

## **Main, high level, research question:**

*How has both discourse and representation in the British parliament changed between 1803 and 2005 from the perspective of gender?*

In this report, we do basic data analysis of the Hansard corpus. We first crawl the data (201 years worth, from 1803-2005, excepting 1816 and 1829 for which we did not find data), and parse it to create three different views of the metadata

1. Parquet files - contain debate and speaker data, separately for each year
2. Jsonline - contain all data parsed and arranged into jsonline per year, good for manual inspection
3. SQLite DB - untested so far - meant to enable keyword search

We first present some statistics around our data.

### **Core Statistics:**

- Years: 201 years (1803-2005)
- Debates: 802,178 total debates
- Average Debate Length: 1,487 words
- Unique Speakers: 89,472 speakers
- Total Corpus: 1.19 billion words (~6.7 GB)
- Coverage: 99% of possible years (only missing 1816 & 1829)
- Chambers: 79.7% Commons, 20.3% Lords
- Speaker Activity: 29.9 average mentions per speaker
- Debates/Year: ~4,000 average

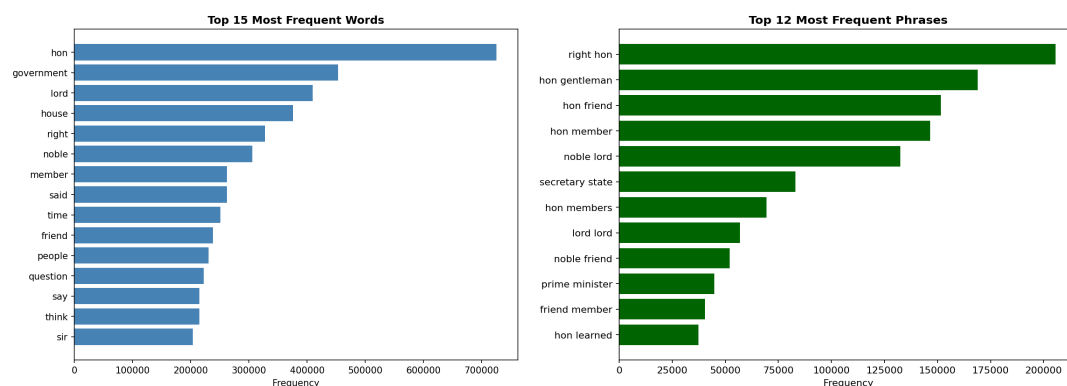
Overall, we see that this is a rich dataset with lots of valuable text data and associated metadata. Below, we attempt to analyze this data from the perspective of gender. We try to give a sense of overall distribution of the kind of data that we have collected, and then we show gender-based analysis.

### **Methodology:**

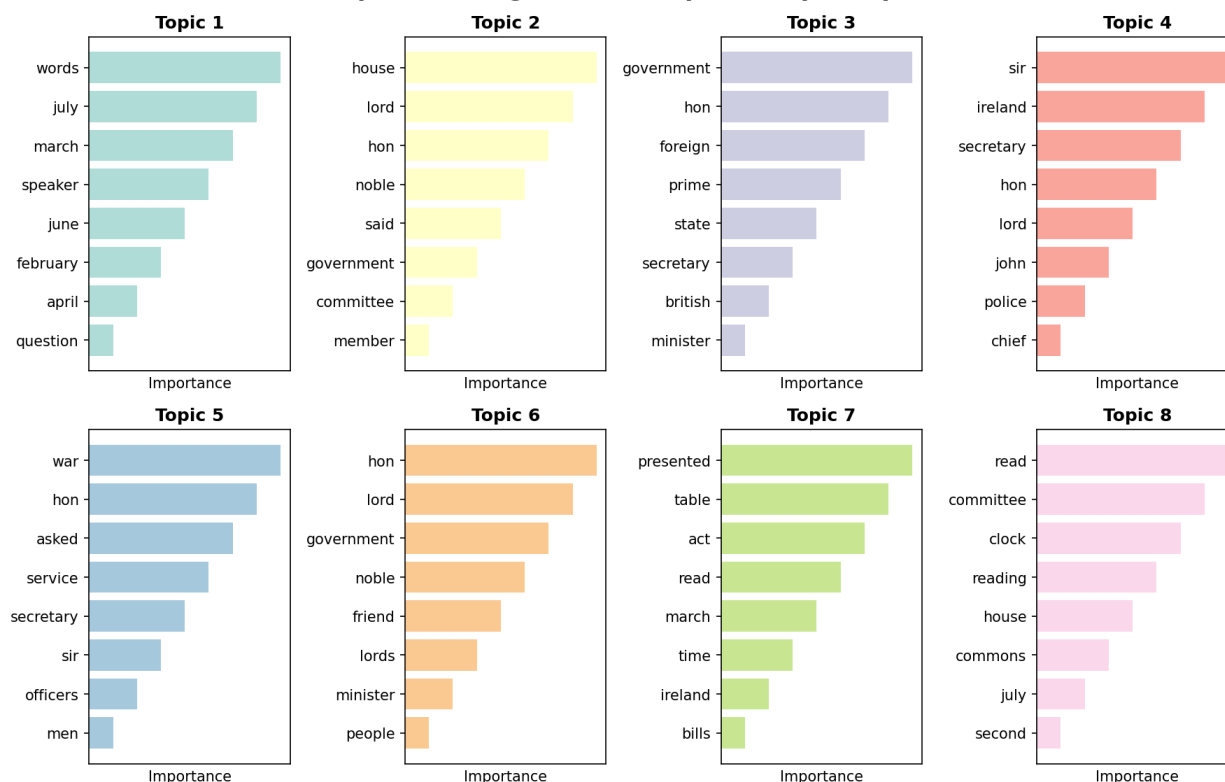
For overall analysis, we sample 100k debates, using stratified sampling to give us a good representation for each year we have availability for, upon which we do simple unigram and bigram analysis, which show, as one would expect, a lot of procedural language which is very common in parliamentary speech. We also run topic modeling (using LDA) over the sample of 100k debates. In addition, we have speaker metadata, for which we use title-based heuristics (ex. Lord - Male, Duchess - Female) to identify gender.

### **Overall Corpus Analysis:**

We see the words “lord”, “government”, “hon” - a contraction of the word honorable - are most common, which seems to follow the pattern that one would expect in parliament, since these are forms of addressing fellow members. Our bigram analysis also shows that formal forms of address are extremely common, again, due to the rule-based nature of parliament, where members address each other in a fairly standardized manner.

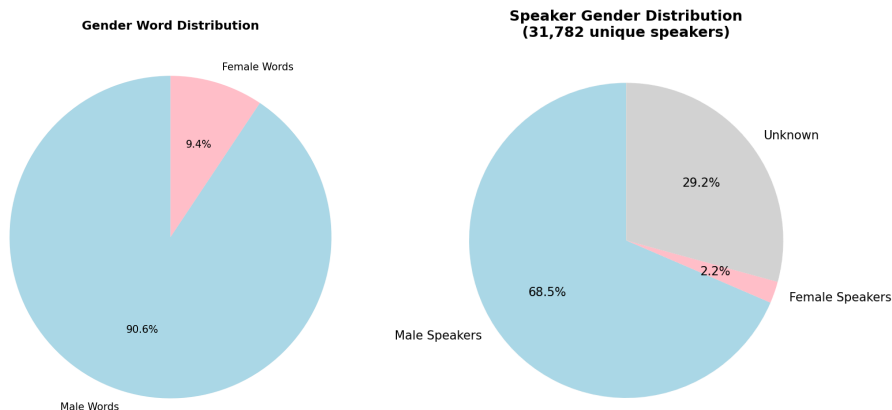
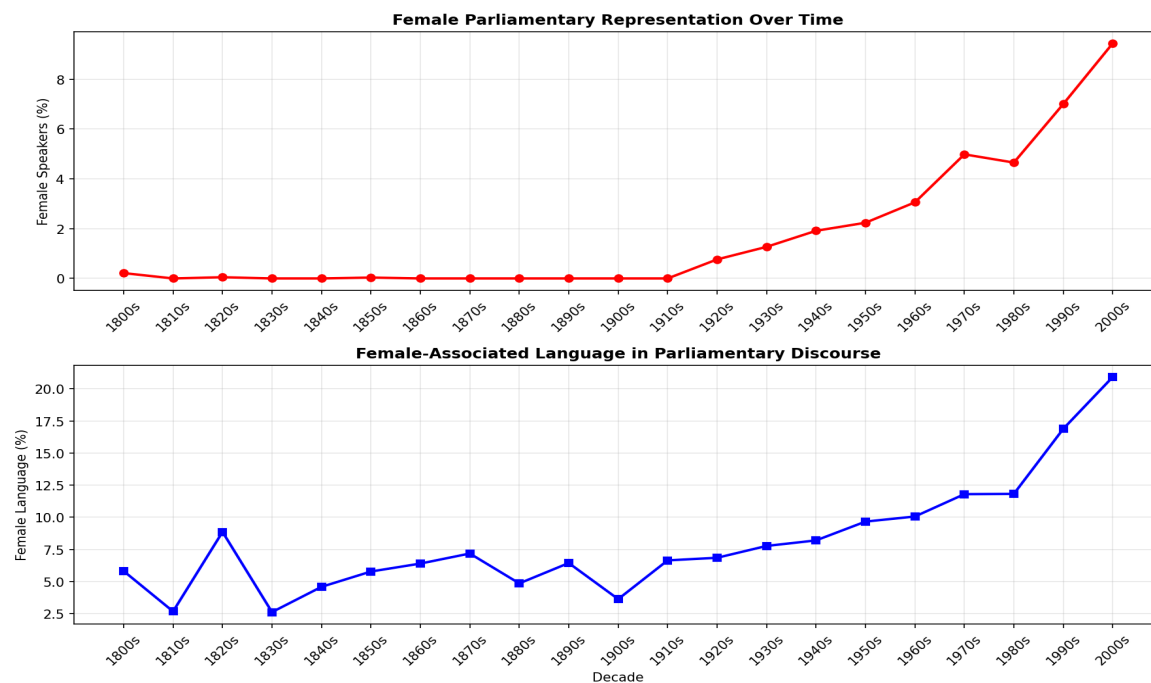


### Topic Modeling Results - Top Words per Topic



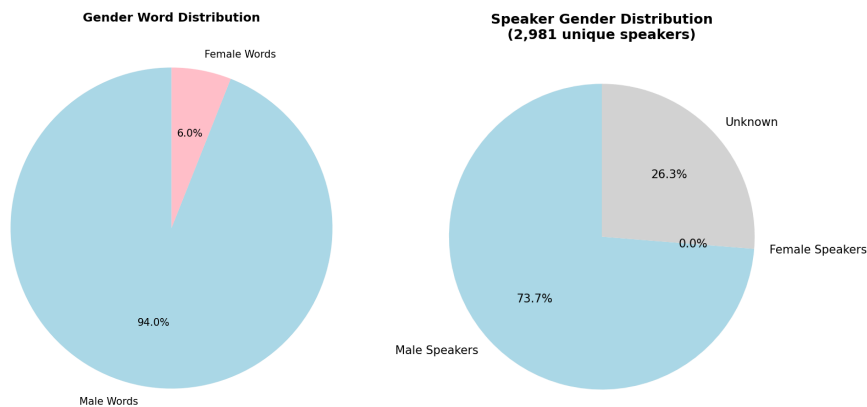
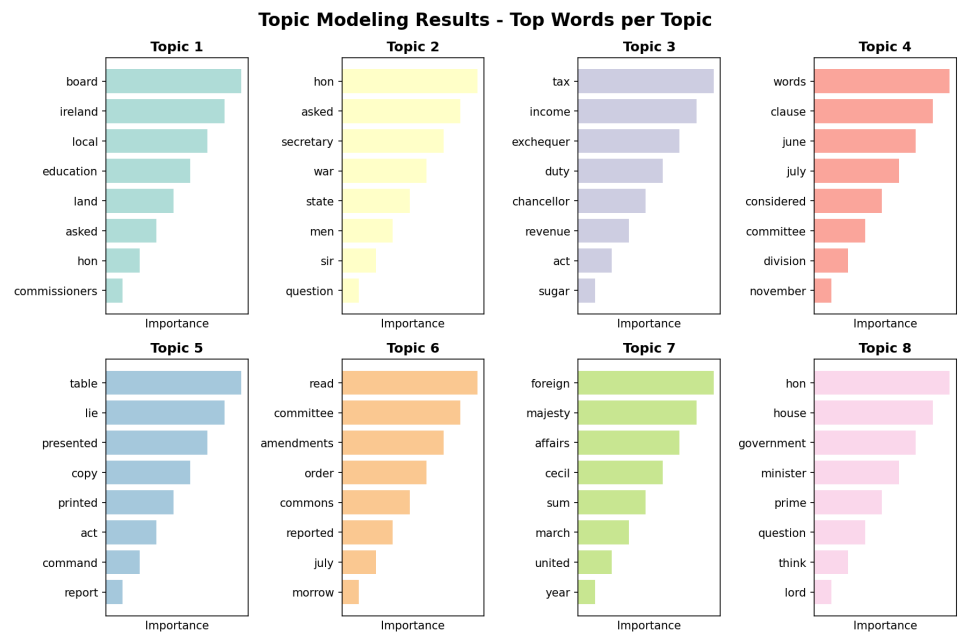
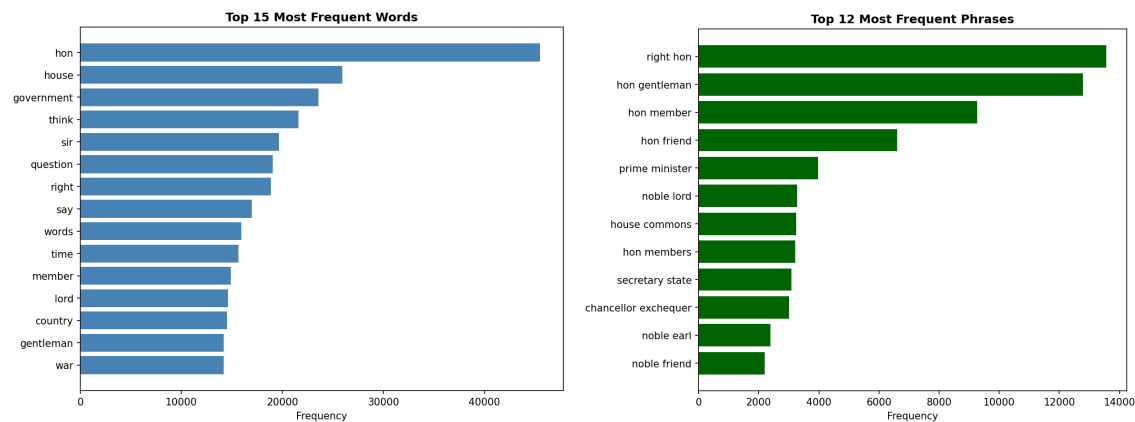
We see that there's a steady increase in representation (starting in 1918), with an interesting drop in 1960 (check if this is noise or corresponds to other social events).

We also collect a set of gendered words (Zhao et. al. -2018) that we use to analyze the prevalence of gendered language over time. The female-associated language graph (the percentage is calculated over all gendered words in our word list) shows that there is an increased proportion of gendered words over time.

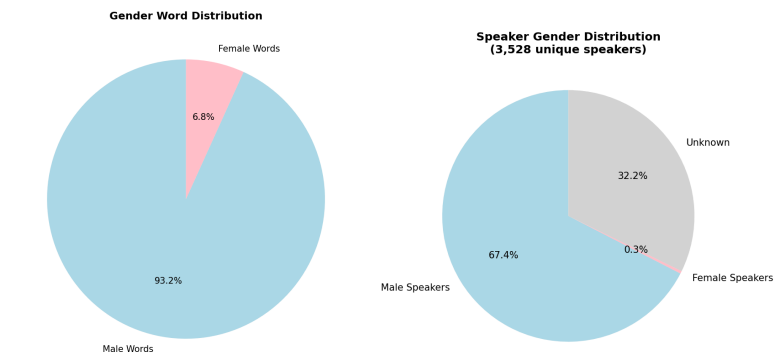
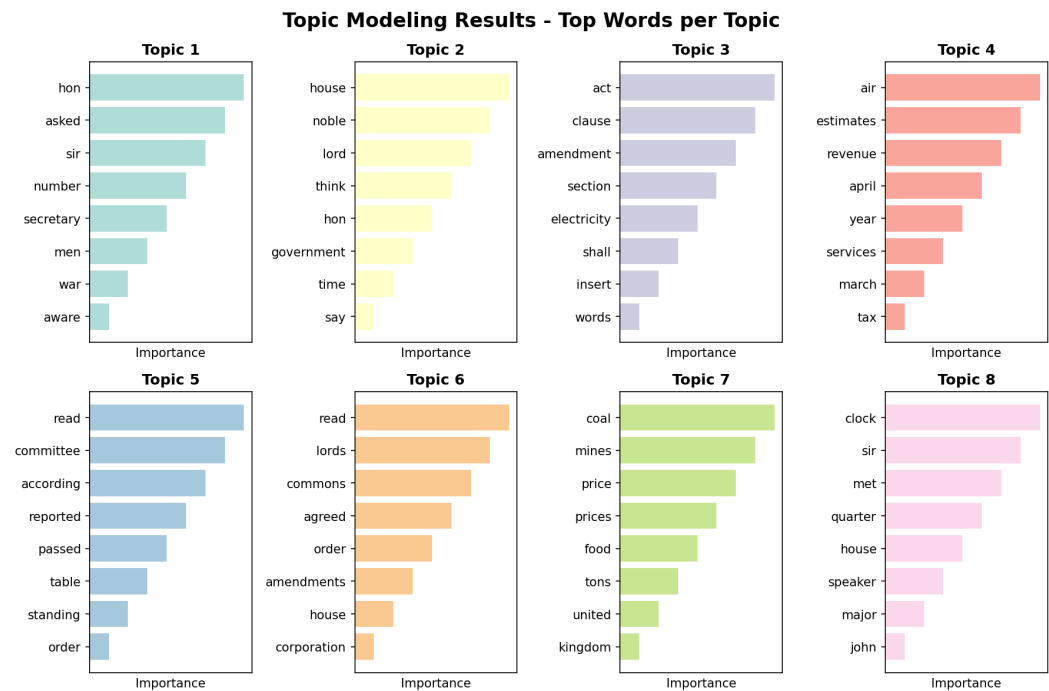
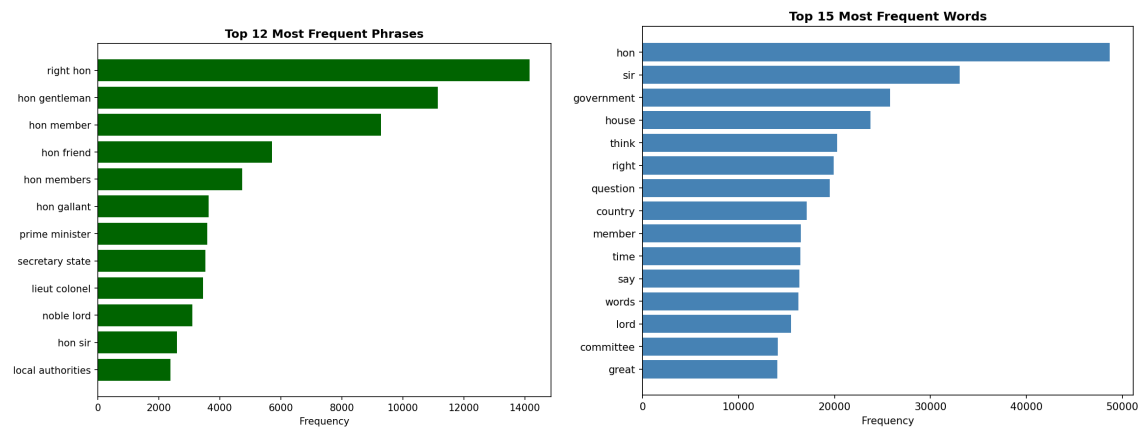


In addition, we do similar analyses for several key milestones in history. We look at 10 years before and after this time period, and in the case where the event lasted multiple years, we separate our analysis into pre/during/post the event under consideration.

1918 - Women in Britain get partial suffrage and women are allowed to be MPs (1918 Representation of the People Act).  
Pre-1918

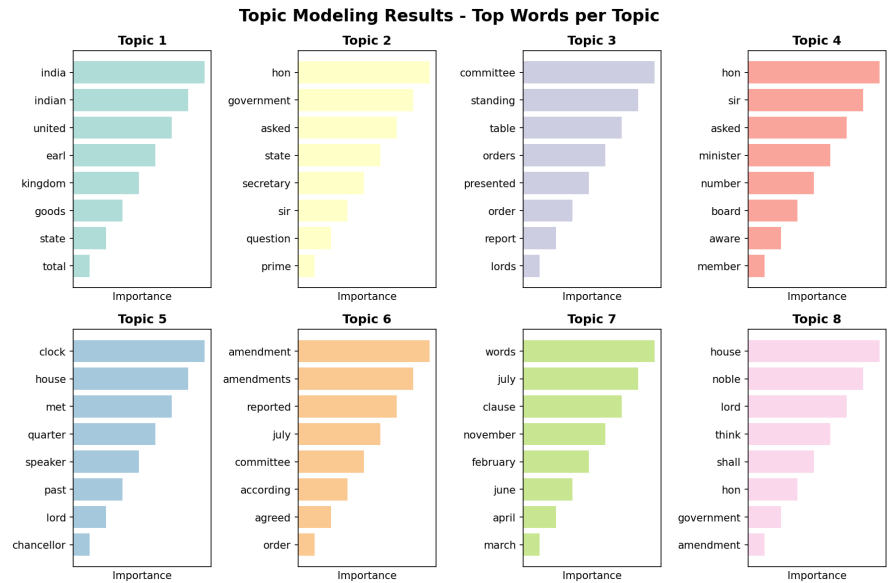
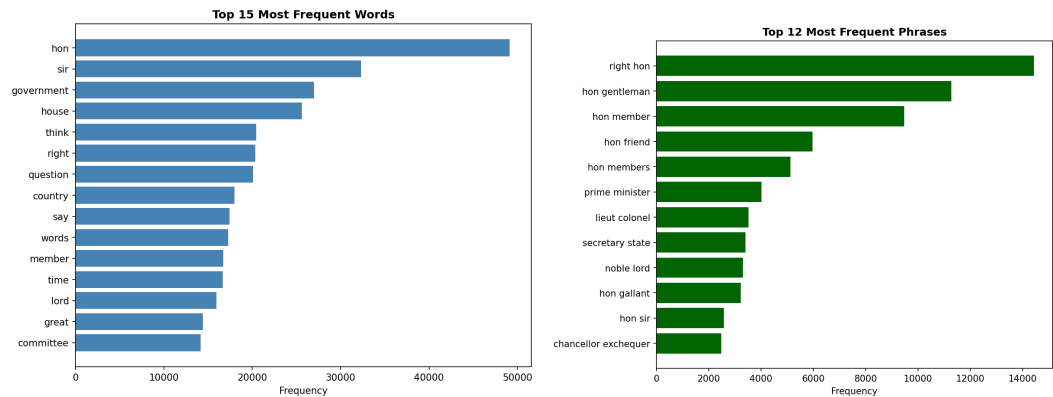


Post-1918

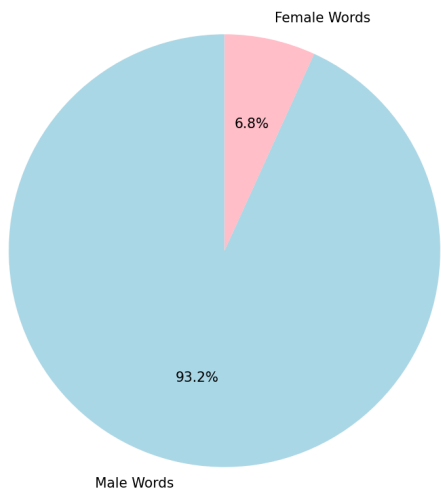


1928 - Women get universal suffrage (1928 Equal Franchise Act)

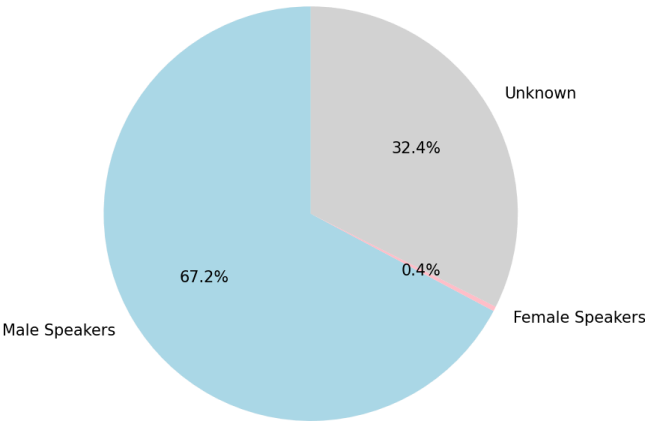
Pre-1928



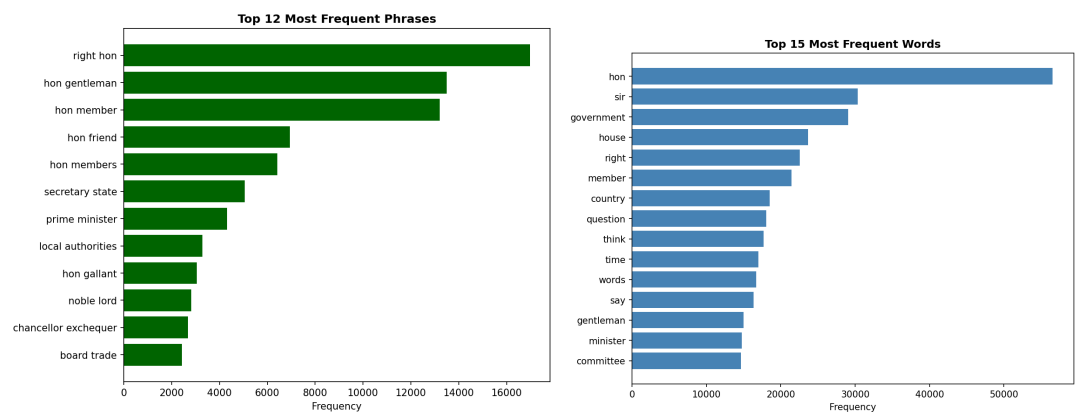
Gender Word Distribution



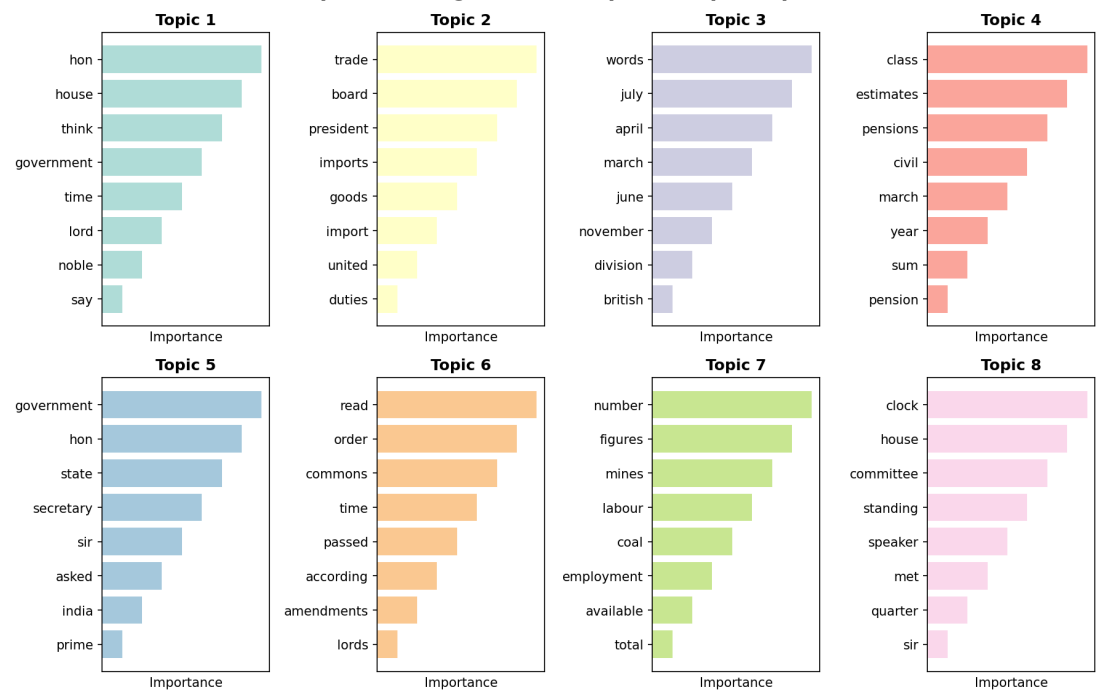
Speaker Gender Distribution (3,490 unique speakers)



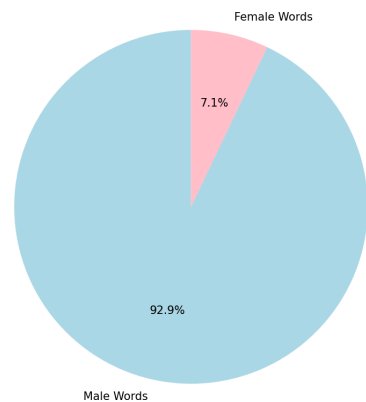
post-1928



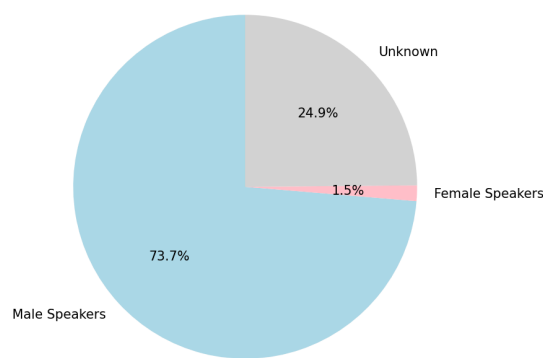
Topic Modeling Results - Top Words per Topic



Gender Word Distribution

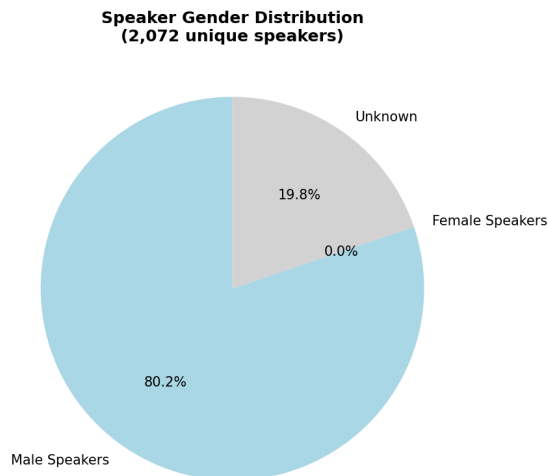
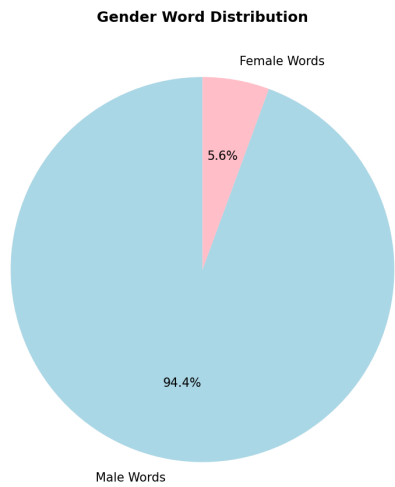
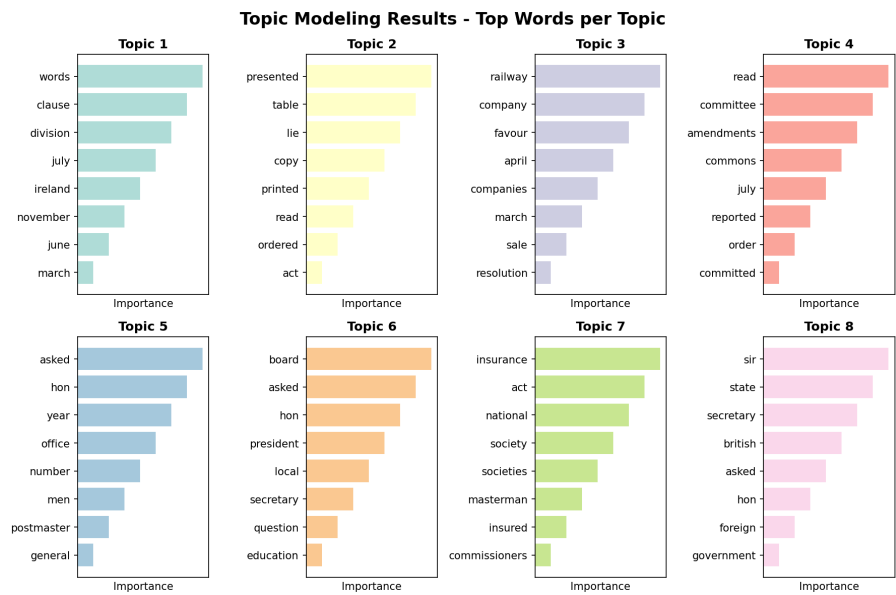
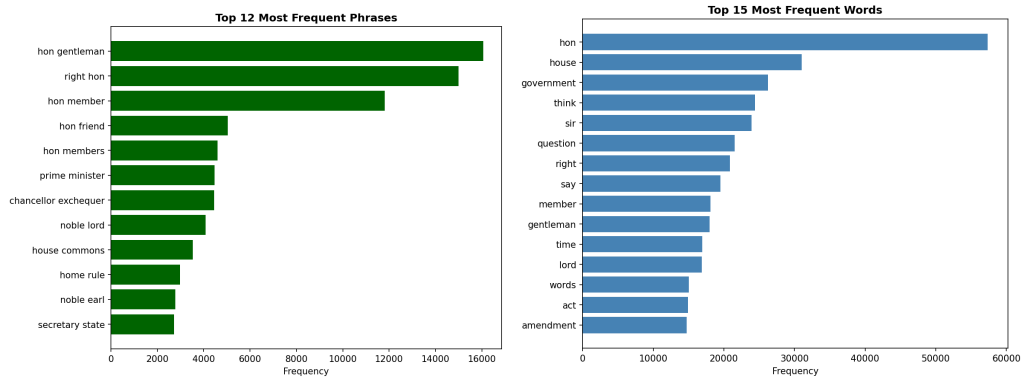


Speaker Gender Distribution (3,475 unique speakers)



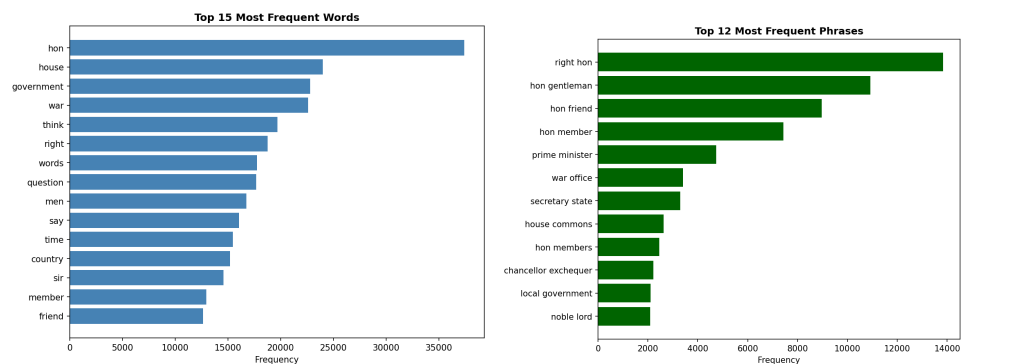
WW1 - 1914 - 1918

Pre-WW1

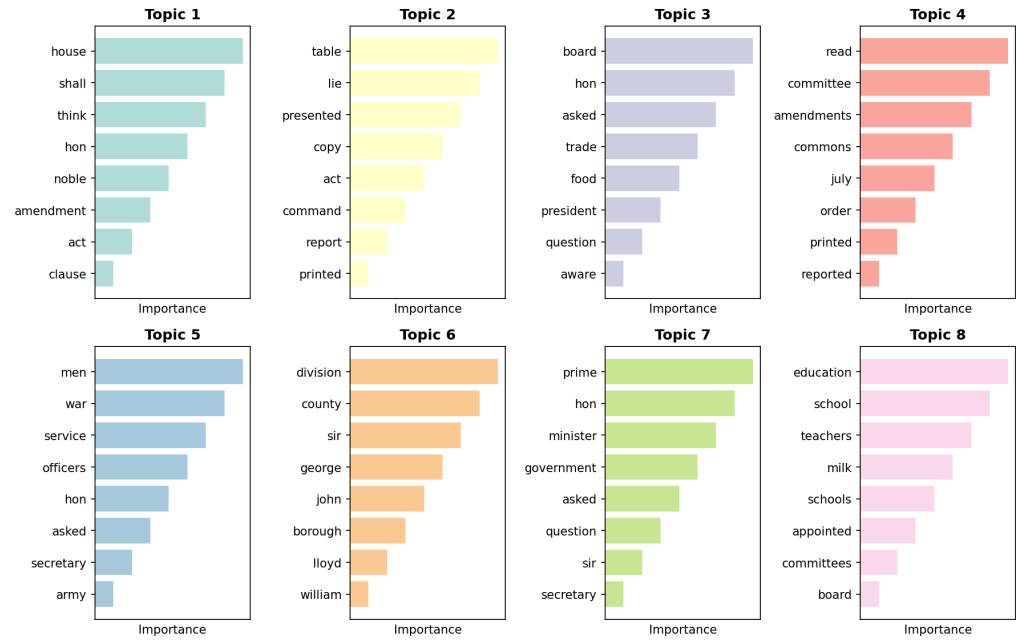




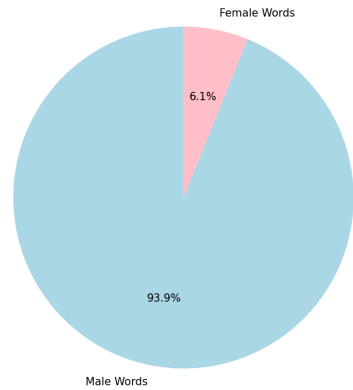
During-WW1



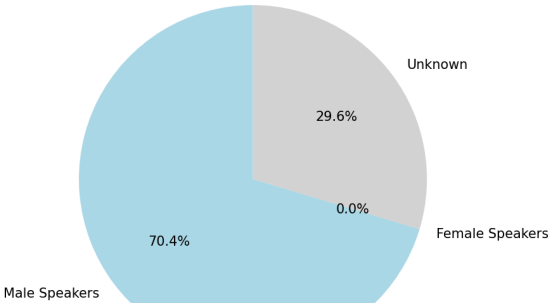
Topic Modeling Results - Top Words per Topic



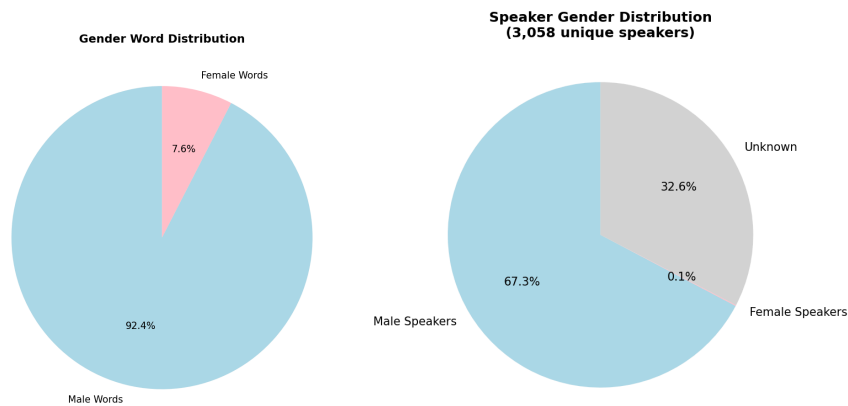
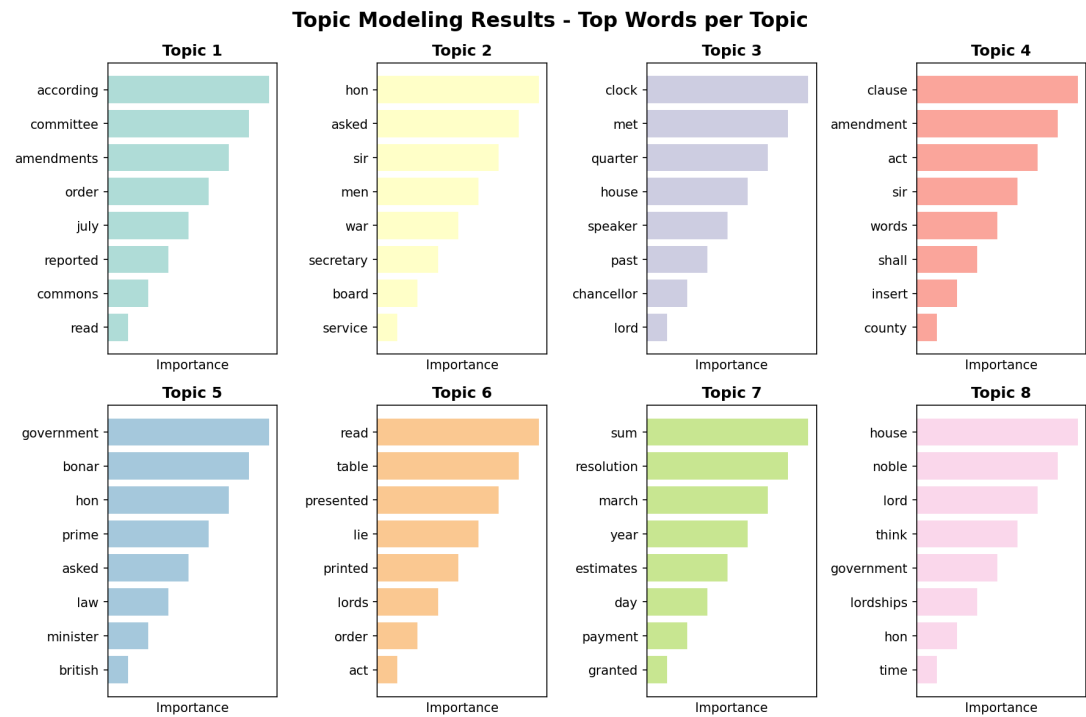
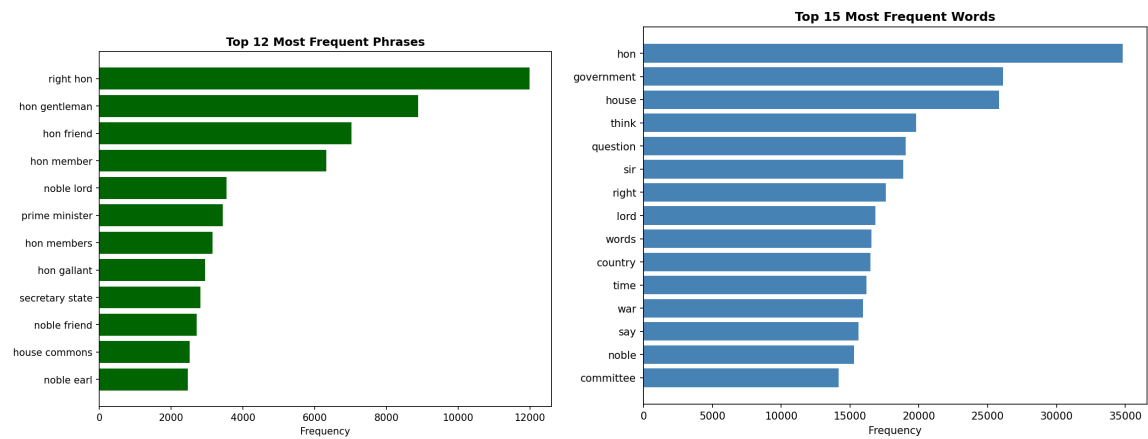
Gender Word Distribution



Speaker Gender Distribution (2,240 unique speakers)

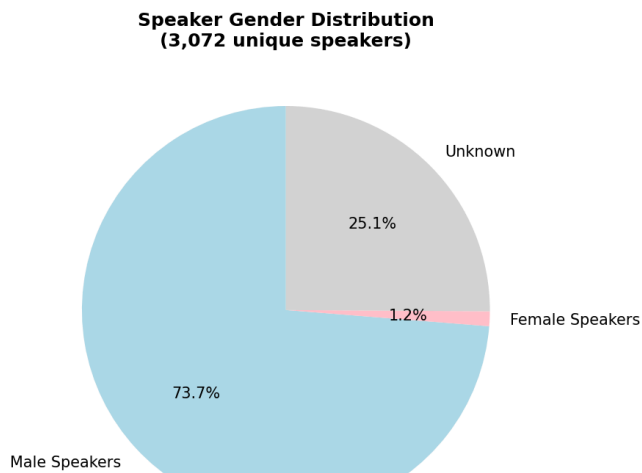
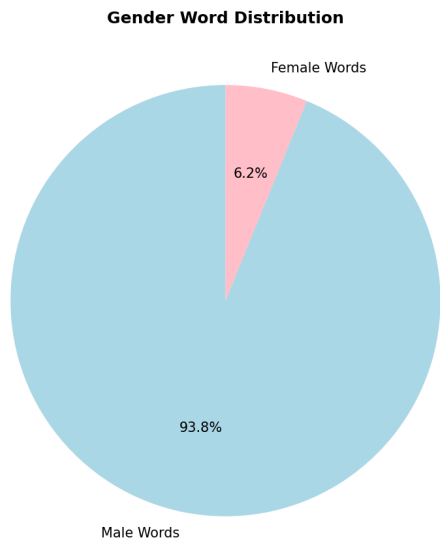
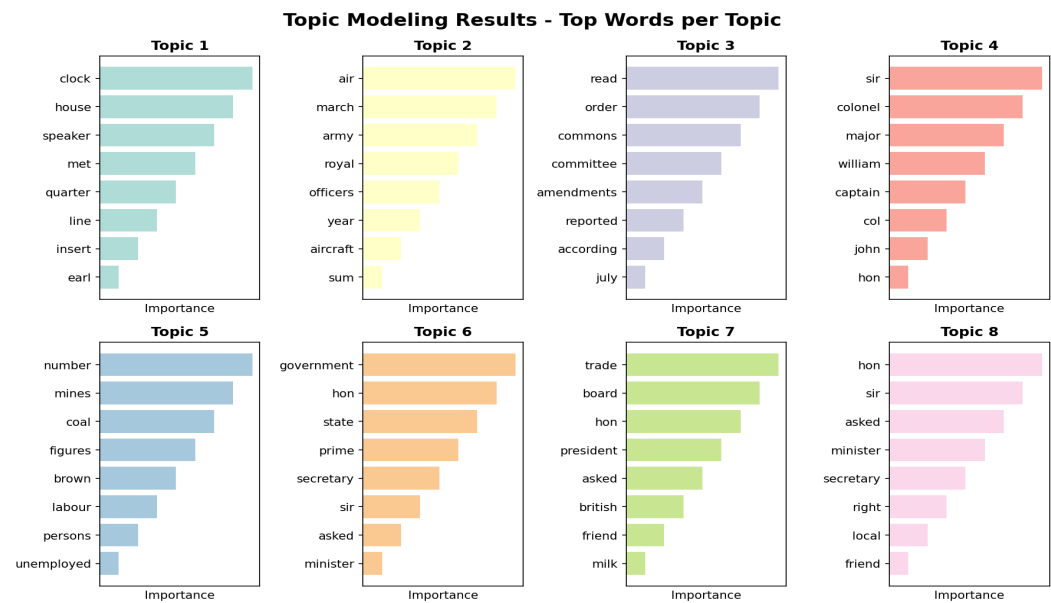
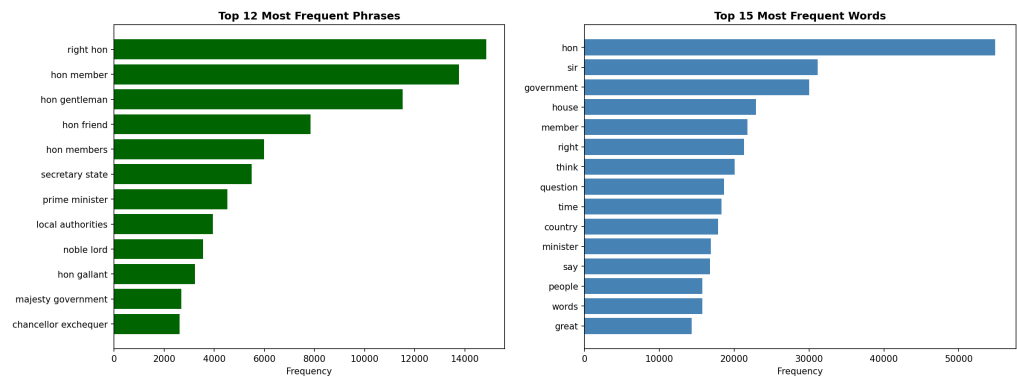


Post-WW1

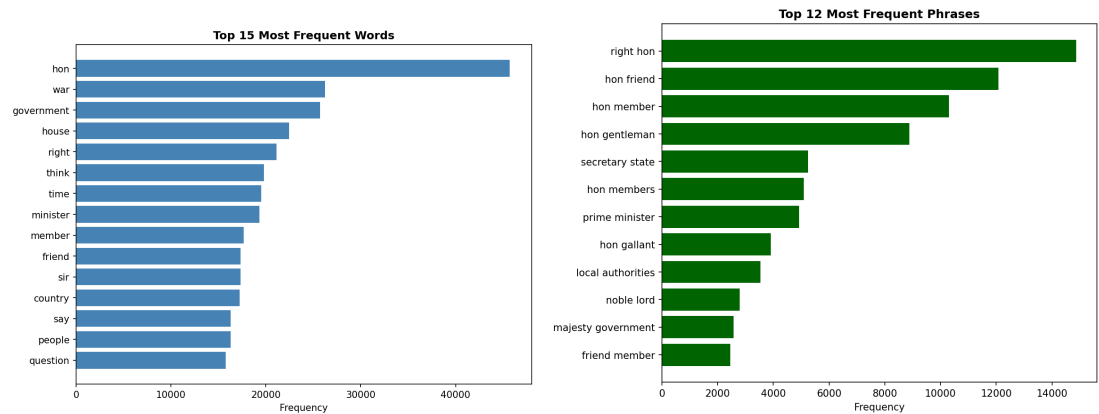


WW2 - 1939-1945

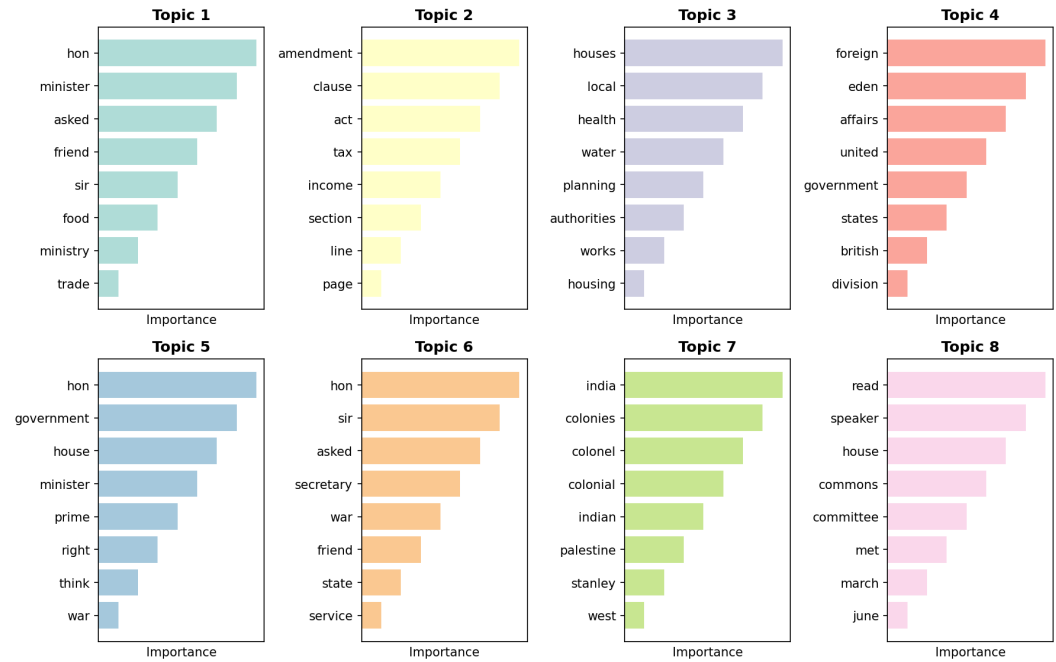
Pre-WW2



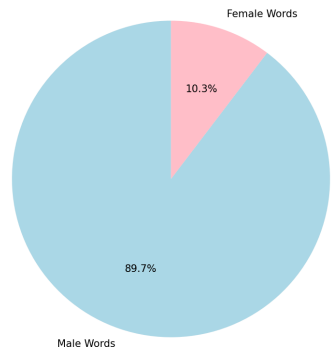
During-WW2



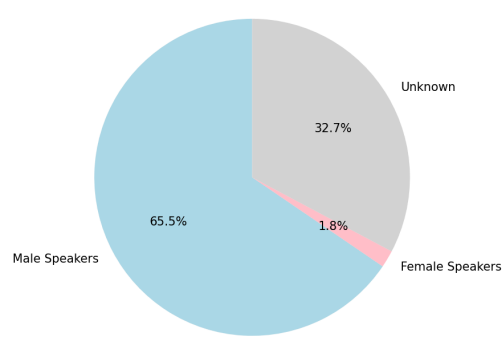
Topic Modeling Results - Top Words per Topic



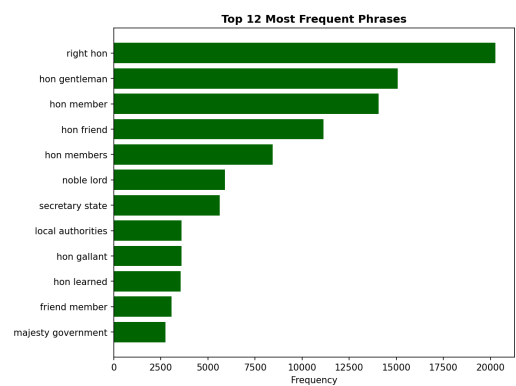
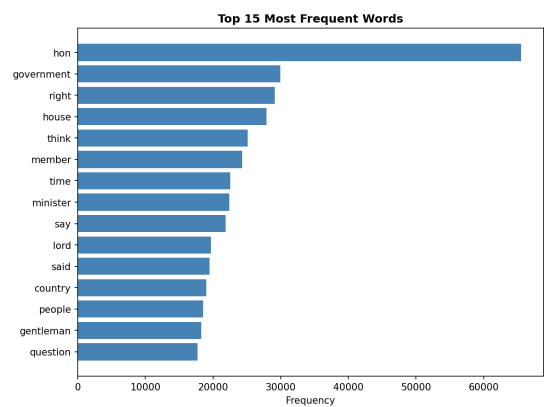
Gender Word Distribution



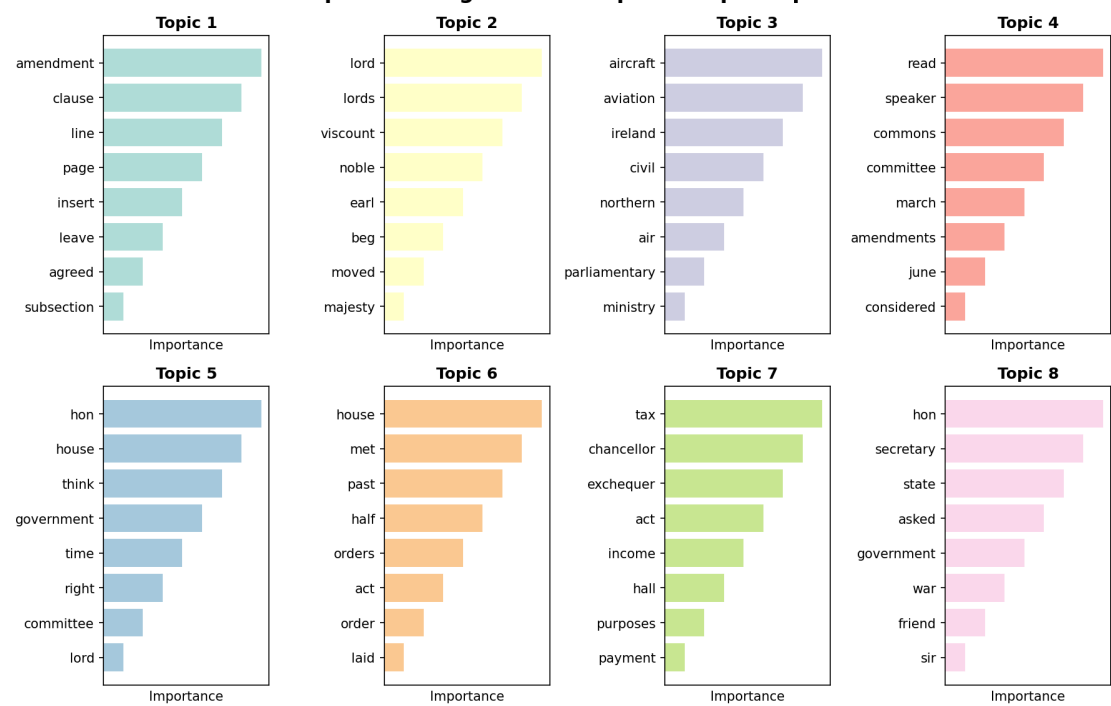
Speaker Gender Distribution (2,783 unique speakers)



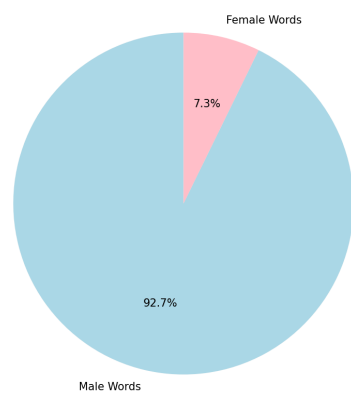
Post-WW2



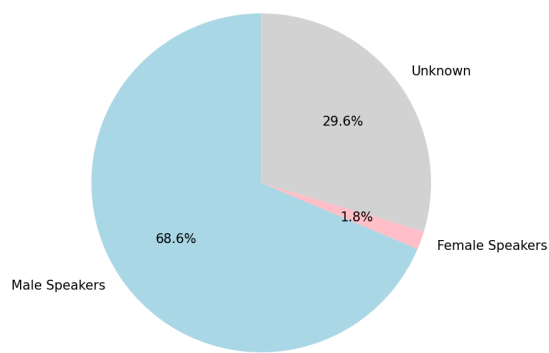
**Topic Modeling Results - Top Words per Topic**



**Gender Word Distribution**

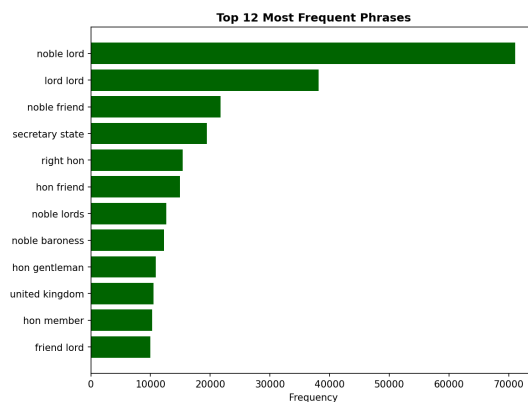
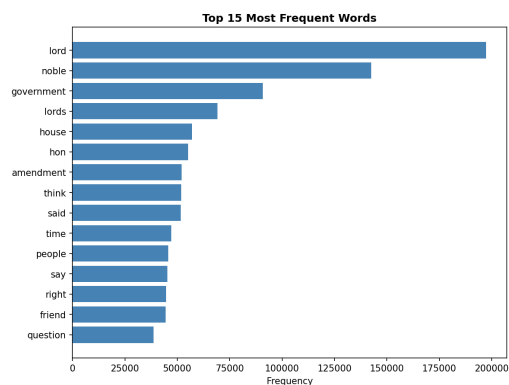


**Speaker Gender Distribution (2,887 unique speakers)**

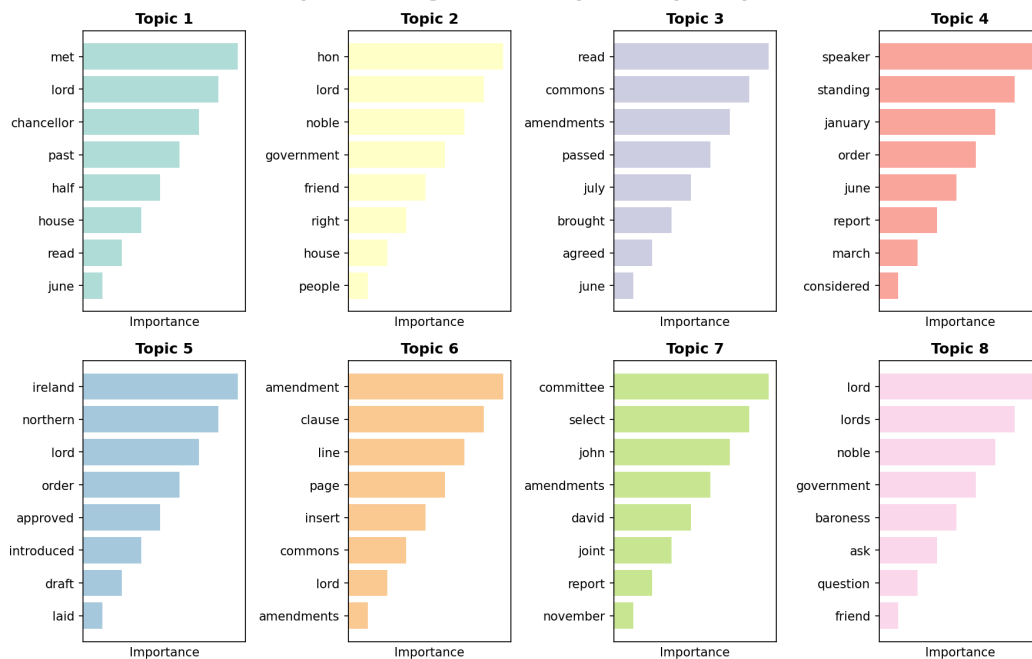


## 1979 - 1990 Margaret Thatcher is PM (first woman PM in the UK)

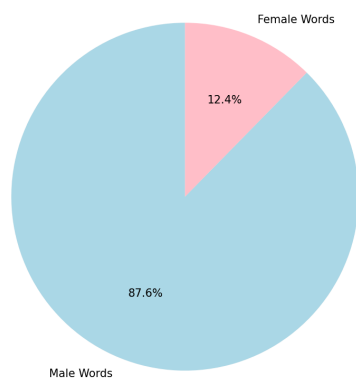
Pre-Thatcher:



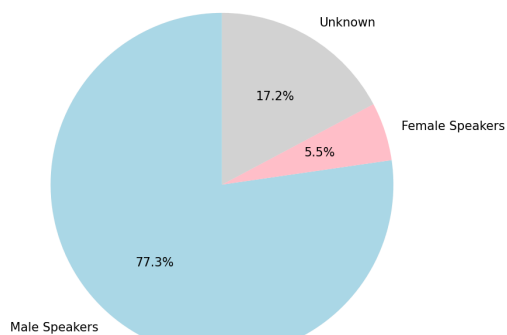
### Topic Modeling Results - Top Words per Topic



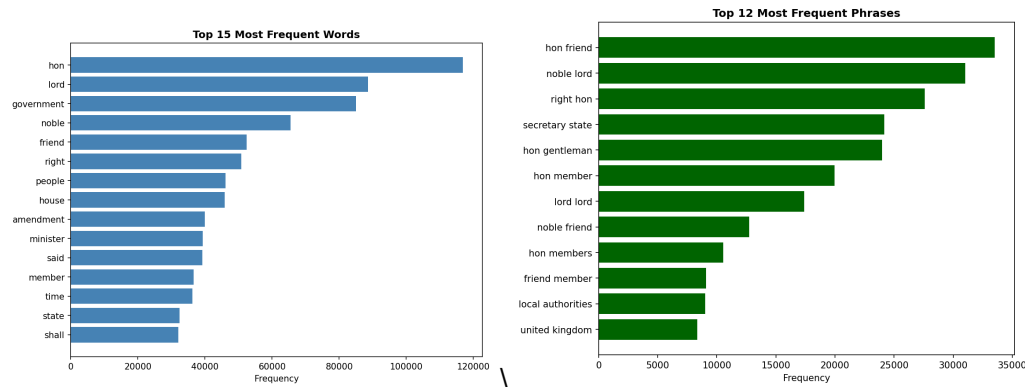
### Gender Word Distribution



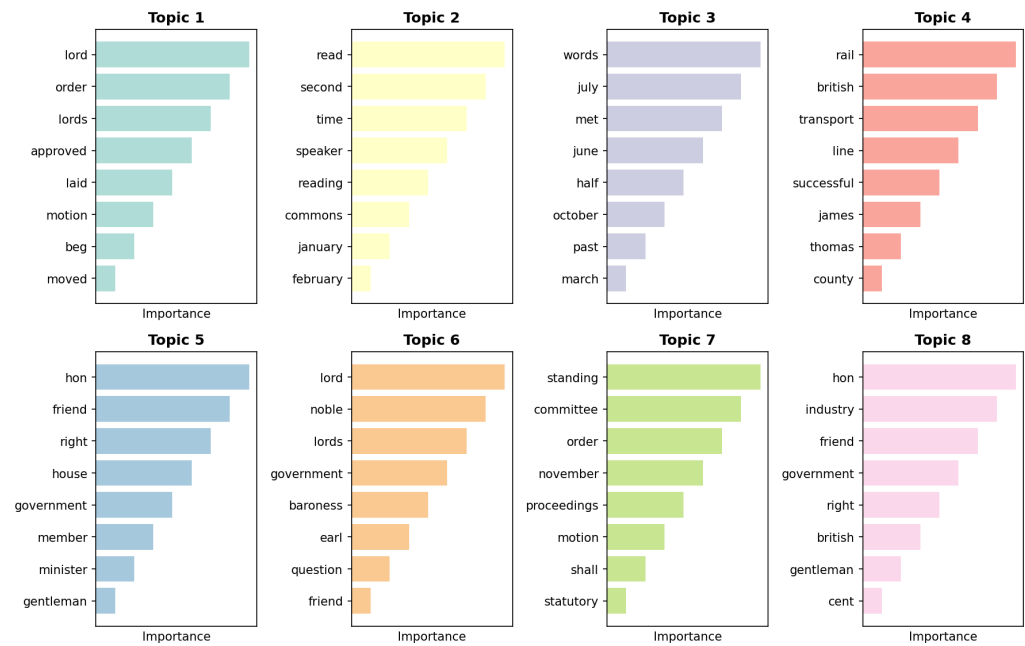
### Speaker Gender Distribution (3,591 unique speakers)



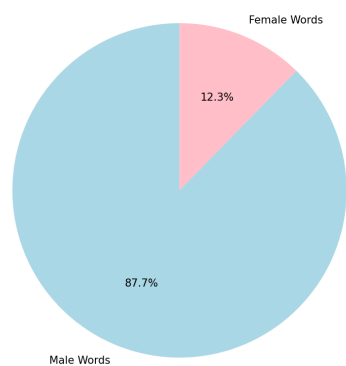
During-Thatcher



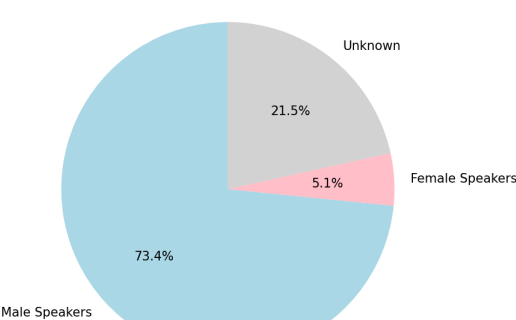
Topic Modeling Results - Top Words per Topic



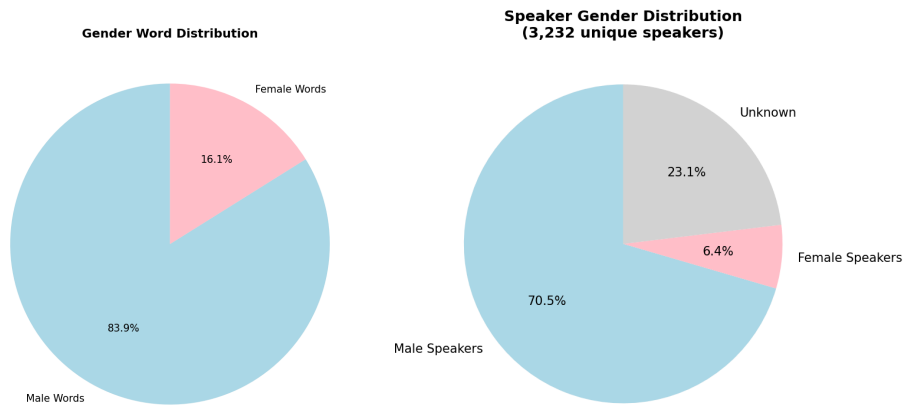
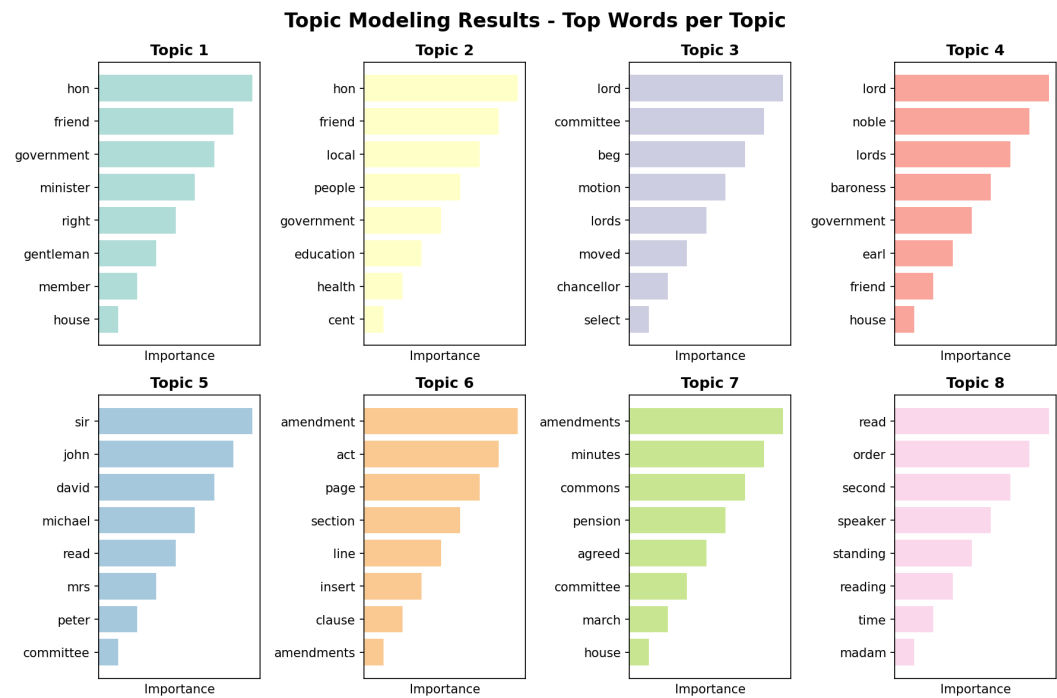
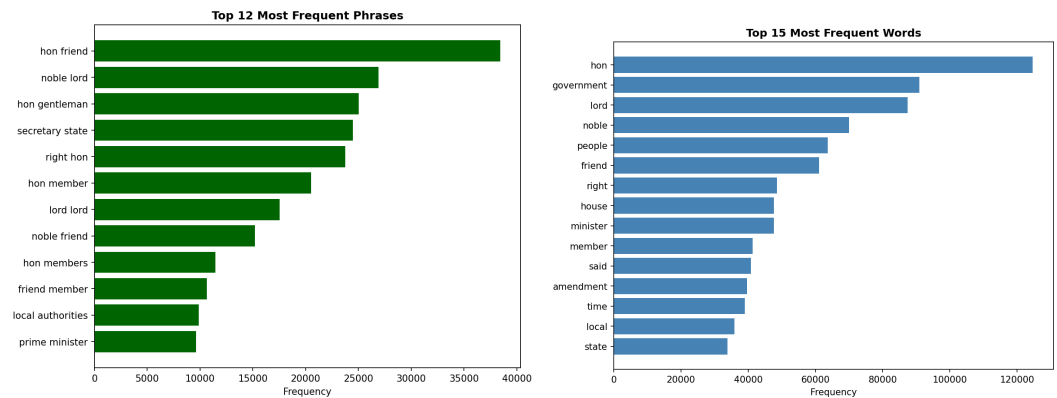
Gender Word Distribution



Speaker Gender Distribution (4,249 unique speakers)



Post-Thatcher





With the analysis we've completed for these different historical periods, we can see that there's a lot of commonalities in bigram/unigram analysis, as well as topics found. We seem to be mainly extracting parliamentary procedures and specific terminology, rather than the substance of the debates. Therefore, our next steps would be:

1. Develop a method to extract debate content versus parliamentary terminology.
2. Current heuristic based method for extracting MP gender is likely somewhat inaccurate, let's find a list of all members of houses in this period and cross-reference against other sources (such as wikipedia) to get a more accurate gender count
3. Develop a DSPy pipeline (potentially, still need to look into this) to use an LLM to extract and clean up debate text.
4. Improve speed of analysis so that we can run it locally