

Pro-TEXT: Annotated Corpus of Keystroke Logs

Aleksandra Miletić¹, Georgeta Cislaru², Christophe Benzitoun³ and Santiago Herrera-Yanez¹

1 Paris 3 Sorbonne Nouvelle University

2 Paris Nanterre University

3 Lorraine University

Language Resources and Evaluation Conference, June 21-23 2022



Evolution of a text:

afin d'aborder un projet de l'Université de Poitiers

in order to address a project of the University of Poitiers

Evolution of a text:

afin d'aborder un projet de l'Université de Poitiers

in order to address a project of the University of Poitiers

~~afin d'aborder un projet de l'Université de Poitiers~~

~~*in order to address a project of the University of Poitiers*~~

Evolution of a text:

afin d'aborder un projet de l'Université de Poitiers

in order to address a project of the University of Poitiers

afin d'aborder un projet de ~~l'Université de Poitiers~~

in order to address a project ~~of the University of Poitiers~~

afin d'aborder un projet **songé par l'Universitté**

*in order to address a project **conceived by the Universitty***

Evolution of a text:

afin d'aborder un projet de l'Université de Poitiers

in order to address a project of the University of Poitiers

afin d'aborder un projet de ~~l'Université de Poitiers~~

in order to address a project ~~of the University of Poitiers~~

afin d'aborder un projet **songé par l'Universitté**

*in order to address a project **conceived by the Universitty***

afin d'aborder un projet songé par l'Universitté

in order to address a project conceived by the Universitty

Evolution of a text:

afin d'aborder un projet de l'Université de Poitiers
in order to address a project of the University of Poitiers

afin d'aborder un projet ~~de l'Université de Poitiers~~
in order to address a project ~~of the University of Poitiers~~

afin d'aborder un projet **songé par l'Universitté**
*in order to address a project **conceived by the Universitty***

afin d'aborder un projet songé par l'Universitté
in order to address a project conceived by the Universitty

afin d'aborder un projet songé par l'Université **de Poitiers**
*in order to address a project conceived by the University **of Poitiers***

Evolution of a text:

afin d'aborder un projet de l'Université de Poitiers
in order to address a project of the University of Poitiers

afin d'aborder un projet ~~de l'Université de Poitiers~~
in order to address a project ~~of the University of Poitiers~~

afin d'aborder un projet **songé par l'Universitté**
*in order to address a project **conceived by the Universitty***

afin d'aborder un projet songé par l'Universitté
in order to address a project conceived by the Universitty

afin d'aborder un projet songé par l'Université **de Poitiers**
*in order to address a project conceived by the University **of Poitiers***

afin d'aborder un projet songé par l'Université ~~de Poitiers~~
in order to address a project conceived by ~~the University of Poitiers~~

Evolution of a text:

afin d'aborder un projet de l'Université de Poitiers
in order to address a project of the University of Poitiers

afin d'aborder un projet ~~de l'Université de Poitiers~~
in order to address a project ~~of the University of Poitiers~~

afin d'aborder un projet **songé par l'Universitté**
*in order to address a project **conceived by the Universitty***

afin d'aborder un projet songé par l'Universitté
in order to address a project conceived by the Universitty

afin d'aborder un projet songé par l'Université **de Poitiers**
*in order to address a project conceived by the University **of Poitiers***

afin d'aborder un projet songé par l'Université ~~de Poitiers~~
in order to address a project conceived by ~~the University of Poitiers~~

afin d'aborder un projet songé par **notre Université**
*in order to address a project conceived by **our University***

Keystroke logs: recordings of the writing process executed through a keyboard using dedicated software (Leijten and Van Waes, 2006; Strömqvist and Malmsten, 1998; Carl, 2012)

- Character additions, deletions, substitutions
- Behavioural information (pause length, production speed)
- → Study of writing itself and of underlying cognitive processes

Linguistic annotation?

- Available for English and Dutch as part of Inputlog (Leijten et al., 2015)
- Remains rare, typically POS-tagging (Serbina et al., 2015; Carl et al., 2011)
- Highly relevant, cf. (Serbina et al., 2017)

- 1 Pro-TEXT Project
- 2 Pro-TEXT Corpus
- 3 Annotation Methodology
- 4 Results
- 5 Future Work

- 1 Pro-TEXT Project
- 2 Pro-TEXT Corpus
- 3 Annotation Methodology
- 4 Results
- 5 Future Work

Pro-TEXT: linguistic analysis of the textualization process, i.e. the real-time progressive construction of a text

- Funded by the French National Research Agency (ANR-18-CE23-0024-01)
- Interdisciplinary project:
 - **psycholinguistic experiments**: T. Olive, S. Bouriga, D. Chesnet, C. Perret, J. Pylouster and C. Bordes (CERCA, Poitiers University)
 - **linguistic description**: G. Cislaru, S. Fleury, F. Lefeuvre, D. Legallois, A. Boyer, Q. Feltgen and A. Miletic (CLESTHIA, Paris 3 University); C. Benzitoun and M. Dagnat (ATILF, Lorraine University); S. Vandaele (University of Montreal)
 - **machine learning modelling**: G. Cabanes, T. Charnois, N. Grozavu, J. Le Roux, P. Rastin, N. Rogovschi and N. Tomeh (LIPN, Paris 13 University)

Pro-TEXT: linguistic analysis of the textualization process, i.e. the real-time progressive construction of a text

- Funded by the French National Research Agency (ANR-18-CE23-0024-01)
- Interdisciplinary project:
 - **psycholinguistic experiments**: T. Olive, S. Bouriga, D. Chesnet, C. Perret, J. Pylouster and C. Bordes (CERCA, Poitiers University)
 - **linguistic description**: G. Cislaru, S. Fleury, F. Lefeuvre, D. Legallois, A. Boyer, Q. Feltgen and A. Miletic (CLESTHIA, Paris 3 University); C. Benzitoun and M. Dagnat (ATILF, Lorraine University); S. Vandaele (University of Montreal)
 - **machine learning modelling**: G. Cabanes, T. Charnois, N. Grozavu, J. Le Roux, P. Rastin, N. Rogovschi and N. Tomeh (LIPN, Paris 13 University)

→ Produce a corpus suitable for this type of research

- 1 Pro-TEXT Project
- 2 Pro-TEXT Corpus
- 3 Annotation Methodology
- 4 Results
- 5 Future Work

Collected Pro-TEXT Corpus:

Subcorpus	Texts	Words	Writers	Genre
Academic	26	70464	MA students	mini-thesis in linguistics
Professional	10	34504	social workers	reports on child protection
Experimental	165	63533	BA students	essays on different subjects
Children	183	20306	pupils (3 rd -6 th grade)	narrative texts and essays
Translation	38	13682	BA students	EN-FR translation of medical texts and original texts in FR
Total	422	202489	-	-

Recorded using Inputlog (Leijten and Van Waes, 2006) and Scriptlog (Strömqvist and Malmsten, 1998)

- 1 Pro-TEXT Project
- 2 Pro-TEXT Corpus
- 3 Annotation Methodology
- 4 Results
- 5 Future Work

Goal: maximize the utility of the corpus while respecting time constraints

- French Treebank tagsets (POS tags and dependency syntax) (Candito et al., 2009)
- Using automatic pre-annotation to facilitate work of human annotators
- Manual validation to ensure annotation quality
- Agile annotation (Voormann and Gut, 2008) to ensure manual annotation coherence
- Annotating **all** produced content (not only final texts)
 - → intermediate versions of each text

Intermediate Versions of a Text

Intermediate versions:

afin d'aborder un projet de l'Université de Poitiers
in order to address a project of the University of Poitiers

afin d'aborder un projet ~~de l'Université de Poitiers~~
in order to address a project ~~of the University of Poitiers~~

afin d'aborder un projet **songé par l'Universitté**
*in order to address a project **conceived by the Universitty***

afin d'aborder un projet songé par l'Universitté
in order to address a project conceived by the Universitty

afin d'aborder un projet songé par l'Université **de Poitiers**
*in order to address a project conceived by the University **of Poitiers***

afin d'aborder un projet songé par l'Université ~~de Poitiers~~
in order to address a project conceived by ~~the University of Poitiers~~

afin d'aborder un projet songé par **notre Université**
*in order to address a project conceived by **our University***

Intermediate Versions of a Text

Intermediate versions:

afin d'aborder un projet de l'Université de Poitiers → v1
in order to address a project of the University of Poitiers

afin d'aborder un projet ~~de l'Université de Poitiers~~ → v2
in order to address a project ~~of the University of Poitiers~~

afin d'aborder un projet **songé par l'Université** → v3
*in order to address a project **conceived by the University***

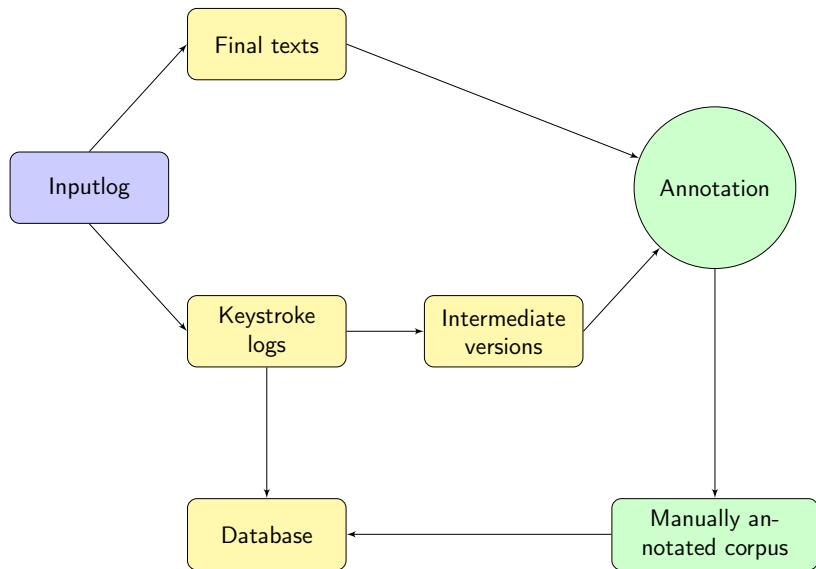
afin d'aborder un projet songé par l'Université → v4
in order to address a project conceived by the University

afin d'aborder un projet songé par l'Université **de Poitiers** → v5
*in order to address a project conceived by the University **of Poitiers***

afin d'aborder un projet songé par ~~l'Université de Poitiers~~ → v6
in order to address a project conceived by ~~the University of Poitiers~~

afin d'aborder un projet songé par **notre Université** → v7
*in order to address a project conceived by **our University***

Annotation Workflow



Annotated intermediate versions:

- CoNLL format:

```
#text_version_id=0
#sentence_id=1
1 C' ce CLS CLS - 2 suj_impers - 1=C__True|2='__True|
2 etait être V P - 0 root - 3=e__True|4=t__True|5=a__True|6=i__True|7=t__True
3 en en P - 2 p_obj - 9=e__True|10=n__True|
4 décembre décembre - NC NC - 3 prep - 12=d__True|13=é__True|14=c__True|1
5 je je CLS ADJ - 0 x - 21=j__True|22=e__True|
212 connaître - VINF VINF - 0 0 - 24=c__True|25=o__True|26=n__True|27=n__Tru
7 ce ce DET DET - 8 det - 35=c__True|36=e__True|
8 garçon garçon NC NC - 0 x - 38=g__True|39=a__True|40=r__True|41=ç__True|42=o__
9 depuis depuis P P - 0 x - 45=d__True|46=e__True|47=p__True|48=u__True|49=i__
10 la le DET DET - 11 det - 52=l__True|53=a__True|
11 marenelle maternelle - NC ADJ - 9 prep - 55=m__True|56=a__True|57=r__True|5
12 il il CLS CLS - 0 x - 65=i__True|66=l__True|
213 etea - ADV ADV - 0 - 68=e__True|69=t__False|70=e__False|71=a__False|
```

```
#text_version_id=1
#sentence_id=1
1 C' ce CLS CLS - 2 suj_impers - 1=C__True|2='__True|
2 etait être V P - 0 root - 3=e__True|4=t__True|5=a__True|6=i__True|7=t__True
3 en en P - 2 p_obj - 9=e__True|10=n__True|
4 décembre décembre - NC NC - 3 prep - 12=d__True|13=é__True|14=c__True|1
5 je je CLS ADJ - 0 x - 21=j__True|22=e__True|
212 connaître - VINF VINF - 0 0 - 24=c__True|25=o__True|26=n__True|27=n__Tru
7 ce ce DET DET - 8 det - 35=c__True|36=e__True|
8 garçon garçon NC NC - 0 x - 38=g__True|39=a__True|40=r__True|41=ç__True|42=o__
9 depuis depuis P P - 0 x - 45=d__True|46=e__True|47=p__True|48=u__True|49=i__
10 la le DET DET - 11 det - 52=l__True|53=a__True|
11 marenelle maternelle - NC ADJ - 9 prep - 55=m__True|56=a__True|57=r__True|5
12 il il CLS CLS - 0 x - 65=i__True|66=l__True|
214 e - ADV ADV - 0 - 68=e__True|
```

Annotated database:

- database-like CSV format:

n_event	st_time	end_time	pause	event_pos	C	doc_len	o	type_o	charID	tokenID	token	lemma	POS	XPOS	ms	governor	function	charStatus	tokenStatus	sentence_id		
1	144098	144597			C	0	0	1	keyboard	1-1=10=11=12=	1-C=0=C1=0=1=ce0=ce1=	1-CLS0=CL9=1-CLS0=CL9=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=	1=
2	146235	146313	1638r		C	1	1	1	keyboard	2-1=10=11=12=	1-C=0=C1=0=1=ce0=ce1=	1-CLS0=CL9=1-CLS0=CL9=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=	1=
3	146641	146719	328r		C	2	2	2	1	keyboard	3-1=20=21=22=	1-etat0=etat1=1=ete0=ete1=	1-V0=V1=V7=1-P0=P1=PY=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
4	146922	146984	203r		C	3	3	3	1	keyboard	4-1=20=21=22=	1-etat0=etat1=1=ete0=ete1=	1-V0=V1=V7=1-P0=P1=PY=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
5	147358	147421	374a		C	4	4	4	1	keyboard	5-1=20=21=22=	1-etat0=etat1=1=ete0=ete1=	1-V0=V1=V7=1-P0=P1=PY=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
6	147577	147702	156i		C	5	5	5	1	keyboard	6-1=20=21=22=	1-etat0=etat1=1=ete0=ete1=	1-V0=V1=V7=1-P0=P1=PY=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
7	147904	147982	202t		C	6	6	6	1	keyboard	7-1=20=21=22=	1-etat0=etat1=1=ete0=ete1=	1-V0=V1=V7=1-P0=P1=PY=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
8	148372	148407	390r		C	7	7	7	1	keyboard	8na	na	na	na	na	na	na	na	na	na	na	na
9	150668	150993	2371r		C	8	8	8	1	keyboard	9-1=30=31=32=	1-en0=en1=1=en0=en1=	1-P0=P1=PY=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
10	151102	151196	109r		C	9	9	9	1	keyboard	10-1=30=31=32=	1-en0=en1=1=en0=en1=	1-P0=P1=PY=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
11	151321	151414	125r		C	10	10	10	1	keyboard	11na	na	na	na	na	na	na	na	na	na	na	na
12	151758	151836	344d		C	11	11	11	1	keyboard	12-1=40=41=42=	1=decembre0=1=decembre0=	1=NC0=NC1P=1=NC0=NC1P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
13	152959	153052	1123f		C	12	12	12	1	keyboard	13-1=40=41=42=	1=decembre0=1=decembre0=	1=NC0=NC1P=1=NC0=NC1P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
14	153364	153427	312r		C	13	13	13	1	keyboard	14-1=40=41=42=	1=decembre0=1=decembre0=	1=NC0=NC1P=1=NC0=NC1P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
15	153379	153817	612r		C	14	14	14	1	keyboard	15-1=40=41=42=	1=decembre0=1=decembre0=	1=NC0=NC1P=1=NC0=NC1P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
16	154456	154612	339m		C	15	15	15	1	keyboard	16-1=40=41=42=	1=decembre0=1=decembre0=	1=NC0=NC1P=1=NC0=NC1P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
17	157795	157920	3183b		C	16	16	16	1	keyboard	17-1=40=41=42=	1=decembre0=1=decembre0=	1=NC0=NC1P=1=NC0=NC1P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
18	158200	158263	280r		C	17	17	17	1	keyboard	18-1=40=41=42=	1=decembre0=1=decembre0=	1=NC0=NC1P=1=NC0=NC1P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
19	158419	158481	156r		C	18	18	18	1	keyboard	19-1=40=41=42=	1=decembre0=1=decembre0=	1=NC0=NC1P=1=NC0=NC1P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
20	159121	159246	640r		C	19	19	19	1	keyboard	20na	na	na	na	na	na	na	na	na	na	na	na
21	165345	165470	6099r		C	20	20	20	1	keyboard	21-1=50=51=52=	1=je0=je1=1=je0=je1=	1=CLS0=CL9=1=ADJ0=ADJ=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
22	166642	166720	172r		C	21	21	21	1	keyboard	22-1=50=51=52=	1=je0=je1=1=je0=je1=	1=CLS0=CL9=1=ADJ0=ADJ=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
23	166938	166932	218r		C	22	22	22	1	keyboard	23na	na	na	na	na	na	na	na	na	na	na	na
24	166250	166344	218r		C	23	23	23	1	keyboard	24-1=60=2121=	1=connaissae=1=connaisse0=	1-V0=VNF0P=1-V0=VNF0P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
25	166749	166843	425r		C	24	24	24	1	keyboard	25-1=60=2121=	1=connaissae=1=connaisse0=	1-V0=VNF0P=1-V0=VNF0P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
26	167092	167186	249r		C	25	25	25	1	keyboard	26-1=60=2121=	1=connaissae=1=connaisse0=	1-V0=VNF0P=1-V0=VNF0P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
27	167514	167592	328r		C	26	26	26	1	keyboard	27-1=60=2121=	1=connaissae=1=connaisse0=	1-V0=VNF0P=1-V0=VNF0P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
28	167694	167872	202a		C	27	27	27	1	keyboard	28-1=60=2121=	1=connaissae=1=connaisse0=	1-V0=VNF0P=1-V0=VNF0P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
29	168808	168949	936i		C	28	28	28	1	keyboard	29-1=60=2121=	1=connaissae=1=connaisse0=	1-V0=VNF0P=1-V0=VNF0P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
30	170275	170368	1325r		C	29	29	29	1	keyboard	30-1=60=2121=	1=connaissae=1=connaisse0=	1-V0=VNF0P=1-V0=VNF0P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
31	171070	171164	702r		C	30	30	30	1	keyboard	31-1=60=2121=	1=connaissae=1=connaisse0=	1-V0=VNF0P=1-V0=VNF0P=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=
32	171414	171507	250r		C	31	31	31	1	keyboard	32=2121=21220=	connaissae0=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=	1=
33	171663	171757	156r		C	32	32	32	1	keyboard	33=0=2121=21220=	connaissae0=	1=	1_0=	1=	1=	1=	1=	1=	1=	1=	1=
34	171928	172022	171r		C	33	33	33	1	keyboard	34na	na	na	na	na	na	na	na	na	na	na	na

- 1 Pro-TEXT Project
- 2 Pro-TEXT Corpus
- 3 Annotation Methodology
- 4 Results**
- 5 Future Work

Annotated Final Texts:

Subcorpus	Texts	Sentences	Tokens	Tok/sent	Lemmas	Types
Children	120	440	11873	27.0	1171	224
Experimental	15	149	5719	38.4	1024	1428
Translation	10	149	3675	24.7	626	887
Academic	2	474	8879	18.7	1523	2013
Professional	0	0	0	0	0	0
TOTAL	147	1212	30146	24.9	3062	5184

Annotated Intermediate Versions:

48 texts

	Sentences	Tokens	Lemmas	Types
Final texts	71	4319	632	1074
Intermediate versions	4621	128518	693	1767

Download from <https://pro-text.huma-num.fr/ressources/> under **CC BY-NC-SA 4.0** license

- 1 Pro-TEXT Project
- 2 Pro-TEXT Corpus
- 3 Annotation Methodology
- 4 Results
- 5 Future Work

Corpus annotation:

- Complete the annotation of intermediate versions
- Expand the annotation of the whole corpus: bootstrapping?

Corpus analysis:

- Linguistic profiling of writing bursts
- Relationship between linguistic structure and pauses
- Behaviour of syntactic dependencies with respect to behavioural factors
- Modelling interactions between behavioural and linguistic information using machine learning techniques

- Marie Candito, Benoît Crabbé, and Mathieu Falco. Dépendances syntaxiques de surface pour le français, 2009.
- Michael Carl. Translog-ii: a program for recording user activity data for empirical reading and writing research. In LREC, volume 12, pages 4108–4112, 2012.
- Michael Carl, Barbara Dragsted, Jakob Elming, Daniel Hardt, and Arnt Lykke Jakobsen. The process of post-editing: A pilot study. Copenhagen Studies in Language, 41:131–142, 2011.
- Mariëlle Leijten and Luuk Van Waes. Inputlog: New perspectives on the logging of on-line writing processes in a windows environment. In Computer key-stroke logging and writing, pages 73–93. Brill, 2006.
- Mariëlle Leijten, Luuk Van Waes, and Eric Van Horenbeeck. Writing(s) at the Crossroads: The Process-Product Interface, chapter Analyzing writing process data: A linguistic perspective, pages 277–302. John Benjamins Publishing Company, 2015.
- Tatiana Serbina, Paula Niemietz, Matthias Fricke, Philipp Meisen, and Stella Neumann. Part of speech annotation of intermediate versions in the keystroke logged translation corpus. In Proceedings of the 9th Linguistic Annotation Workshop, pages 102–111, 2015.

- Tatiana Serbina, Sven Hintzen, Paula Niemietz, and Stella Neumann. Empirical modelling of translation and interpreting, volume 3, chapter Changes of word class during translation—Insights from a combined analysis of corpus, keystroke logging and eye-tracking data, pages 177–208. Language Science Press, Berlin, 2017.
- Sven Strömqvist and Lars Malmsten. Sriptlog pro 1.04: User's manual. Technical report, University of Göteborg, 1998.
- Holger Voormann and Ulrike Gut. Agile corpus creation. Corpus Linguistics and Linguistic Theory, 4(2):235–251, 2008.