# Frantext POS annotation

**bertrand.gaiffe@atilf.fr**

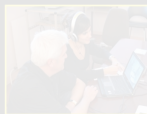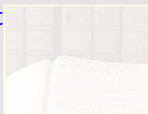march 16th, 2021

# Frantext ?

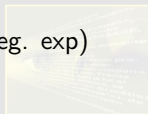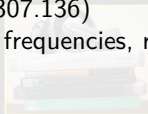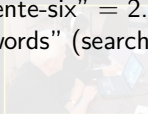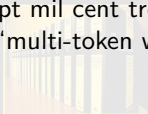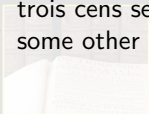- Corpus of around 5500 works (read "book")
- mainly 20$^{th}$ and 19$^{th}$ century
- but, all periods of french language covered (not evenly)
- the whole corpus is tagged in Parts Of Speech and lemmas
- interrogation through a CQL based query language

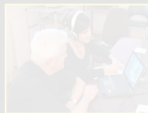Frantext

– a group of persons (including S. Ollinger, C. Benzitoun and L. Berrard)
  - ✓ defined a set of categories
  - ✓ built a gold corpus
  - ✓ learned a model with Talismane
– some particularities :
  - ✓ numbers are represented as a single token (even "deux millions trois cens sept mil cent trente-six" = 2.307.136)
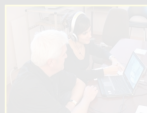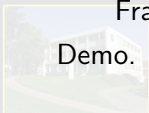  - ✓ some other "multi-token words" (search frequencies, reg. exp)

# Talismane

- We use a rather old version of Talismane (A. Uriely)
- takes "sentences" as entry and produces a conll style output
- takes some time...
- in fact, Talismane is a dependency parser, we use only the POS tagging stage.

- I. Clément (and G. Toubiana) scan and OCR the texts
- they encode them in TEI (Text Encoding Initiative) using an XML editor
- they use a script that tag the texts using Talismane in the background
- and finally they transform the resulting annotated TEI into Frantext format.

Demo.

# The problems

- Talismane (the version we use) relies too much on its sentence splitting
- Natural Language Processing progressed a lot
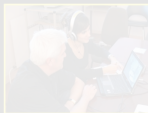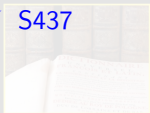- the annotation process was designed for texts pots 1850 only

# An attempt at Bert'izing the annotation

– the gold corpus is divided (as usual) in train, dev and test

– learned a "camembert" based POS tagger

– 97% on test (talismane does slightly better but was trained on test also)

– comparison between the two annotations
  - ✓ R467
  - ✓ S437

– tagging persNames and placeNames ?
– design models for other periods of the french language
  ✓ {vieux Camem | Hau}bert a bert model for middle age french.

Contraints :

– have to scale to Frantext (5500 texts),
– have to provide a tool for the people that enter new texts into Frantext.