



**HAL**  
open science

## The #Idéo2017 platform

Julien Longhi

► **To cite this version:**

Julien Longhi. The #Idéo2017 platform. Journée d'étude CORLI: Traitements et standardisation des corpus multimodaux et web 2.0., May 2018, Paris, France. hal-01873815

**HAL Id: hal-01873815**

**<https://hal.science/hal-01873815v1>**

Submitted on 13 Sep 2018

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



General context: the CoMeRe project

- ❖ Aims to build a kernel corpus of computer-mediated communication (CMC) genres in French
- ❖ Mono and multimodal interactions stemming from networks including the Internet and telecommunications that may be synchronous or asynchronous
- ❖ Members had previously collected and structured different types of CMC corpora within their local teams (in a variety of formats with disparities in corpus compilation choices)
- ❖ Corpora are structured and referred to in a *uniform* way in order that they may form part of the forthcoming *French National Reference Corpus*



Project website

First project: the Polititweet corpora

- ❖ Development of the Interaction Space (IS) model to model CMC interaction (Chanier & Jin, 2013).
- ❖ Includes descriptions of time, set of participants, online location(s) defined by the properties of the set of environments used by participants.
- ❖ Description of the IS within the TEI header and messages and turns encoded in the TEI body using a common <post> element

```
<post xmlns="cmr-polititweets-a392322730999046144" xmlns:cmr="cmr-polititweets-g13121164"
xmlns="2013-10-21T18:13:22" xmlns:tei="http://www.tei-c.org/ns/1.0" type="text" id="1"
xmlns:xsi="http://www.w3.org/2001/XMLSchema-instance" xsi:schemaLocation="http://www.tei-c.org/ns/1.0 cmr-polititweets-g13121164.xsd">
  <text>
    <span>pourrait voir le jour en Poitou-Charentes.</span>
  </text>
  <trailer>
    <cmr:medium>
      <cmr:settingweb/>
    </cmr:medium>
    <cmr:tweetcount>
      <cmr:numeric value="1"/>
    </cmr:tweetcount>
  </trailer>
</post>
```

Openess

- ❖ Corpora released as open-data – paves way for scientific examination, replication and cumulative analyses
- ❖ Released on ORTOLANG (French equivalent of DARIAH the European infrastructure for Humanities)
- ❖ Bibliographic reference created for each corpus and given in <titleSmt> of TEI header. e.g. Longhi J., Marinica C., Borzic B., Alkhouli A., 2014: Polititweets, corpus de tweets provenant de comptes politiques influents. In Chanier T. (ed) Banque de corpus CoMeRe. Ortolang.fr : Nancy. [cmr-polititweets- tei-v1]. https://repository.ortolang.fr/api/content/comere/v3.3/cmr-polititweets.html



Corpus depository

Interests and limits

- ❖ The analysis of political tweets during the election campaigns, or specific events, is increasing and can be seen as a specific type of political discourse (Longhi, 2013)
- ❖ Extensive literature on the analysis of political tweets: but these works are difficult to gather because they come either from the computer sciences, either from the humanities and social sciences.
- ❖ Despite the unquestionable interest in outlining political facts, citizens do not have access to these results.

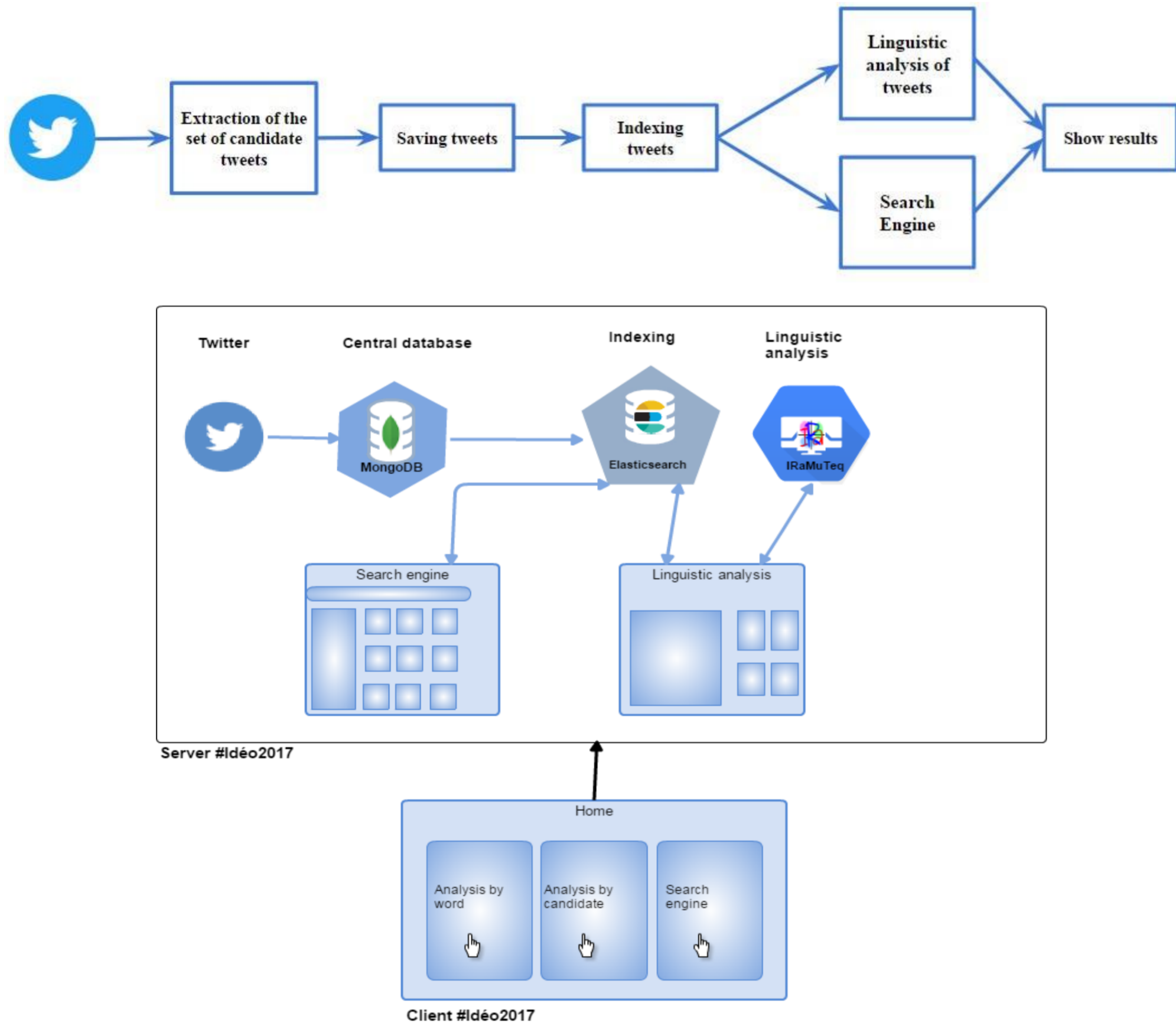
Challenges of #Idéo2017

- ❖ #Idéo2017 is a web platform available online allowing the analysis of the messages, posted on Twitter, related to political news.
- ❖ Its objective is to make available on the web for average citizens a set of statistical analyses and data visualization tools applied to the Twitter messages. #Idéo2017 allowed French electors to analyze the discourse of the candidates by means of their tweets.
- ❖ Citizens can make their own queries (based on linguistic and textometric criteria, more precisely, the most commonly used words by political personalities, analyses of similarities, ALCESTE algorithm, etc.) and obtain comprehensible results.



Tools and analysis

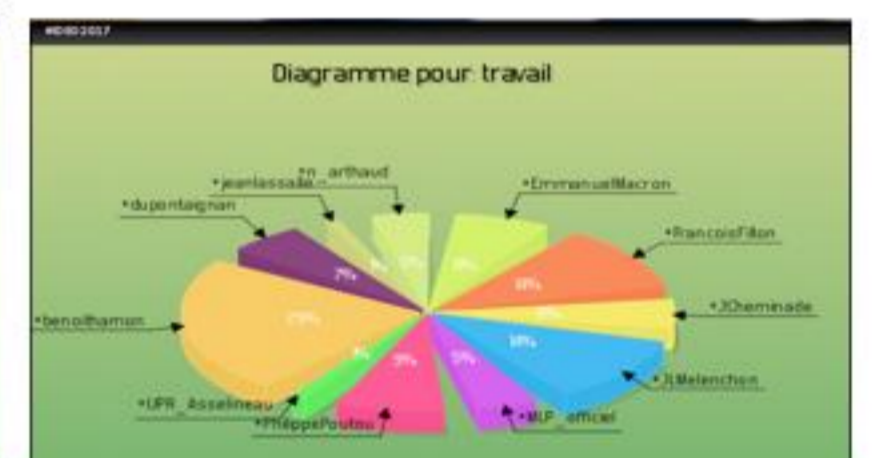
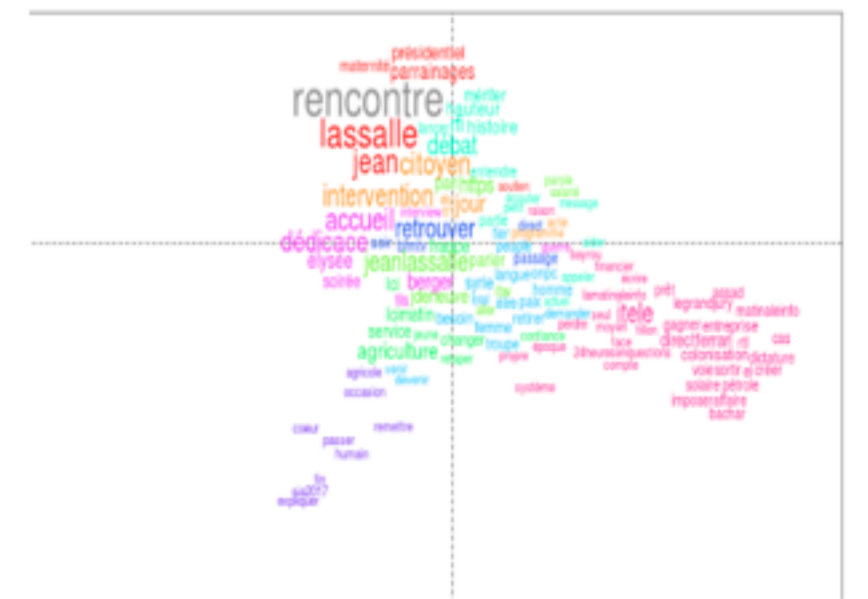
Tool Development



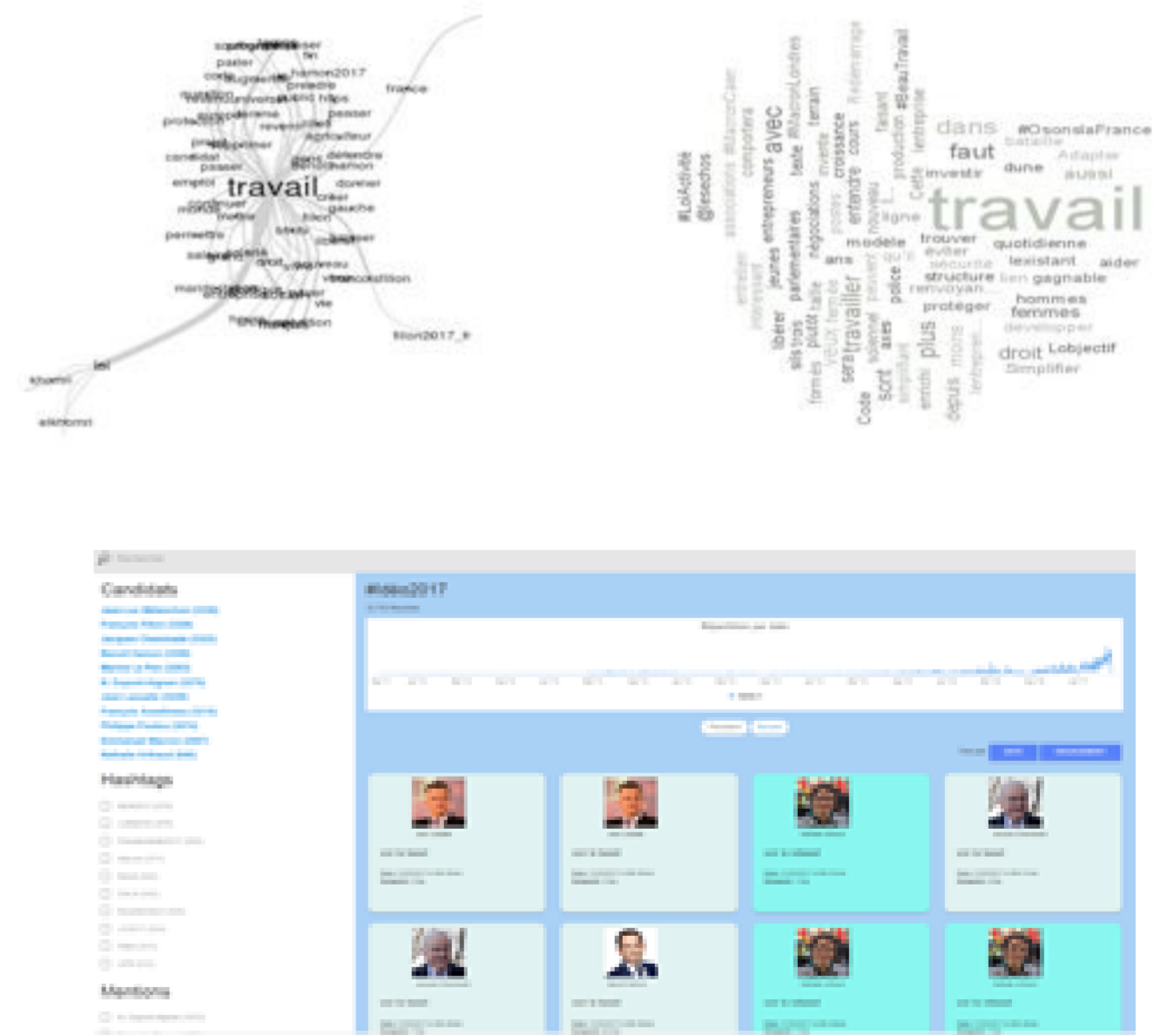
Analysis

The #Ideo2017 platform proposes:  
 The Analysis "I Analyze the Tweets that Contain the Word [Word]"  
 The Analysis "I Analyze the Tweets of [Candidate]"  
 The Search Engine

Forme	Fréquence	Type
présidentielle2017	88	Non reconnue
jeanlasselalle	74	Non reconnue
lassalle	37	Non reconnue
jean	34	Nom
france	34	Non reconnue
retrouver	30	Verbe
grand	29	Adjectif
lele	26	Non reconnue
candidat	25	Nom
français	23	Adjectif
politique	21	Adjectif



The screenshots show the user interface of the #Idéo2017 platform. On the left, there are filters for 'J'analyse les tweets qui contiennent le mot:' with options like France, Démocratie, État, République, Sécurité, Immigration, etc. The middle section shows 'J'analyse les tweets de:' with a list of candidates such as Nathalie Arthaud, François Asselineau, Jacques Cheminade, etc. The right section shows 'Je navigue dans tous les tweets par filtre:' with a search bar and a 'Naviguez' button. Below these are visualizations of tweet data.



Conclusions

- ❖ #Idéo2017 combines different technologies and inputs, which give citizens the opportunity to grasp a part of the discursive issues of the election.
- ❖ This development, which can be enriched, allows citizen to easily use a set of features usually accessible by software requiring different transformations of the data.
- ❖ For the period from September 1st 2016 to May 7th 2017, 42290 tweets were extracted for the 11 candidates. These tweets were gathered in a collection published in a TEI corpus in the standards of the Ortolang platform, thanks to the Corli Consortium.

CoMeRe Repository (2014). Repository for the CoMeRe corpora [website]. <http://hdl.handle.net/11403/comere>  
 Burnard, L. & Bauman, S. (2013). TEI P5: Guidelines for electronic text encoding and interchange. TEI consortium, *tei-c.org*. <http://www.tei-c.org/release/doc/tei-p5-doc/en/Guidelines.pdf>  
 Chanier, T., Poudat, C., Sagot, B., Antoniadis, G., Wigham, C.R., Hriba, L., Longhi, J. & Seddah, D. (2014). The CoMeRe corpus for French: structuring and annotation heterogeneous CMC genres, in BeiBwenger, M., Oostdijk, N., Storrer, A & van den Heuvel, H. Building and Annotating Corpora of Computer-Mediated Discourse: Issues and Challenges at the Interface of Corpus and Computational Linguistics, *Journal of Language Technology and Computational Linguistics* (special issue), pp1-31. [http://www.jlcl.org/2014\\_Hef2/Hef2-2014.pdf](http://www.jlcl.org/2014_Hef2/Hef2-2014.pdf)  
 DCMI (2014). Dublin Core Metadata Initiative. <http://dublincore.org/>  
 Laboratoire AGORA (AGORA), Equipes Traitement de l'Information et Systèmes (ETIS) (2017). *Présidentielle2017: corpus des tweets de la présidentielle2017* [Corpus]. ORTOLANG (Open Resources and TOOLS for LANGuage) - [www.ortolang.fr](http://www.ortolang.fr). <https://hdl.handle.net/11403/corpus-presidentielle2017/v1>  
 Longhi J., Marinica C., Borzic B., Alkhouli A., 2014: Polititweets, corpus de tweets provenant de comptes politiques influents. In Chanier T. (ed) Banque de corpus CoMeRe. Ortolang.fr : Nancy. [cmr-polititweets- tei-v1]: <https://repository.ortolang.fr/api/content/comere/v3.3/cmr-polititweets.html>  
 Longhi, J., Wigham, C. (2015). Structuring a CMC corpus of political tweets in TEI: corpus features, ethics and workflow. Poster presented in Corpus Linguistics 2015, Jul 2015, Lancaster, United Kingdom, available in: <https://halshs.archives-ouvertes.fr/halshs-01176061>.  
 Longhi J. (dir.), 2017 : #Idéo2017, plateforme d'analyse de tweets politiques en campagne électorale. OLAC. (2008). Best Practice Recommendations for Language Resource Description. *Open Language Archives Community*. University of Pennsylvania. <http://www.languagearchives.org/REC/bpr.html>  
 ORTOLANG (2013). Open Resources and TOOLS for LANGuage [website]. ATILF / CNRS - Université de Lorraine: Nancy. <http://www.ortolang.fr>