

# Creating Semantic Spaces Using Large Document Clustering

Team: Tristan Weger, Yahya Emara, & Ryan Rubadue

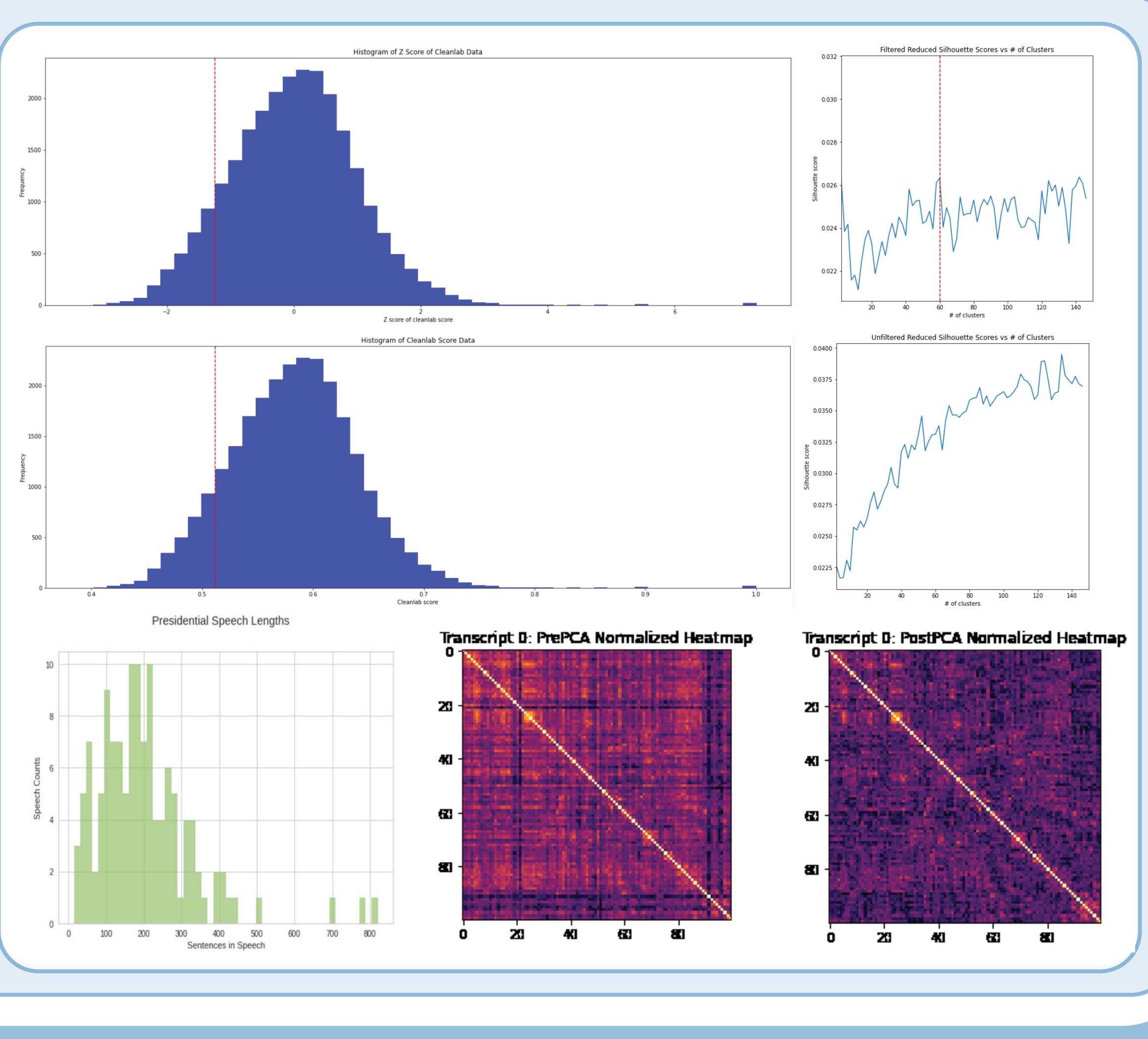
Advisor: Professor Ali Minai

## TASK

1. Analyze long text documents, where text for each doc is not confined to a single subject area
2. Create a representation of natural language that captures meaning (semantic space) using long documents
3. Compare performance of algorithm on filtered (using cleanlab) and unfiltered raw data

## ANALYSIS

- Obtain various metrics
- Compare
  - Segment vs sentence
  - Filtered vs unfiltered
  - Reduced vs not reduced
  - Optimized vs original
- Analyze
  - Cleanlab score
  - Silhouette score
  - Z score
  - Coherence score
  - Cosine Similarity

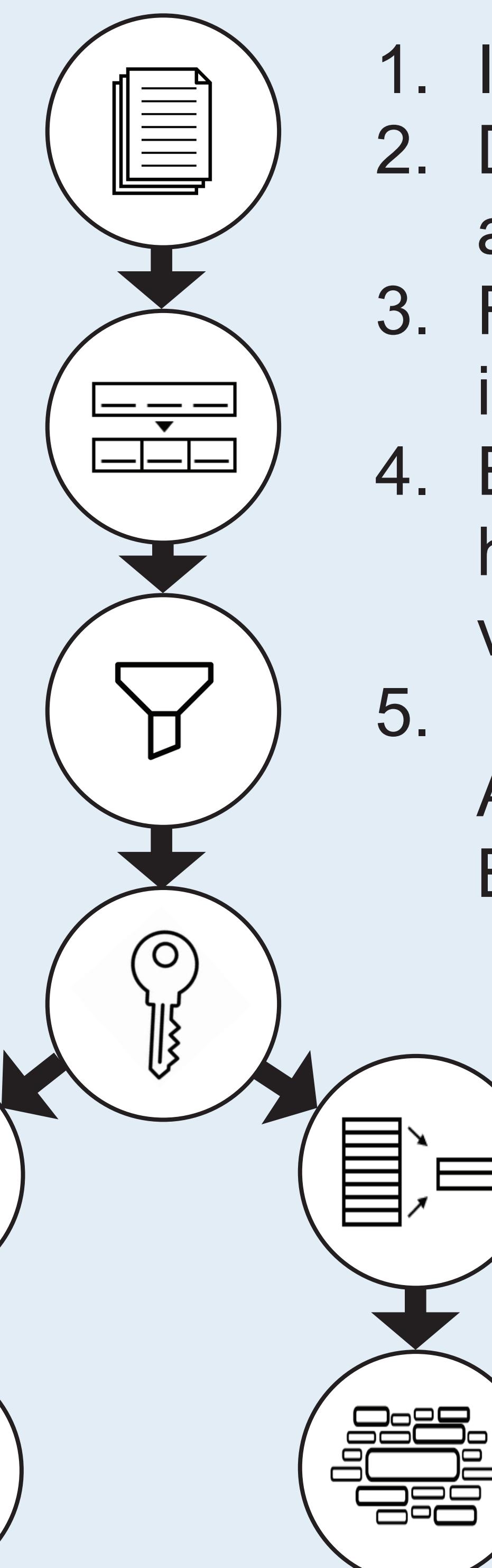


## TECHNOLOGY

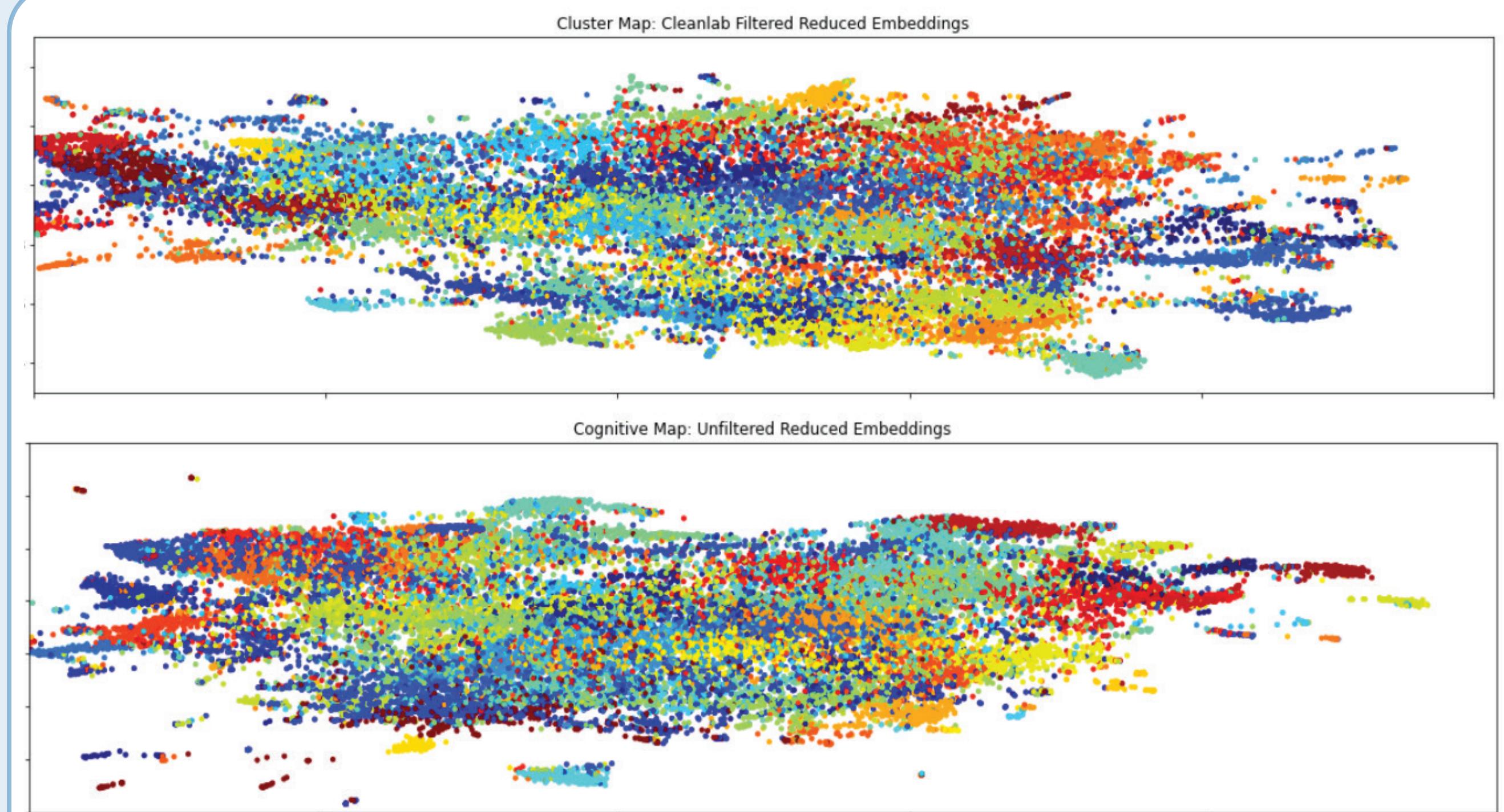
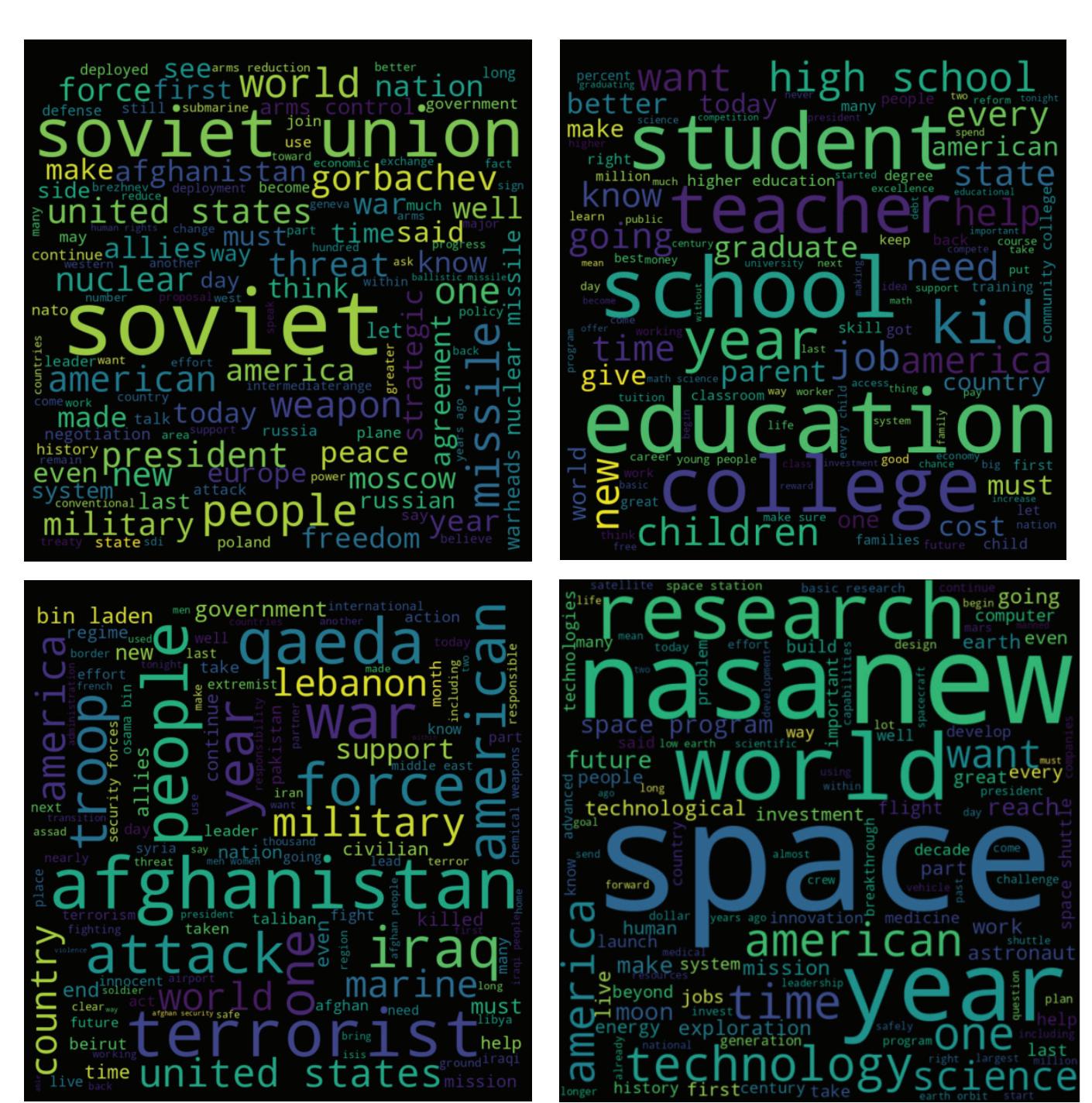


- Python
- NLTK
- SBERT
- UMAP
- Cleanlab
- Numpy
- Pandas
- Sklearn
- Torch
- Matplotlib

## ALGORITHM OVERVIEW



## RESULTS



## CONCLUSION

- Filtering Documents at Sentence Level is EFFECTIVE
- Observed X percent of sentences in presidential speeches is 'relevant'
- \*Splitting documents at sentence level is BETTER than at paragraph
- Reducing embedding dimensions resulted in more coherent outputs
- Creation of a functional semantic space was achieved