

# Parliamentary corpora

Treasure trove for the Digital Humanities  
Challenge for infrastructure development



Christian Mair (University of Freiburg & CLARIN-D F2)

## 1. What makes parliamentary data attractive for research?

- natural database for **interdisciplinary research**, bringing together (**historical**) **linguistics**, **history**, **political science** and **cultural studies** under the **Digital Humanities** umbrella
- **self-replenishing** source of data with potentially considerable **historical time-depth**
- potential for **comparative research** across national parliamentary traditions
- potential (as yet largely unexplored) for research on **multimodality**
- significant amounts of **multilingual** and **translation data** (e.g. Belgium, European Parliament)

## 2. Case studies from the Hansard Corpus (CLARIN-UK): "Distant Reading" (Moretti 2013) of 1.6 billion words, 1803-2005

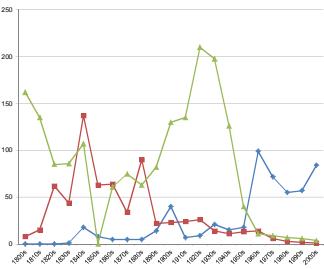


Fig. 1: Immigration, emigration, Empire (1803-2005, n/pmw)

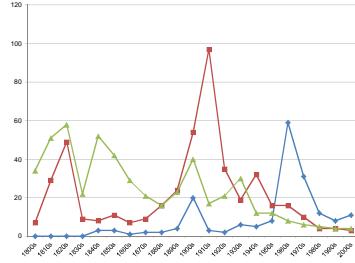


Fig. 2: Immigrants, aliens, foreigners (1803-2005, n/pmw)

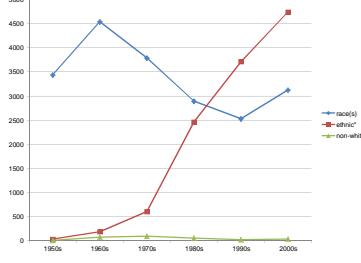
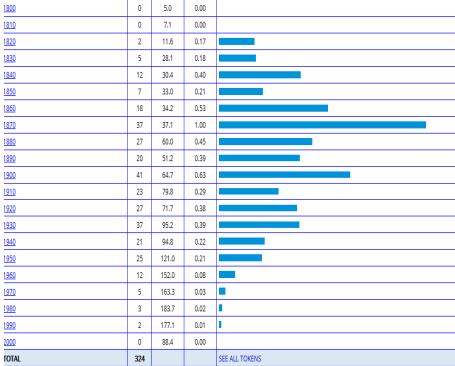
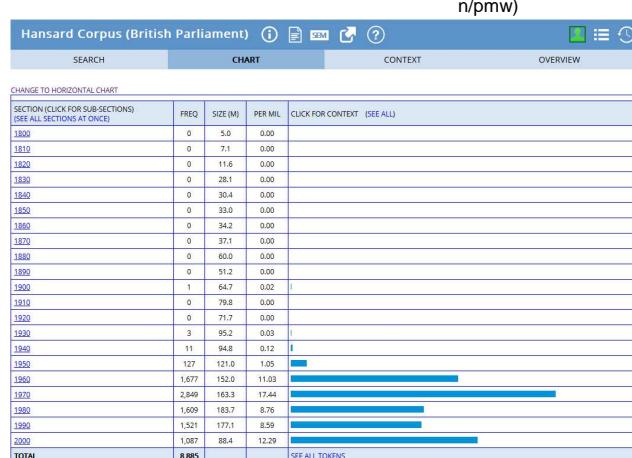


Fig. 3: From 'race' to 'ethnicity' (1951-2005, n/pmw)



↑Fig. 4: dominant|inferior|superior|backward|advanced RACE\_n\* (1803-2005)



→Fig. 5: race relations|discrimination|equality (1900-2005)

## 3. Challenges for infrastructure development in the European CLARIN context

- existing parliamentary corpora (cf. Fišer & Lenardič 2017) extremely variable in size and historical time depth: e.g. Dutchparl 800m, 1814-2014; PTPARL 1m, 1970-2008
- enormous differences in type and depth of annotation

### Perspectives for the future:

Given CLARIN's helpful but generally modest role in the various national projects, comprehensive standardisation of data and tools on CLARIN terms is unrealistic. What will work is a **one-stop easy-mode web interface**, hosted by CLARIN and providing access to European parliamentary corpora on the following terms:

- easy-to-use full-text search (with POS tagging, collocational statistics, and named entity recognition) for all corpora
- standardisation across corpora and tools, to attract large numbers of novice users across the disciplines
- long-term and systematic support to help regular users develop into a visible international research community.

## 4. Traditional "close reading"

1848: Other plans had been tried: European colonists had been sent out, but **European immigration** had, to a great extent, failed: **Coolie immigration** had been tried, and that had likewise to a considerable degree failed [...] (context: post-Abolition West Indies)

1905: [...] they must bear in mind that some of the undoubted evils which had fallen upon portions of the country from an **alien immigration** which was largely **Jewish**, gave those of them who, like the right hon: Baronet and himself, condemned nothing more strongly than the manifestation of the **anti-Semitic** spirit, some reason to fear that this country might be, at however great a distance, in danger of following the evil example set by some other countries [...]

1967: The hon: Member for Barnsley (Mr. Mason), the Minister of Defence for Equipment, who is grinning at this, should go back to Barnsley and tell his voters that his policy is to have more **black people** in this country than **white** and see what they do to him: The Home Office says that by 1985 there will be 3½ million **coloured people** in this country. It should know more about it than the hon: Gentleman does: By doing nothing, the Government are betraying this country and compelling the English [to] **commit race suicide**: [...] This is true; it is no good hon: Members tut-tutting: For all these reasons I find the Address bitterly disappointing and the Government a cowardly failure (Commons, Osborne, Labour)

### References

- Alexander, Marc, & Andrew Struan. 2017. Digital Hansard: The politics of the uncivil. In *Digital Humanities 2017*, Montréal, Conference Abstracts. [https://dh2017.adho.org/abstracts/DH2017\\_abstracts.pdf](https://dh2017.adho.org/abstracts/DH2017_abstracts.pdf).
- Baker, Helen, Vaclav Brezina, and Tony McEnery. 2017. Ireland in British parliamentary debates, 1803-2005. In Tanja Säily, Arija Nummi, Minna Palander-Collin and Anita Auer, eds. *Exploring future paths for historical sociolinguistics*. Amsterdam: Benjamins. 83-108.
- Fišer, Darja, and Jakob Lenardič. 2017. Parliamentary corpora in the CLARIN infrastructure. *Selected papers from the CLARIN Annual Conference 2017*, Budapest, 18 - 20 September 2017. Conference Proceedings published by Linköping University Electronic Press at [www.ep.liu.se/ecp/content.asp?issue=147](http://www.ep.liu.se/ecp/content.asp?issue=147).
- Gabrielatos, Costas, Tony McEnery, Peter Diggle, and Paul Baker. 2012. The peaks and troughs of corpus-based contextual analysis. *International Journal of Corpus Linguistics* 17: 151-175.
- Mair, Christian. 2019. Brexitiness: The Ebbs and Flows of British Eurosceptic Rhetoric since 1945. *Open Library of Humanities*, 5(1): 50, 1-26. DOI: <https://doi.org/10.16995/oh.424>.
- Molin, Sandra. 2007. The Hansard hazard: Gauging the accuracy of British parliamentary transcripts. *Corpora* 2: 187-210.
- Moretti, Franco. 2013. *Distant reading*. London: Verso.