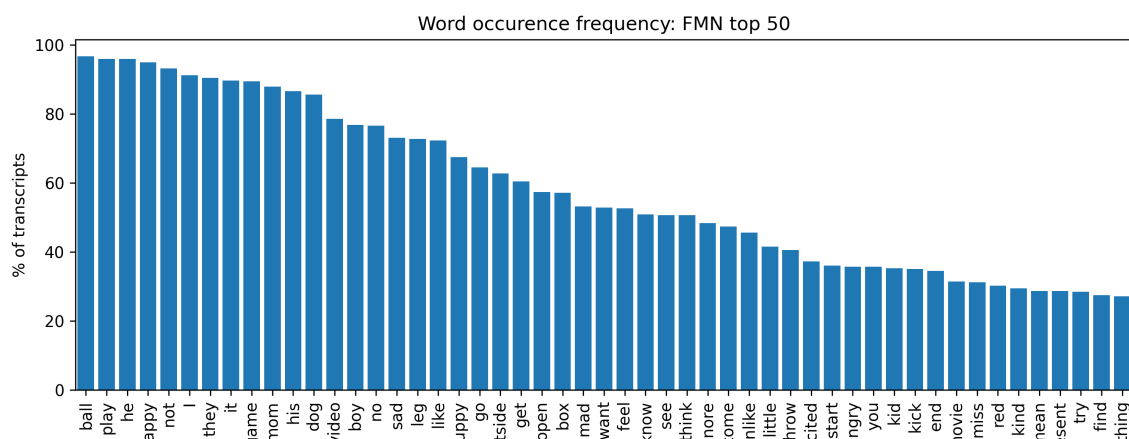


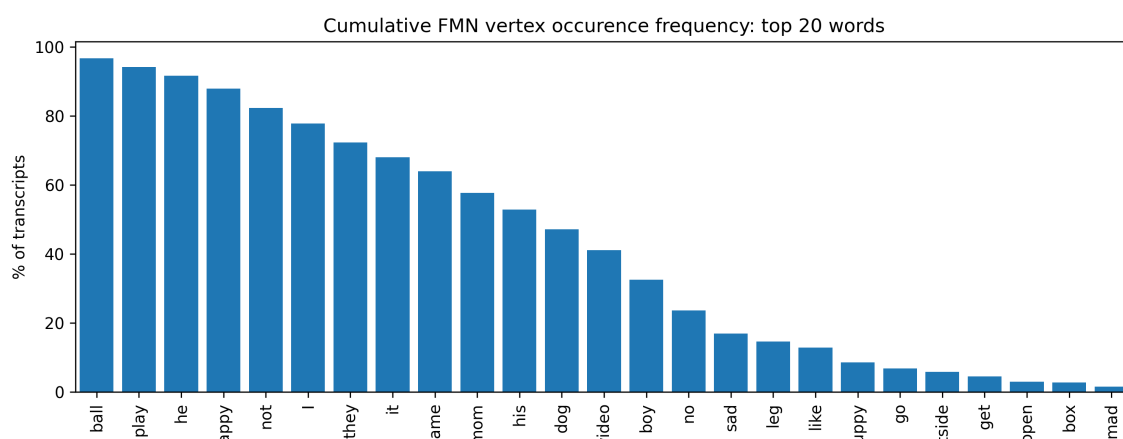
Transcript analysis report

Mutually included words:

No words occur consistently across all transcripts. And when we consider only the words which are useful in forming formant networks (FMN), ignoring non-informative words such as link-words, the most commonly occurring words across all transcripts is as follows:



Selecting a subset of these words which mutually appear in all transcripts quickly reduces the number of usable transcripts. The following plot shows the cumulative effect of including the most common words on the % of usable transcripts:

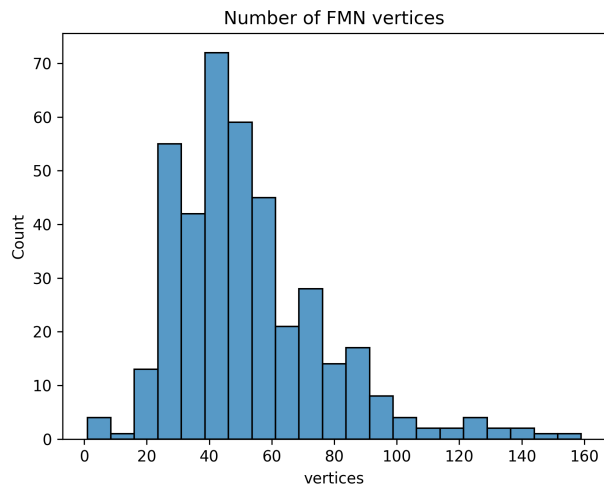


It therefore seems that a model which requires mutual inclusion of specific words will not be viable, as using even the 5 most popular words reduces the usable transcripts by ~20%, reducing the number of usable transcripts from 389 -> 318.

This would also cause the predictive ability of the model on new data to be dependent on inclusion of the specified words.

General data statistics:

The average number of FMN usable words in each transcript is 53, with the a small number of transcripts containing very few words:



Most sentences are short, with fewer than 5 words in them:

