé	était	ét	ait	étét	ét	ét
équivalen	faire	fai	re	fère	fé	re
normé	sais	sai	s	sais	sai	s
normé	dit	di	t	dit	di	t
normé	était	ét	ait	était	ét	ait
normé	est	es	t	est	es	t



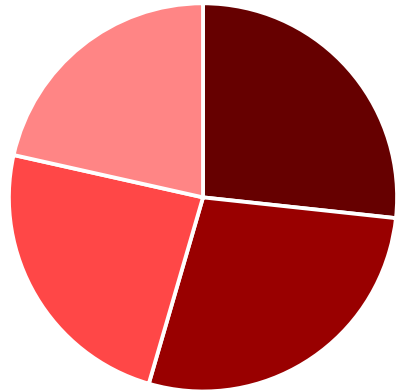
# Analysis of verbal morphography

I – *Scoedit* Project

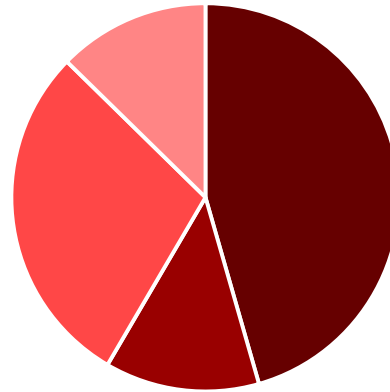
II – Alignment/Comparison approach

III – Perspectives

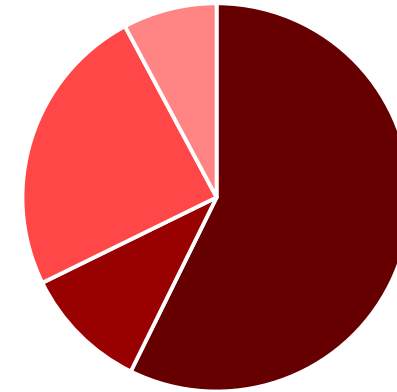
1st year



2nd year



3rd year



- Normed
- stem only
- termination only
- stem and termination

- Normed
- stem only
- termination only
- stem and termination

- Normed
- stem only
- termination only
- stem and termination

- Brissaud C., Totereau C., Ponton C., Wolfarth C. (to be published in 2018). Usage d'un corpus longitudinal, le cas de la morphologie verbale. *Revue Repères* : collecter, interpreter, enseigner l'écriture. *Analyses linguistiques des écrits d'élèves*.



# III – Perspectives

- Method development
  - Treatment chain : evaluation des différents modules
  - quality of the data : Checking of the transcriptions and normalisations
- New modules
  - graphemes - phonemes correspondences
  - Agreement in noun groups
- E-Calm Project
  - Textual cohesion and coherence
  - Inflectional morphology and derivational morphology
- Purpose : describe linguistic uses in these corpora and their evolutions



# Bibliography

- BEAUFORT R., ROEKHAUT S. (2011). AUTOMATION OF DICTATION EXERCISES. *COMPUTATIONAL LINGUISTICS IN THE NETHERLANDS JOURNAL*, (1), 1-20.
- 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.
- 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*, SCéren, CRDP de Versailles.
- JUEL, C. (1988). Learning to read and write: A longitudinal study of 54 children from first through fourth grades. *Journal of educational Psychology*, 80(4), 437.
- 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.
- SMITH, N., & MCENERY, T. (1998). Issues in Transcribing a Corpus of Children's Handwritten Projects. *Literary and linguistic computing*, 13(4), 217-225.

