# Bible Semantic Similarity Explorer

**Student Name:** [Your Name Here] **Date:** January 2026 **Project:** NLP Practical Project

## 1. Abstract

The **Bible Semantic Similarity Explorer** is a web-based application designed to explore semantic relationships within the Romanian Bible text. By leveraging advanced Natural Language Processing (NLP) techniques — specifically Word2Vec and Doc2Vec — the application allows users to discover synonymous words, analyze analogies, and find semantically similar Bible passages. This project demonstrates the practical application of embedding-based models to capture and visualize the semantic richness of religious texts.

## 2. Introduction

Natural Language Processing enables computers to understand specific nuances in human language. In the context of religious study, finding connections between different passages or understanding the semantic field of specific theological terms can be challenging due to the size and complexity of the text.

### Objectives

- **Semantic Search:** Enable users to find words and passages based on meaning rather than just keyword matching.
- **Analogy Resolution:** Solve linguistic analogies (e.g., "Father is to Son as Mother is to ?") to uncover structural relationships in the text.
- **Quantifiable Similarity:** Provide mathematical metrics (cosine similarity) to quantify how closely related two concepts are.

### Dataset

The project utilizes the **Romanian Bible Paraphrase Dataset** (sourced from HuggingFace `andyP/ro-paraphrase-bible`).

- **Size:** Approximately 247,000 training pairs.
- **Content:** Various Romanian translations and paraphrases of Biblical texts.
- **Preprocessing:** The dataset was deduplicated and cleaned to ensure high-quality training data.

## 3. Methodology & Architecture

The system follows a standard NLP pipeline: Data Preprocessing $\rightarrow$ Model Training $\rightarrow$ Interactive Interface.

### 3.1 Data Preprocessing

Raw text cannot be directly fed into machine learning models. We implemented a robust preprocessing pipeline using **spaCy** (Romanian model `ro_core_news_sm`).

**Key Steps:**

1. **Tokenization:** Splitting text into individual units (words).
2. **Lemmatization:** Converting words to their dictionary root (e.g., "iubind" $\rightarrow$ "iubi").
3. **Noise Removal:** Eliminating punctuation, special characters, and stopwords.
4. **Normalization:** Converting all text to lowercase.

## 3.2 Model Architecture

Two distinct embedding models were trained using the **Gensim** library:

**Word2Vec (Continuous Bag of Words - CBOW)**

- **Purpose:** Captures the meaning of individual words based on their context.
- **Configuration:** 100-dimensional vectors, window size of 5, minimum frequency of 2.
- **Mechanism:** Predicts a target word given its surrounding context words.

**Doc2Vec (Distributed Memory - PV-DM)**

- **Purpose:** Captures the meaning of entire sentences or paragraphs.
- **Configuration:** 100-dimensional vectors, window size of 5, 20 training epochs.
- **Mechanism:** Learns a unique vector for each document that, when combined with context words, predicts the next word in the sequence.

## 3.3 System Design

The application is built using **Python** and structured as follows:

- `preprocessing.py`: Handles text cleaning and normalization.
- `models.py`: Manages the loading of the heavy ML models.
- `services.py`: Contains the core logic for similarity computations and analogy math.
- `ui.py`: A **Gradio** web interface that connects the backend logic to a user-friendly frontend.
- `app.py`: The entry point that orchestrates the application launch.

# 4. Features & Results

The application provides four main interactive features. Below are case studies demonstrating the system's capabilities.

## 4.1 Similar Words Discovery

Finds words that share semantic contexts with a query word.

- **Query:** "iubire" (love)
- **Results:**
    1. *dragoste* (0.82) - Synonymous
    2. *chemare* (0.79) - Contextually related
    3. *devoțiune* (0.75) - Theologically related

*Interpretation: The model correctly identifies "dragoste" as the closest synonym, capturing the nuanced usage of "love" in the Bible.*

## 4.2 Word Pair Similarity

Calculates a direct similarity score between two specific words.

- **Comparison:** "Dumnezeu" (God) vs. "Domnul" (Lord)
- **Score: 0.87** (Very High Similarity)
- **Comparison:** "Dumnezeu" (God) vs. "păcat" (sin)
- **Score: 0.32** (Low Similarity)

*Interpretation: The models successfully encode the semantic distance between divine concepts and their opposing theological concepts.*

## 4.3 Word Analogies

Solves vector arithmetic problems: $\vec{A} - \vec{B} + \vec{C} \approx \vec{D}$.

- **Query:** "tatăl" (Father) is to "fiul" (Son) as "mama" (Mother) is to...
- **Result:** *fiica* (Daughter) - Score: 0.89

*Interpretation: The model has learned gender and familial relationships implicitly from the text structures.*

## 4.4 Sentence Similarity

Retrieves Bible passages that match the semantic meaning of a user's input, even if they don't share exact keywords.

- **Query:** "Dumnezeu este iubire" (God is love)
- **Top Match:** "Cine nu iubește, n-a cunoscut pe Dumnezeu; pentru că Dumnezeu este iubire." (1 John 4:8)

# 5. Technical Discussion

## Performance

- **Inference Speed:** Word queries are near-instantaneous (<20ms). Sentence queries take slightly longer (~200ms) due to the inference step in Doc2Vec.
- **Memory Footprint:** The application requires approximately 400MB of RAM to hold both models and the spaCy pipeline in memory.

## Technologies Used

- **Gradio:** For rapid UI prototyping and interactive visualization.
- **Gensim:** Use-optimized library for training vector embeddings.
- **spaCy:** Industrial-strength NLP for preprocessing.

# 6. Conclusion

The Bible Semantic Similarity Explorer successfully bridges the gap between raw text and semantic understanding. By processing nearly 250,000 Bible phrases, it offers a tool that can aid theological study, translation comparison, and linguistic analysis. Future work could include visualizing these embeddings in 2D space or expanding the dataset to include commentaries and cross-references.

---

## 7. References

1. **Mikolov et al. (2013).** *Efficient Estimation of Word Representations in Vector Space.* (Word2Vec)
2. **Le & Mikolov (2014).** *Distributed Representations of Sentences and Documents.* (Doc2Vec)
3. **Gradio Documentation:** https://gradio.app/
4. **HuggingFace Datasets:** andyP/ro-paraphrase-bible