

# Optimizing A Minimal Core Vocabulary For Language Representation

This research investigates identifying a minimal core vocabulary that maintains expressive power in English. The study constructs a representative vocabulary using computational techniques, develops an algorithm to generate sentences with the core vocabulary, and measures its effectiveness. By integrating advanced language models, dimensionality reduction, and clustering methods, we optimize vocabulary size and evaluate efficacy. We analyze the relationship between similarity scores and vocabulary size using Cosine and Jaccard similarity metrics to assess semantic coherence.



by Roey Feingold,Yohai Haddad

# Vocabulary Building Process

## **Corpus-Based Selection**

Extract most frequent words from corpora like Brown, Reuters, and Gutenberg. Filter out stop words and non-alphabetic characters.

**WordNet Enhancer**

Expand vocabulary using WordNet to include hypernyms, covering broader semantic concepts.

## Embedding Filterin

Apply GloVe and BERT embeddings to refine vocabulary, pruning words with high semantic overlap.



# Sentence Generation Algorithm

- 1
- 2
- 3

## Preprocessing

Tokenize input, remove stop words, identify out-of-vocabulary words

## Word Substitution

Find closest words using cosine similarity of embeddings

## Sentence Reconstruction

Rebuild sentence using substitutions while maintaining structure

# Semantic Similarity Metrics

## Cosine Similarity

Measures angle between word vectors in semantic space. Range: -1 to 1. Higher values indicate greater similarity. Useful for comparing word embeddings.

## Jaccard Similarity

Measures overlap between word sets. Defined as intersection divided by union. Useful for evaluating lexical overlap between sentences.

## Coverage and Compression Ratios

Coverage: proportion of words representable by core vocabulary.  
Compression: reduction in vocabulary size while retaining meaning.

# Evaluation Methodologies

## Embedding-Based Substitution

Use word embeddings to find semantically similar words in reduced vocabulary. Calculate cosine similarity between original and substituted sentences.

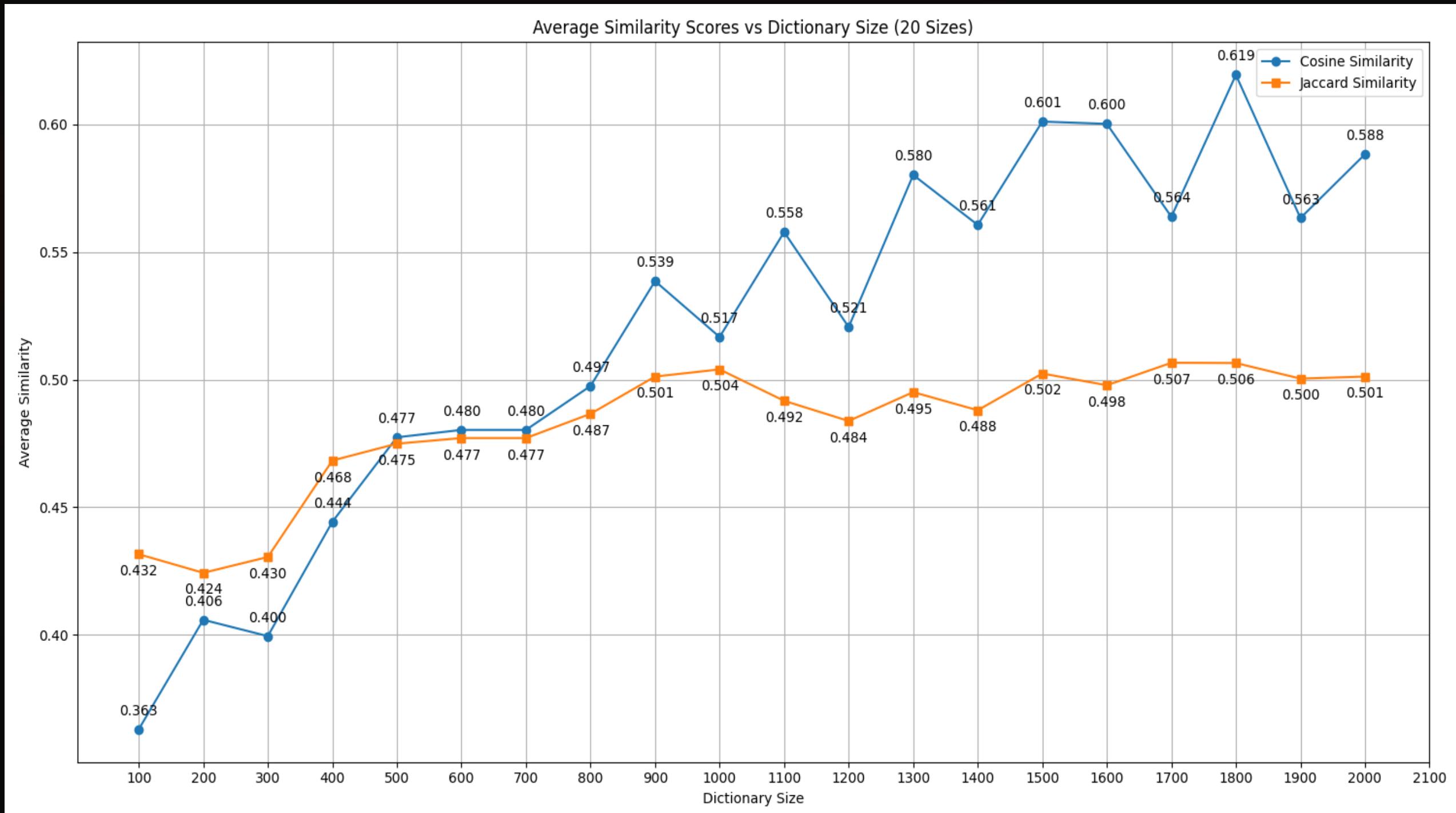
## Back-Translation Analysis

Translate sentence to another language and back. Assess how well reduced vocabulary captures nuances across languages.

## Human Evaluation

Have evaluators assess intelligibility and naturalness of sentences constructed using reduced vocabulary.

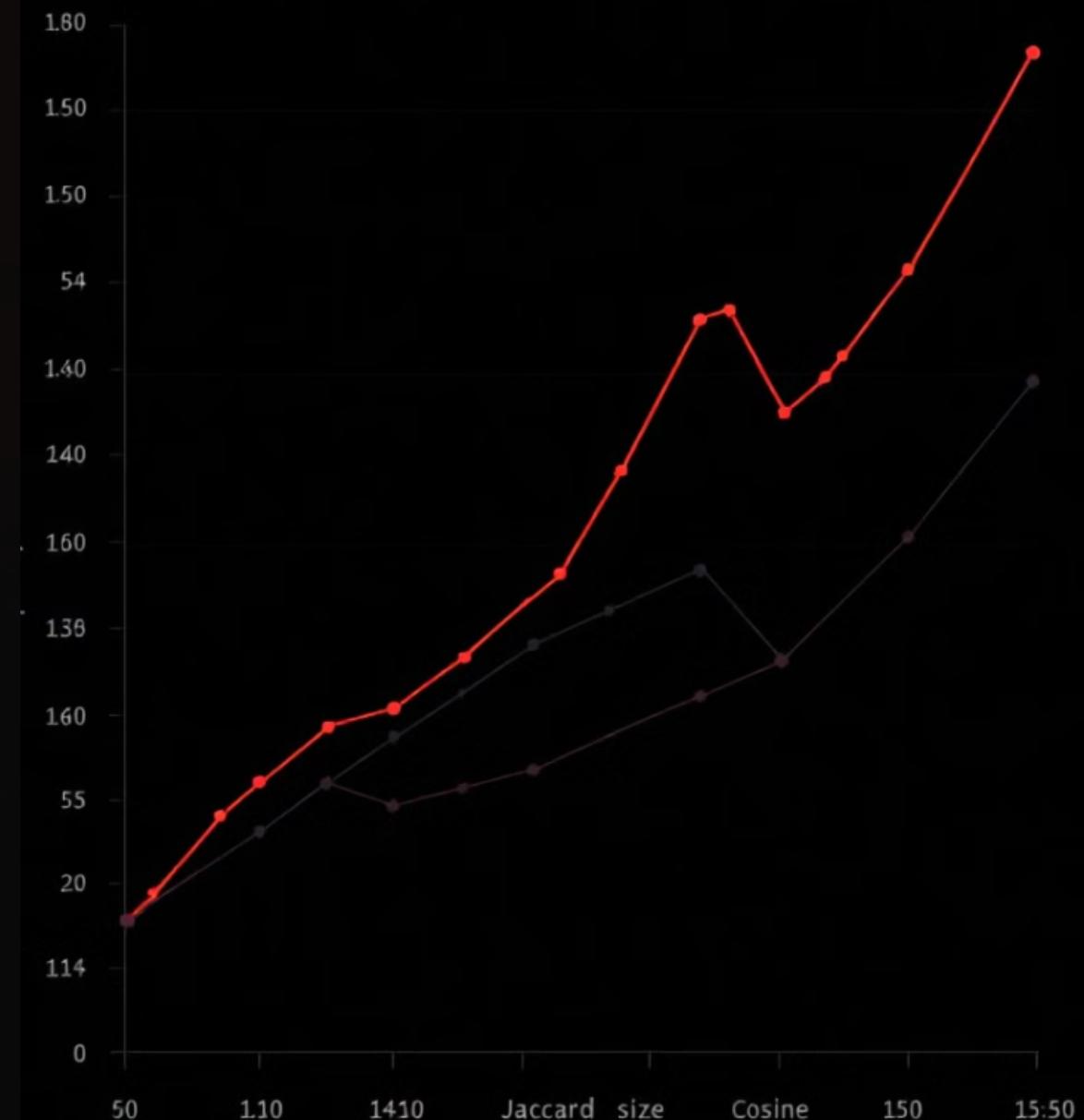


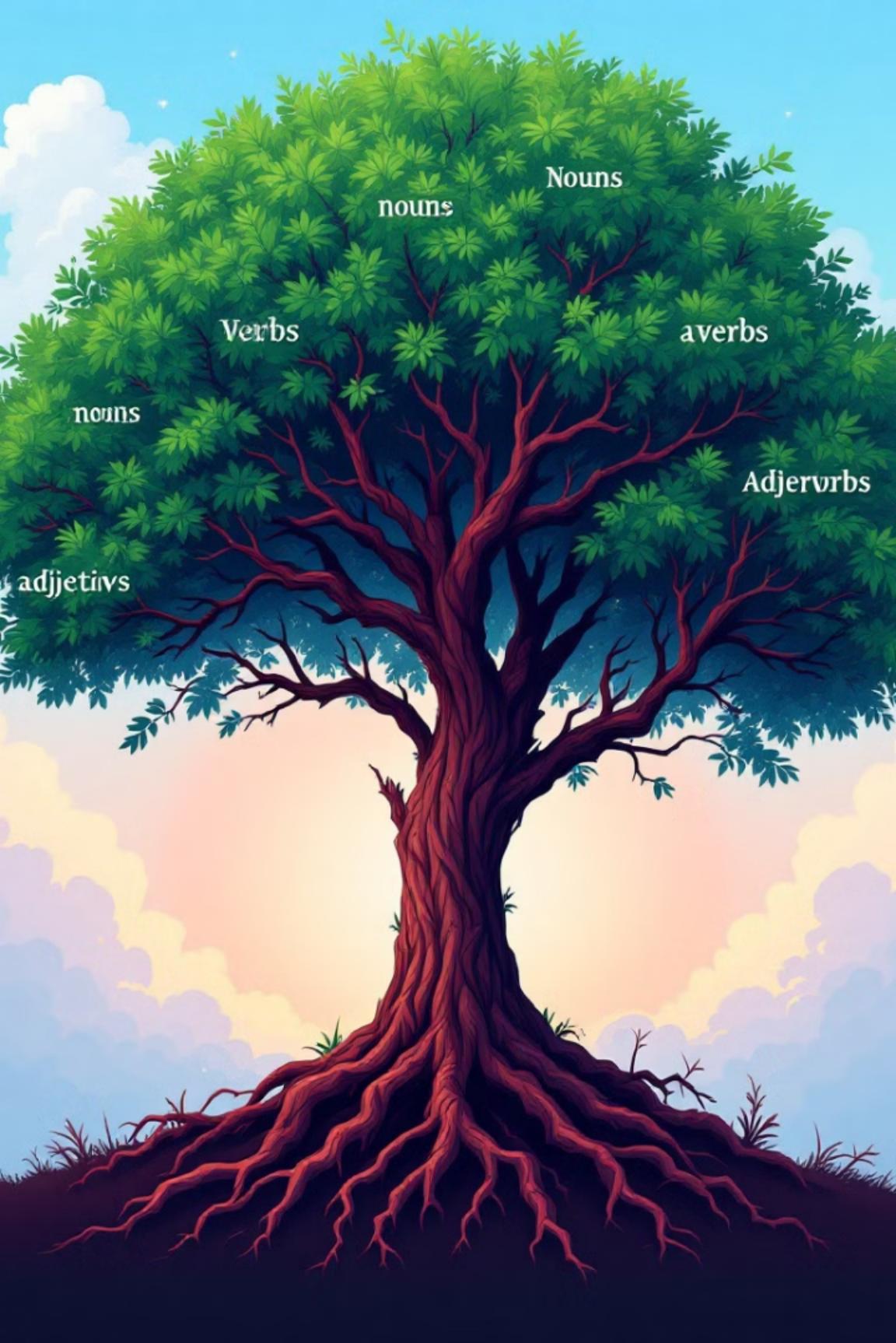


# SICK Dataset Results

Metric	SICK Creators	Our Approach
Cosine Similarity	0.6009	~0.619
Jaccard Similarity	0.4884	~0.501

Our approach using a limited vocabulary achieved comparable or slightly better results than the SICK dataset creators, who used unrestricted vocabulary. This demonstrates the effectiveness of our core vocabulary in maintaining semantic fidelity.





# Optimal Vocabulary Size

## 1 1000 Words

Lower bound for retaining significant semantic information. Suitable for simplified language models and beginner language learning.

## 2 1500 Words

Balances common terms and some specialization. Ideal for intermediate learners and systems requiring richer language understanding.

## 3 2000 Words

Upper range showing consistent improvement in similarity scores. Accommodates more nuanced expression while maintaining efficiency.

# Conclusion and Future Directions



## Refinement

Further refine core vocabulary selection using advanced techniques like contextualized embeddings.



## NLP Applications

Develop more efficient language models, especially for low-resource languages.

This research demonstrates the feasibility of reducing language to a core vocabulary while retaining communicative power. Future work should address limitations in specialized domains and subjectivity of similarity measures.



## Cross-Linguistic

Apply approach to other languages to explore universality of findings.



## Language Learning

Create simplified learning approaches focusing on mastering essential words first.