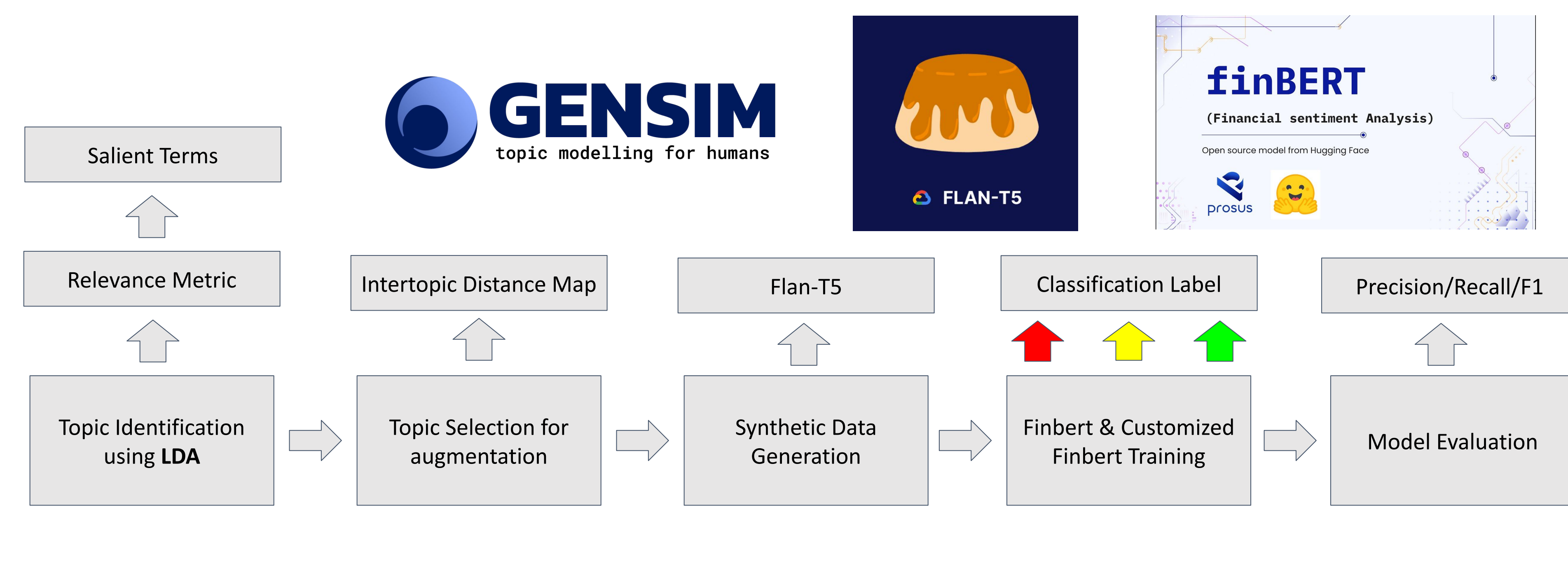
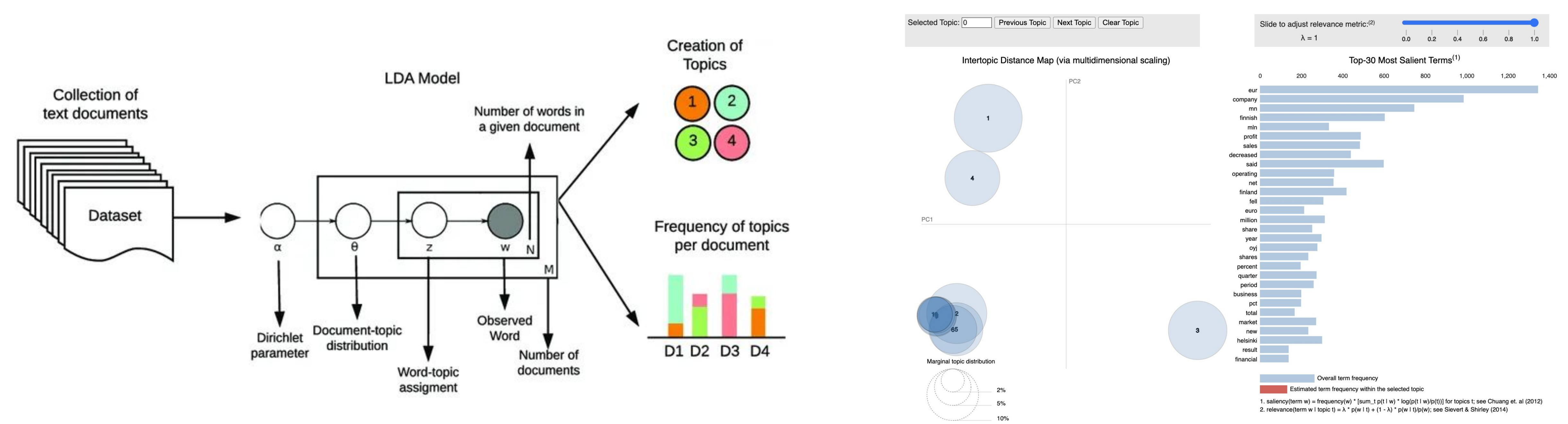


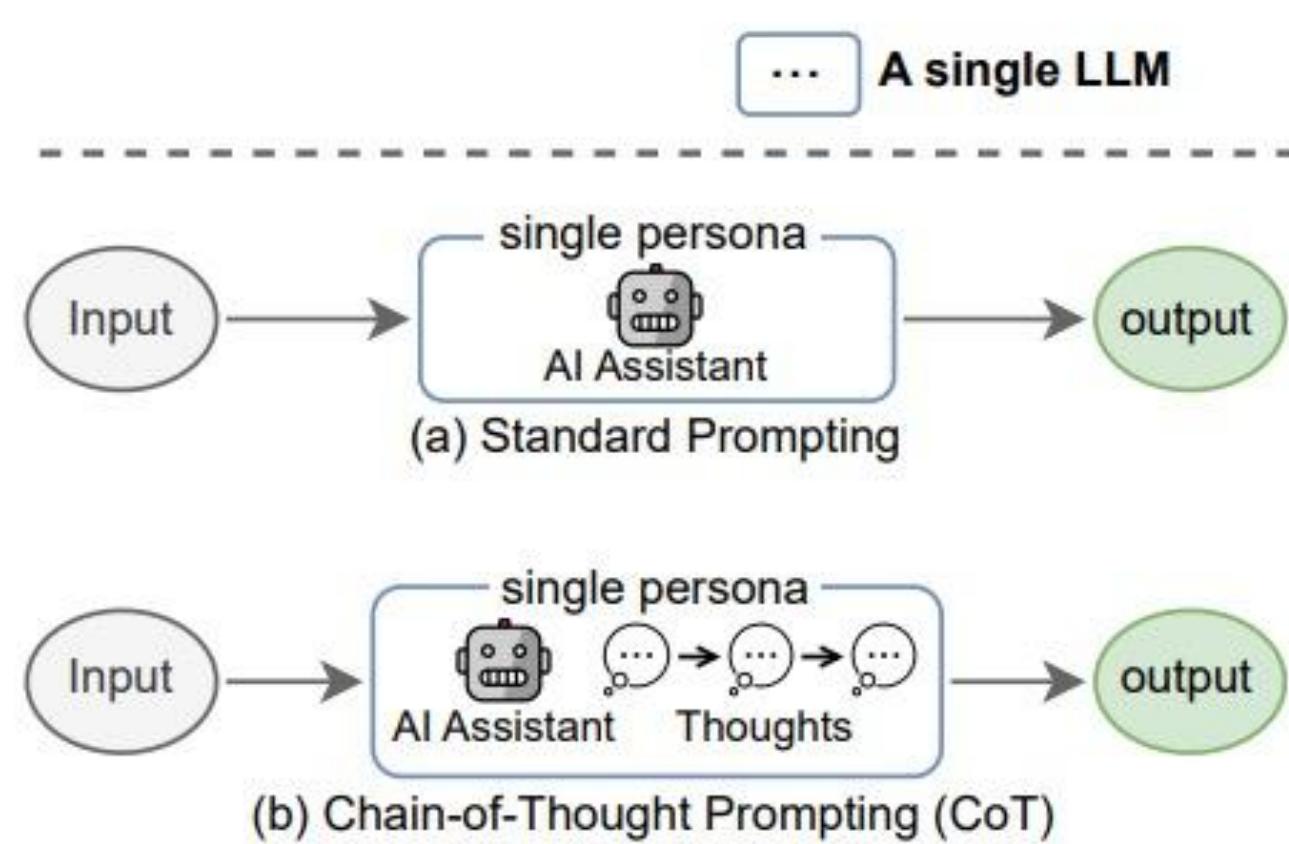
Process Flowchart



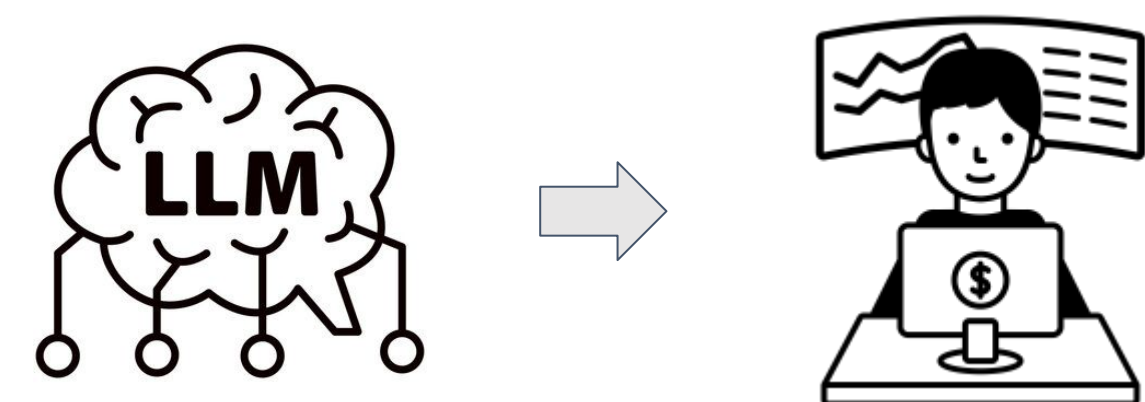
Latent Dirichlet Allocation



Chain of Reasoning



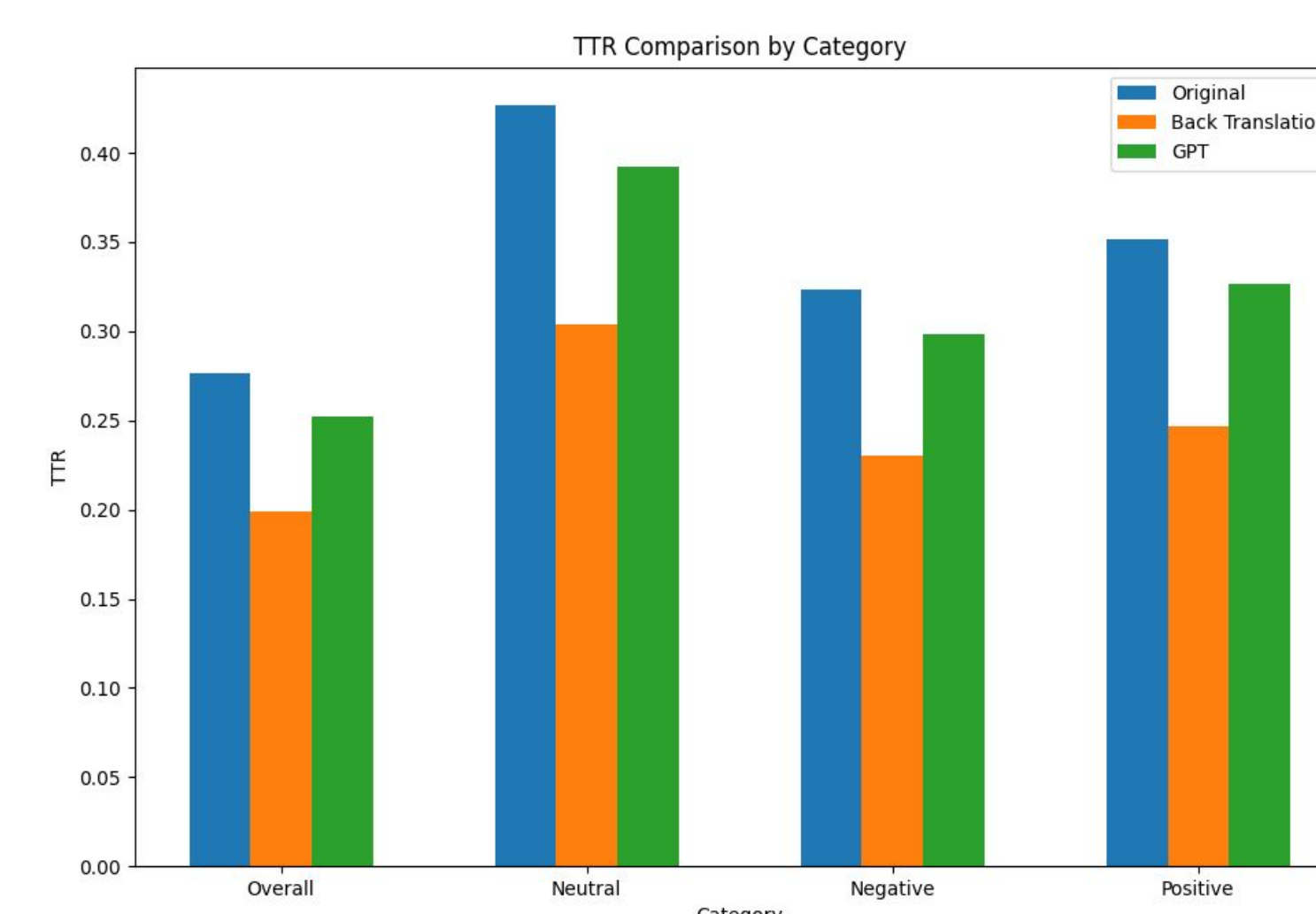
Persona



The above words are the resultant keywords for topic 3 when performing LDA. Act as a financial sentiment analysis expert. Create five positive, five negative, and five neutral financial headlines as prompts for future synthetic data generation, and explain why you present these headlines using the keywords for this topic.

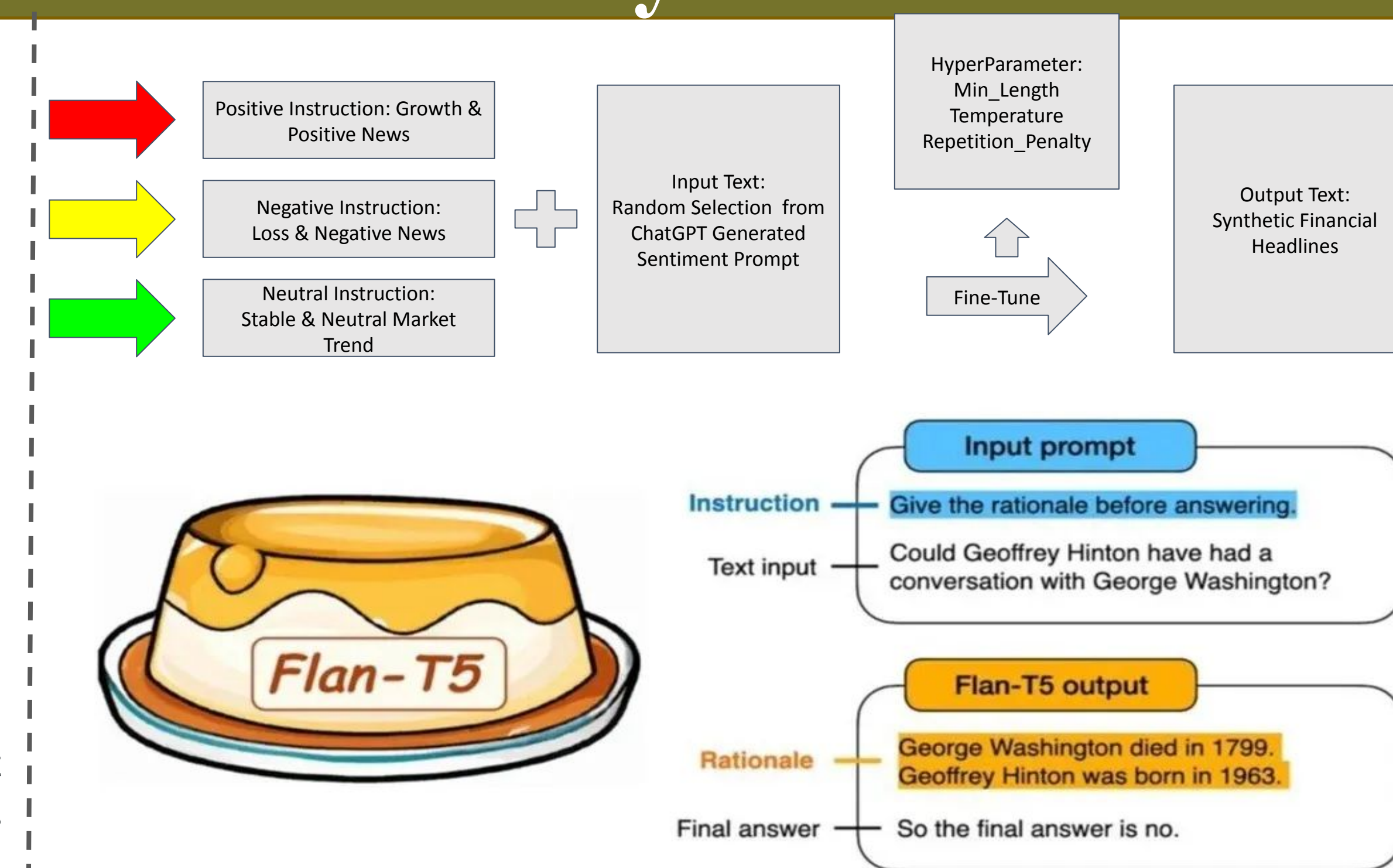
Positive Financial Headlines:

- "EUR 500mn Profit Reported as Sales Increase in Latest Quarter"**
Reason: This headline indicates a profit of EUR 500mn and increased sales, reflecting a positive financial outcome for the company.
- "Operating Profit Surges in EUR as Net Sales Jump 15% Year-over-Year"**
Reason: Words like "surges" and "jump" suggest strong financial performance, making this a positive headline.
- "EUR 200mn Net Profit Outpaces Expectations in a Strong Q1 Performance"**
Reason: "Outpaces expectations" and "strong performance" indicate better-than-expected results, creating a positive sentiment.
- "Sales and Profit Both Increase in Latest EUR Quarterly Report"**
Reason: An increase in both sales and profit is always a positive financial indicator.
- "Company Reports EUR 300mn in Net Profit, Highlighting Strong Year-End Growth"**
Reason: This headline emphasizes a significant profit figure and highlights growth, which are positive indicators for investors.

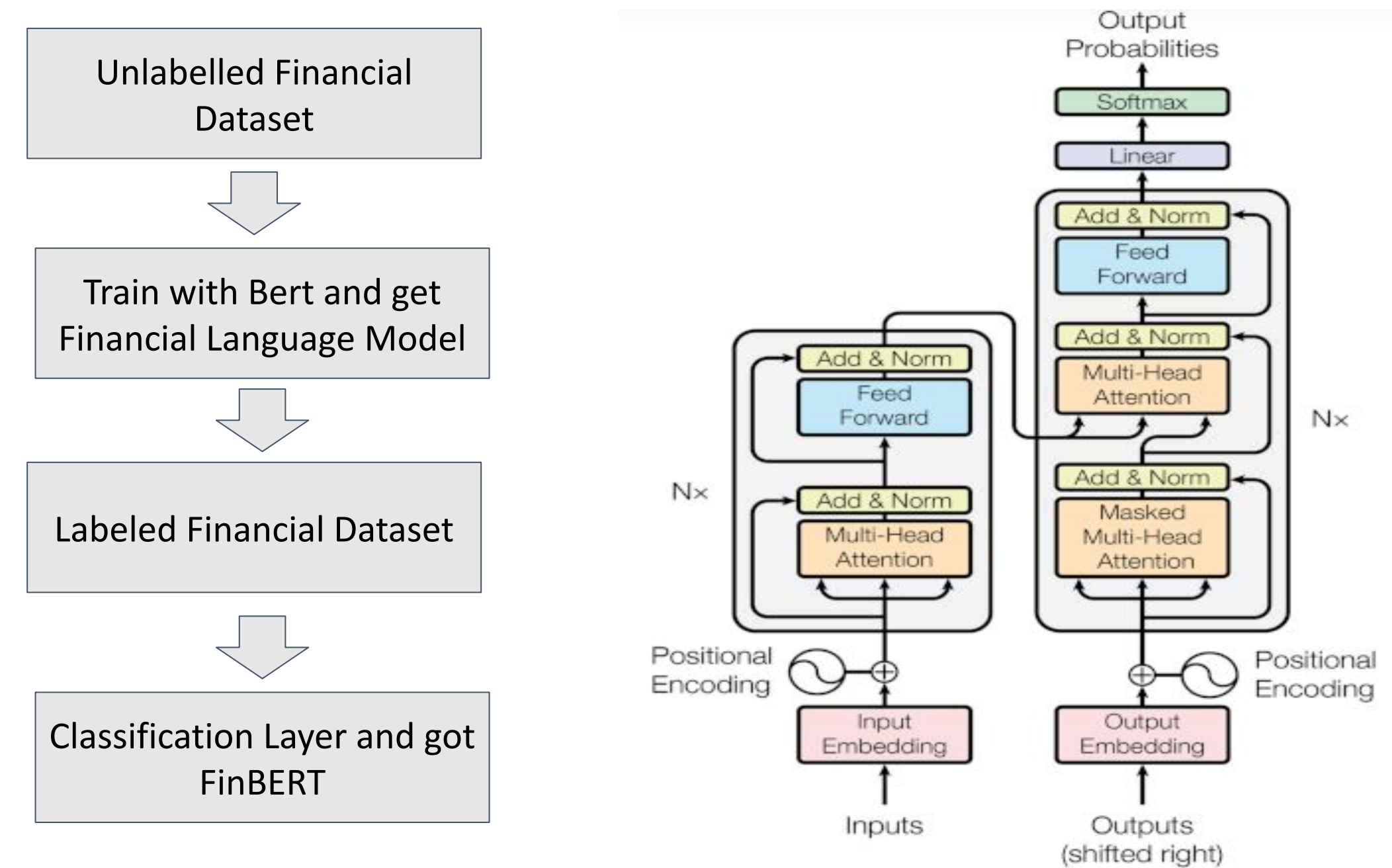


Measure of Lexical Diversity & Higher TTR result suggests greater variety in word usage and is beneficial for sentiment analysis.

Flan-T5 Synthetic Data



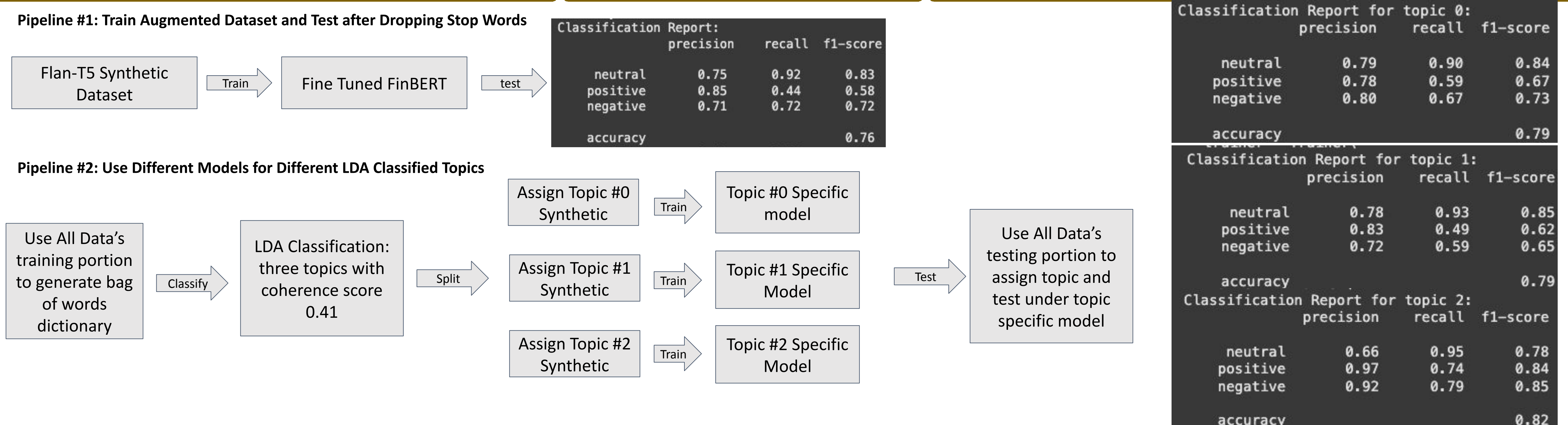
FinBERT



Base Architecture:

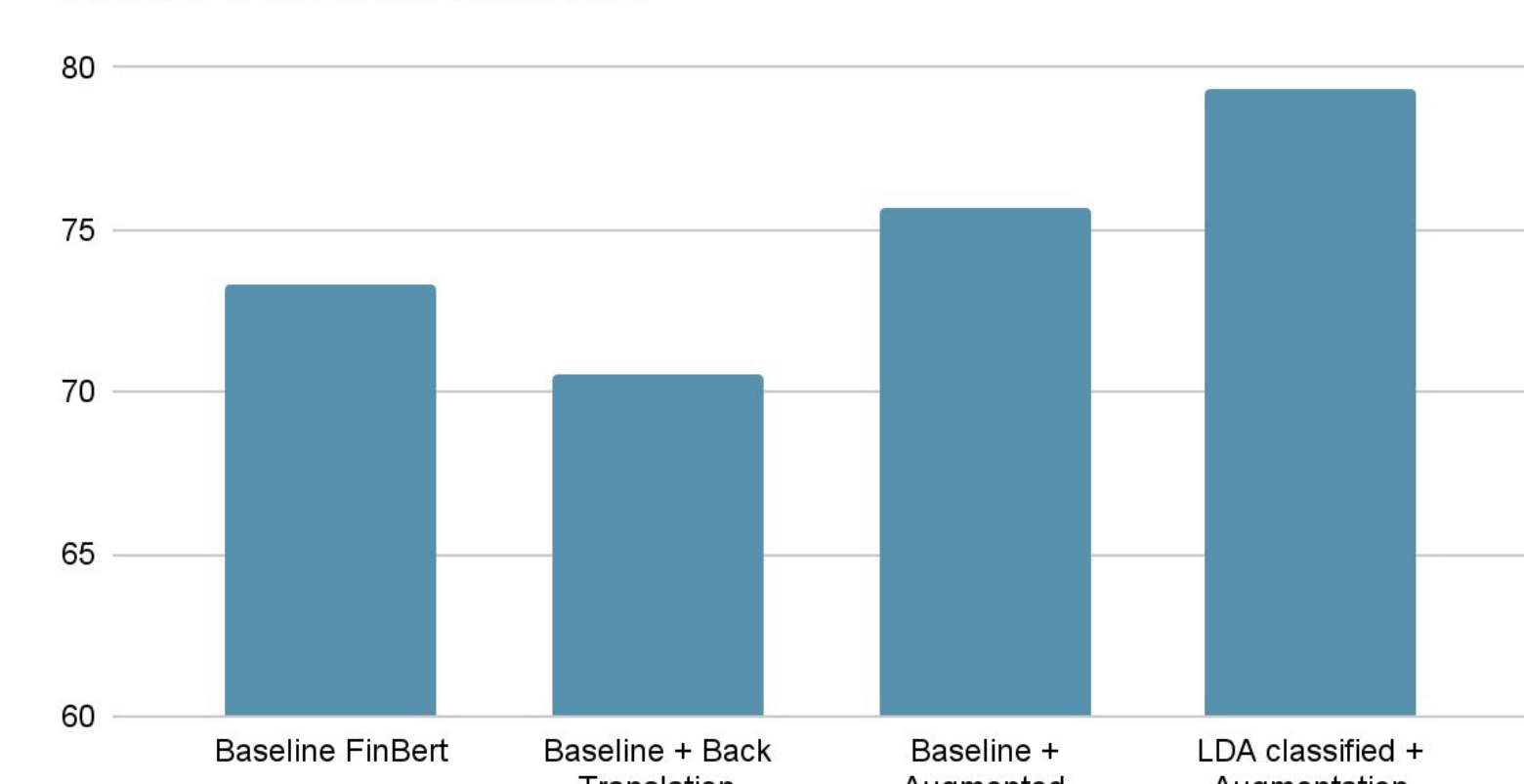
- BERT Core:** 12 layers, 12 attention heads per layer, 768 hidden dimensions, positional embeddings.
- Pre-trained BERT Model:** FinBERT starts from BERT's pre-trained model and fine-tunes it for financial sentiment tasks, retaining BERT's structural advantages.
- Fine-tuning Process:**
 - Financial Dataset:** Unlike generic datasets, FinBERT is trained and fine-tuned using financial texts like **company filings**, **earnings call transcripts**, and **financial news**.
 - Sentiment-Specific Output:**
 - Classification Layer:** While the core transformer architecture is identical to BERT, FinBERT introduces a specialized **classification head** fine-tuned for sentiment analysis tasks.

Implementation Pipelines



Result Comparisons

Test Accuracies



Under the augmentation steps that we took, our synthetic dataset used for training out-performed the Baseline FinBERT by 2.5% and the clustered training approach's average testing accuracy outperformed Baseline FinBERT by 7%. The Back Translation Approach did not do well as it lowers the variability of training dataset.

Potentials and Future Work

- Goal: Utilizing CleanLab to identify potential label issues in dataset. Result from the pretrained model is used to compare with the actual label of the data point, a larger difference indicates a lower confidence.
- 3 Trials CleanLab:
- With built-in Linear Regression Model
 - Issue labels identified, but the model is not well-trained.
 - With baseline Bert Model on augmented data
 - Top 20% data points with low confidence are identified with potential label issues.
 - After manual inspection, identified label issues are actually correct, proving the quality of the augmented data.
 - With small batch of kaggle data (10% of the original dataset)
 - Finbert used as the pretrained model with 83.51% of accuracy.
 - 0 label issues was identified.

Teacher Student Network and Sentence BERT

- Teacher Student Network is able to fine-tune a model based on incoming data and a baseline teacher model
- We can potentially apply Sentence BERT on the Student Network to let it understand more about the sentence embedding instead of only the tokens embedding. Pairs of words can be understood better

References

- [FinBERT: Financial Sentiment Analysis with Pre-trained Language Models](#) - Dogu Tan Araci
- [Conditional BERT Contextualize Augmentation](#) - Xing Wu
- [Confident Learning: Estimating Uncertainty in Dataset Labels](#) - Curtis G. Northcutt
- [Learning with Confident Examples: Rank Pruning for Robust Classification with Noisy Labels](#) - Curtis G. Northcutt
- [What is Latent Dirichlet Allocation?](#) - IBM
- [Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks](#) - Nils Reimers