

# ROSETTA STONE

PRESENTED BY ELENA DEL ROSAL

# DATASET STRUCTURE

[illegible]

# LIBRARIES USED

- pandas – For data manipulation and DataFrame operations.
- numpy – For mathematical operations, normalization, and array handling.
- sentence-transformers – To generate sentence embeddings using models like LaBSE.
- torch (PyTorch) – Backend framework used by sentence-transformers.
- scikit-learn (sklearn) – For:
  - Evaluation metrics (accuracy, MSE)
  - Confusion matrix and classification report
  - Optional scaling or normalization
- deep-translator – To automatically translate sentences using GoogleTranslator.
- Parrot – For generating paraphrased versions of English sentences.
- spaCy – To compute semantic similarity using pre-trained GloVe word embeddings.
- matplotlib – For creating scatter plots, bar charts, and data visualizations.
- seaborn (optional) – For visualizing the correlation matrix (e.g., heatmaps).
- time – To pause between translation attempts and handle retries safely.

# MAIN PURPOSE:

analyzing the semantic similitude between the sentences

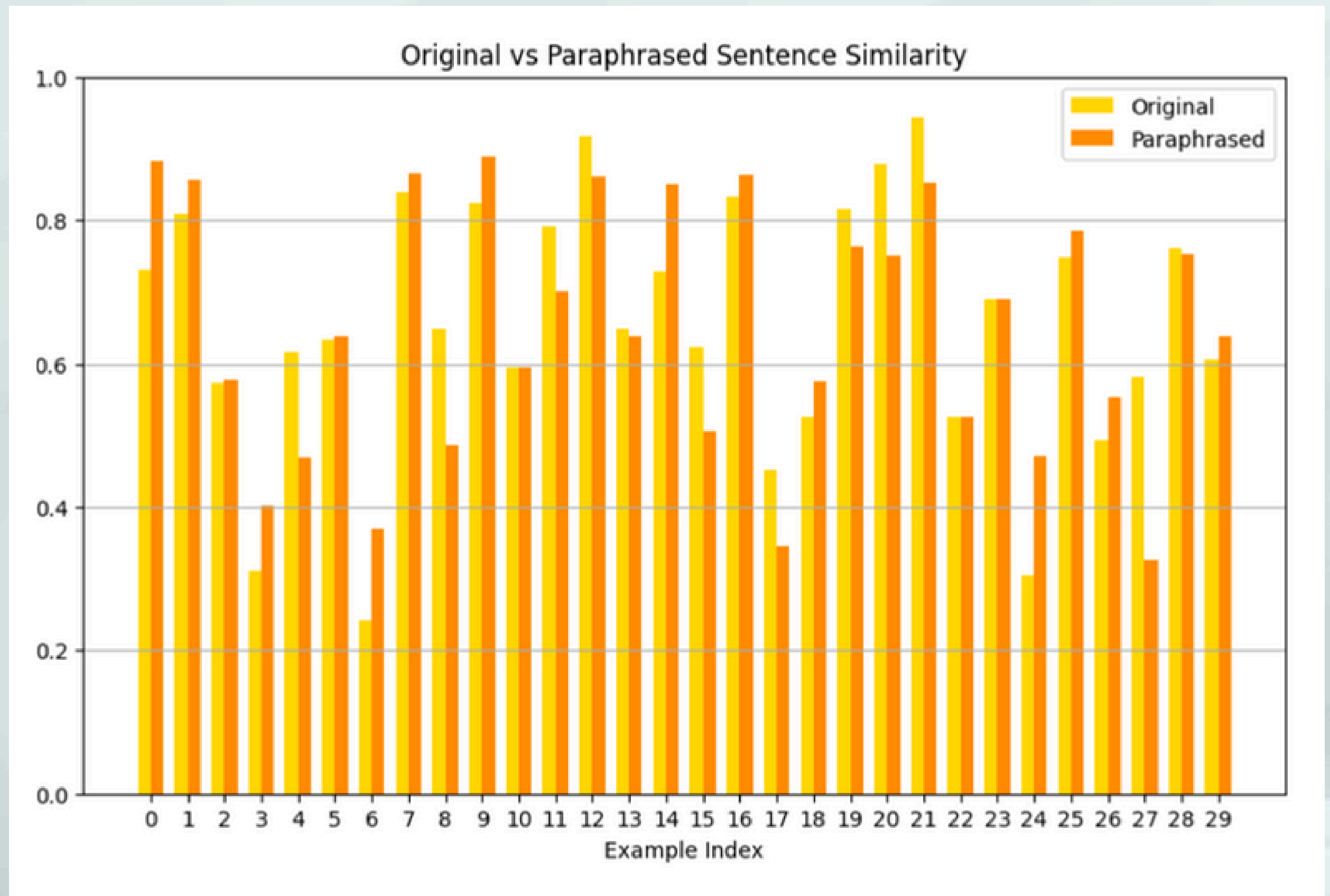
# OBJECTIVE:

comparing two different models to see how they perform with the same task and to see their differences in a practical way.

# PROCEDURE

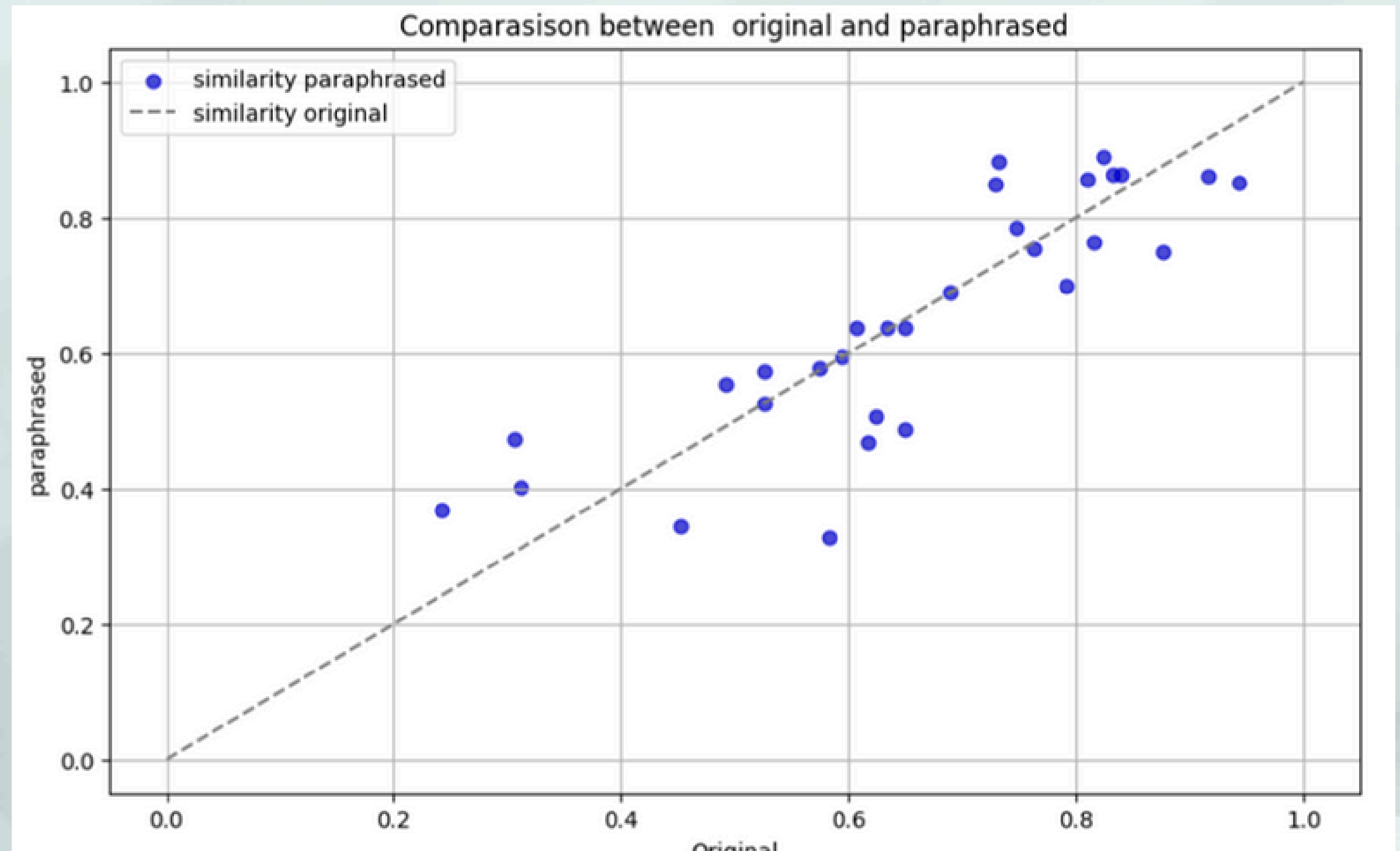
- 1 Data processing: Stemming
- 2 Translation
- 3 Data augmentation: Paraphrasing
- 4 Transformers model
- 5 SpaCy model

IS  
PARAPHRASING  
AS IMPORTANT  
AS WE THINK?

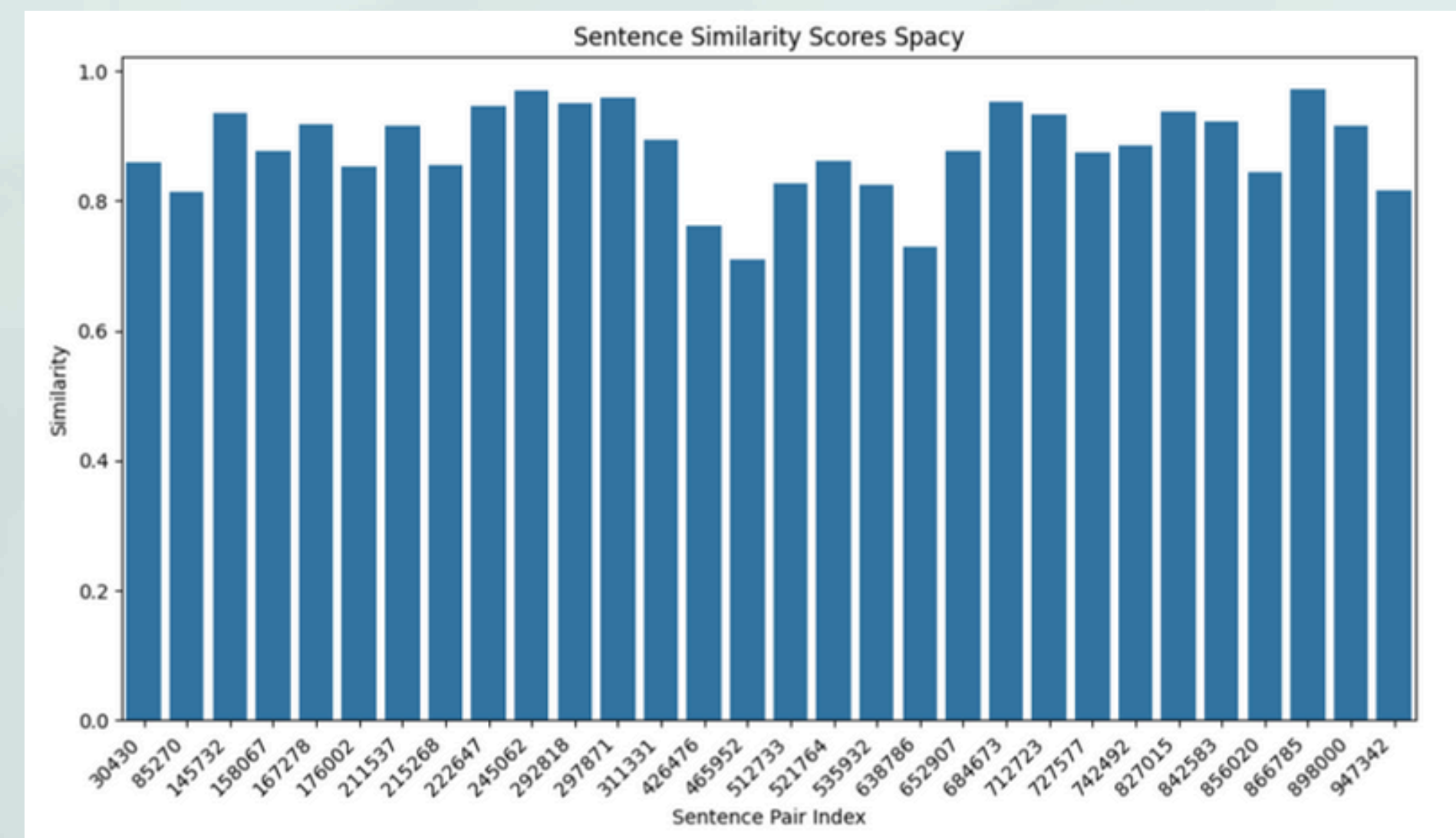
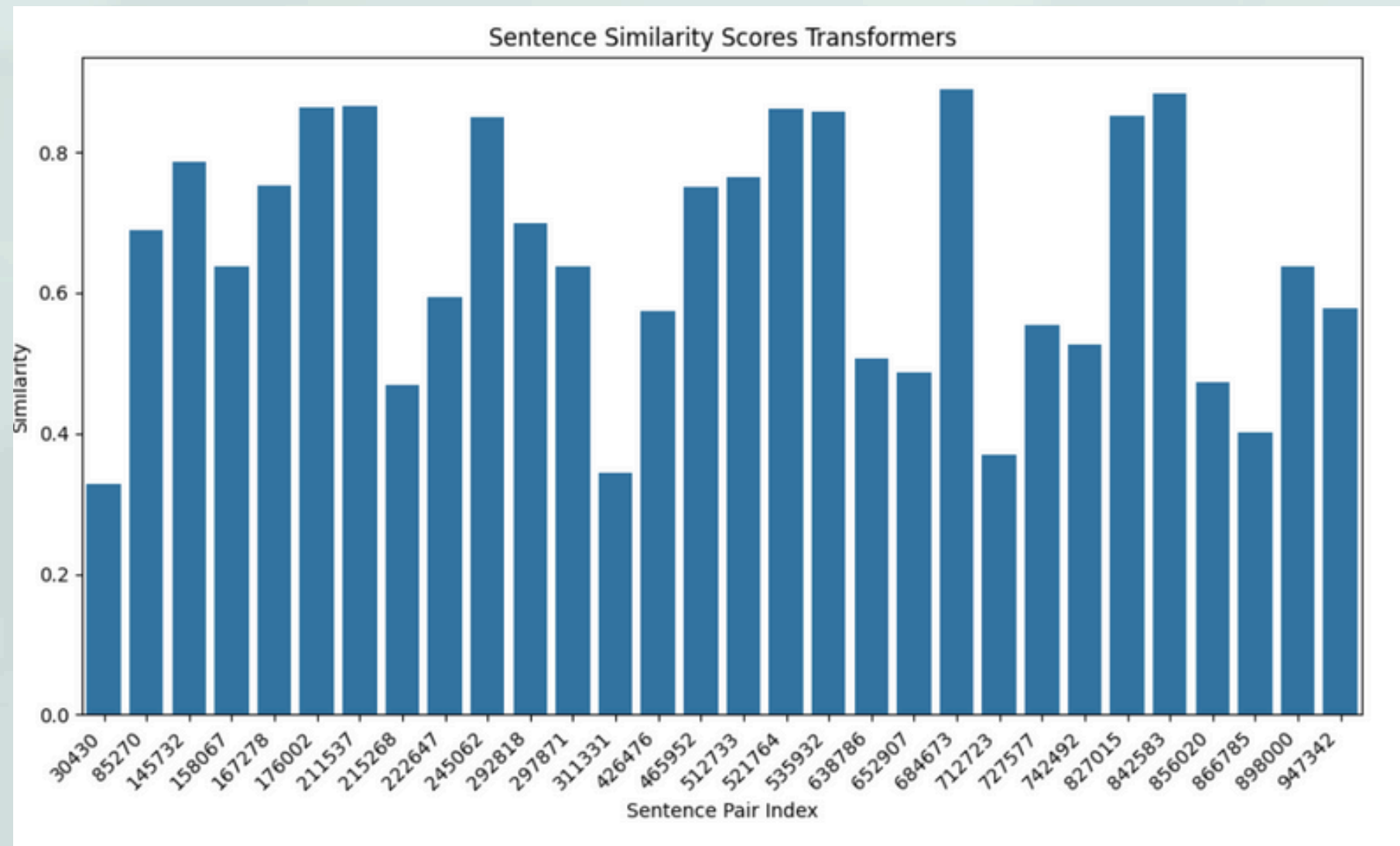


# SCATTER PLOT

compares original similarity with the paraphrased sentences similarities.



# TRANSFORMERS VS SPACY

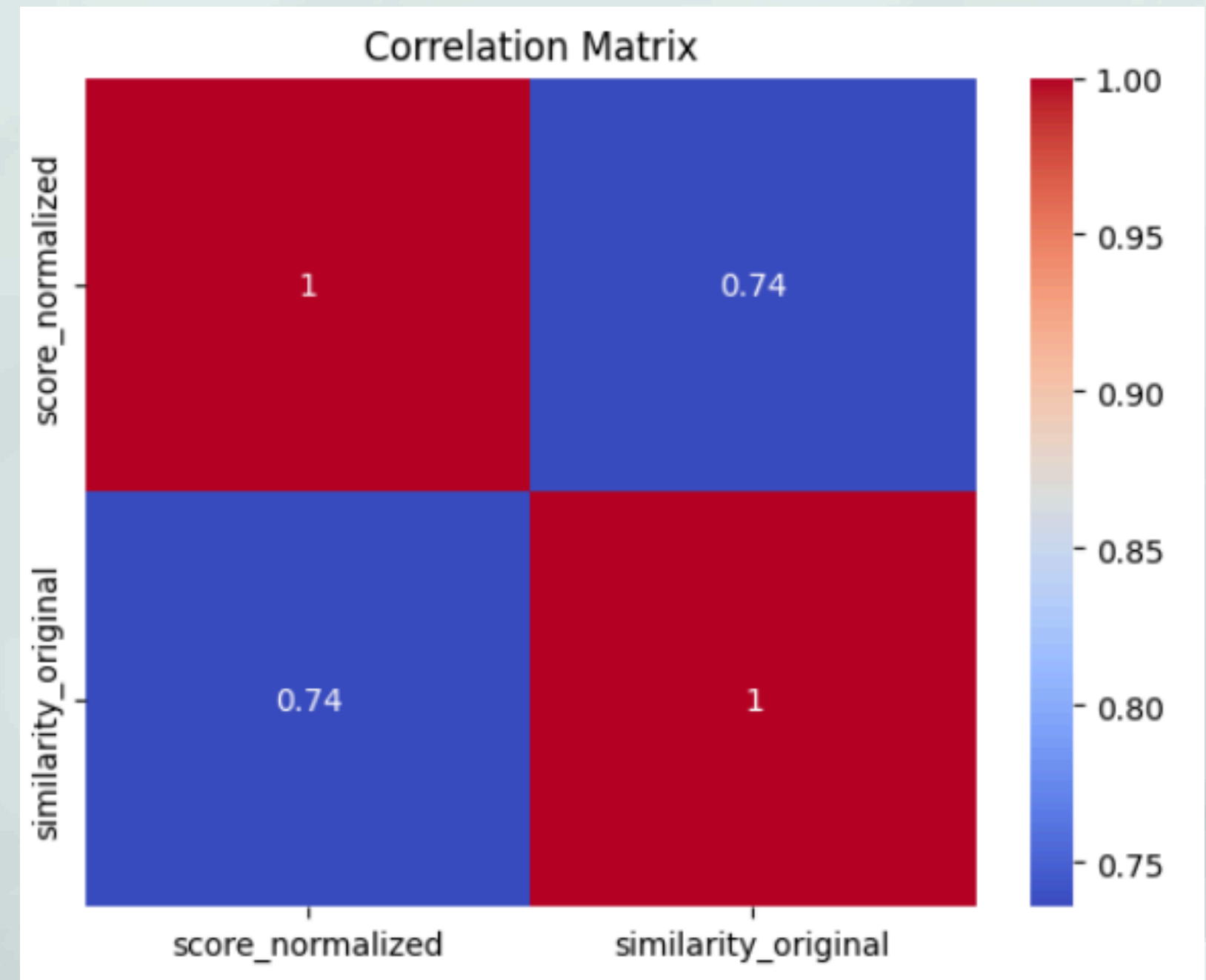


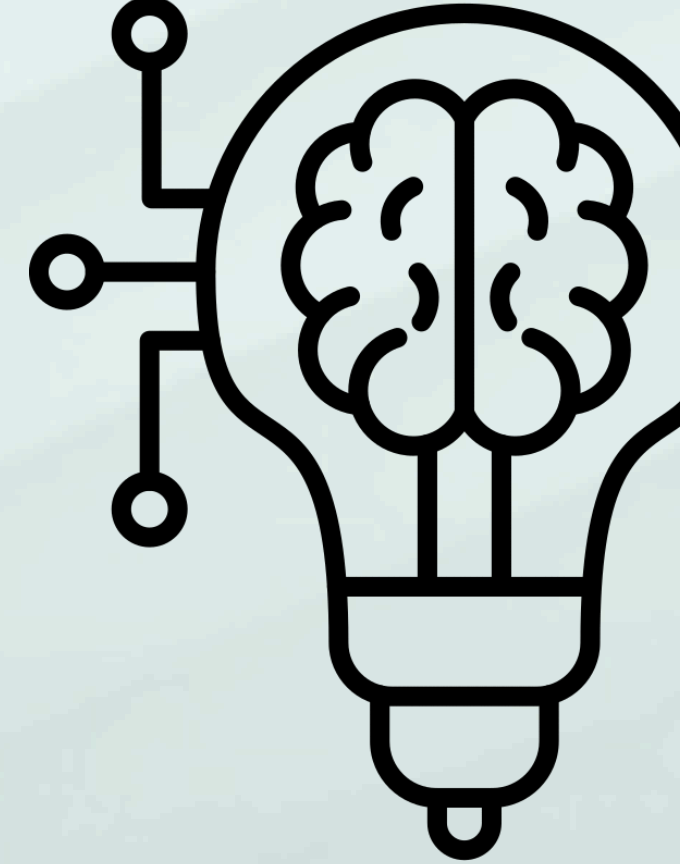
As we can see, SpaCy similarities are higher, which makes sense because it doesn't take into account the context, this is due to its reliance on word-level embeddings, which can overestimate similarity based on vocabulary overlap. Transformer-based models, such as LaBSE, typically provide lower but more accurate scores by capturing the full contextual meaning of sentences.



# PERFORMANCE OF MY TRANSFORMER MODEL

As we see, performance is really good. This correlation matrix is indicating a strong positive relationship. This suggests that the transformer-based model effectively captures semantic equivalence in line with the reference scores provided in the dataset.





# THANK YOU!

