



Goal : to analyze the performance of data classification using LLMs with traditional deep learning/machine learning models.

Scenarios 2 types:
Multi class Classification-->Classifying employee work locations based on job reviews.
Binary Classification -->Fake news detection in news articles

Testing models : Various LLMs with different sizes, quantization, and architectures + traditional ML models

Evaluation Metrics : Weighted F1 score

Key analysis :
Impact of prompt engineering techniques on LLM performance.
Trade-off between performance (F1 score) and time (inference time)

Key results
Differences in model responses based on prompting strategies.
Superior performance of LLMs (Llama3, GPT-4) in complex classification tasks (but with increased inference time).
Better performance-time balance of ML models in simple binary classification tasks.

Three main architectures

- Encoder-only (e.g., BERT)**: Uses only transformer encoder layers. Employs Masked Language Modeling (MLM) during training to predict masked tokens based on bidirectional contextual embeddings.
- Decoder-only (e.g., GPT)**: Uses only transformer decoder layers. Processes input sequentially and predicts the next token based on previous tokens.
- Encoder-decoder (e.g., T5)**: The encoder processes the input into an encoded representation, and the decoder reconstructs the output sequence step-by-step.

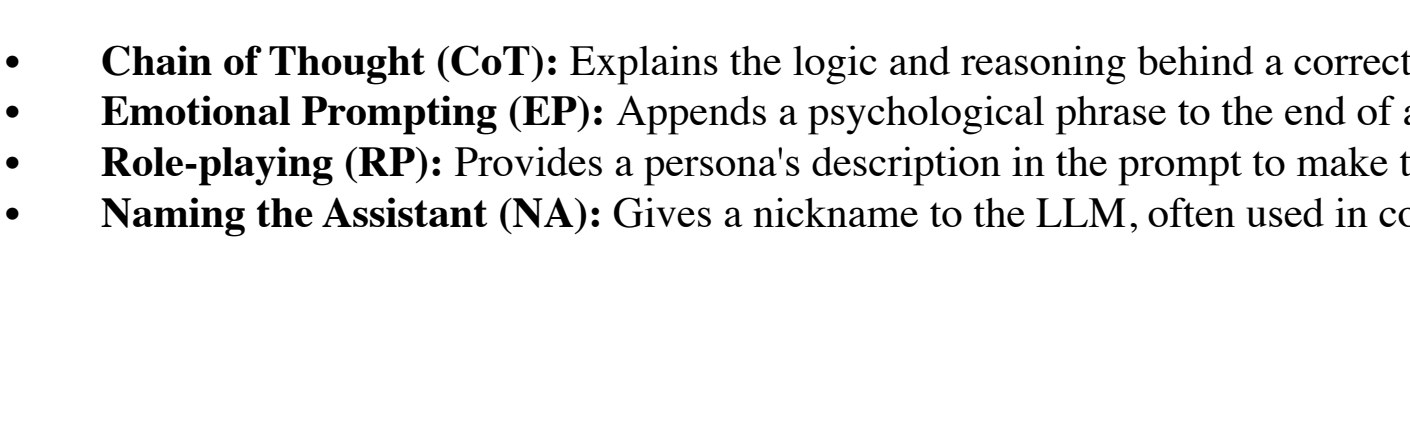
- Temperature**: Controls randomness (lower = predictable, higher = creative).
- Top-k Sampling**: Limits word choices to the top k most likely options.
- Max Tokens**: Sets the maximum number of tokens the model can generate.

Decoder only models :

- 1 Mistral-7B OpenOrca (Mistral-OO)
- 2 OpenHermes 2.5 Mistral-7B (Mistral-OH)
- 3 zephyr-7B-beta (Zephyr)
- 4 Nous-Hermes Llama2-13B (Llama2)
- 5 Xwin-MLewd 13B v0.2 (Xwin)
- 6 Gemma 2 9B (Gemma)
- 7 Meta Llama 3 70B (Llama 3 70B)
- 8 Meta Llama 3 8B (Llama 3 8B)
- 9 Mistral 8x7B (Mistral)
- 10 gpt-4-turbo (Gpt4-turbo)

Encoder only models : RoBERTa (Robustly optimized BERT approach)

Prompt engineering:Decoder only models : construct using zero-shot and few-shots



Used prompt engineering techniques:

- Chain of Thought (CoT)**: Explains the logic and reasoning behind a correct response, providing a step-by-step thought process to the model.
- Emotional Prompting (EP)**: Appends a psychological phrase to the end of a prompt to act as emotional stimuli for the LLM.
- Role-playing (RP)**: Provides a person's description in the prompt to make the LLM adopt a specific character and perspective.
- Naming the Assistant (NA)**: Gives a nickname to the LLM, often used in combination with Role-playing by appending the name to the beginning of the instruction.

Used traditional ML models:

- Naive Bayes (NB)**
- Support Vector Machines(SVM)**

Configurations and Implementations

- Hardware:**
- Groq models (Mistral, Llama3 70B/8B, Gemma) used Groq LPU via API.
 - Gpt4-turbo accessed via OpenAI API (hardware details undisclosed).
 - AWQ-quantized models (Llama2, Xwin, Mistral-OO/OH, Zephyr) loaded on NVIDIA Tesla T4 GPU in Google Colab.
 - RoBERTa and ML algorithms also trained on T4 GPU in Google Colab.
- Hyperparameters:**
- temperature set to 0 to eliminate randomness.
 - Outputs requested once per prompt

RoBERTa Training:

- Transformers library and PyTorch used for fine-tuning.
- Adam optimizer used with 5-fold cross-validation.
- Learning rate: 1e-5, Batch size: 32.

NB and SVM Trainings:

- Scikit-learn library used for classifiers and feature extraction.
- TfidfVectorizer optimized via GridSearchCV (5-fold cross-validation).
- Feature count varied by dataset and model (FakeNewsNet: up to 5000/10000, Employee Reviews: up to 10000/5000).

LLM Querying

- vLLM library used for batching requests for HuggingFace models (improved throughput).
- Groq prompts sent sequentially (no batching, rate limits applied).
- GPT4-turbo accessed with batched requests to optimize inference and handle rate limits.

Used Datasets

Use two datasets for binary and multiclass classification:

Table 1: Dataset Details			
Model	Employee Reviews	FakeNewsNet	
Total Records	776 (776)	Fake News (107)	
Classes	Working Regularly (776) (776)	Not Working Regularly (107) (107)	
Avg Length	108 (108)	108 (108)	

FakeNewsNet Dataset:used for a binary classification task, where the model must categorize each article as either "fake" or "real". The challenge lies in the model's ability to detect subtle clues such as sensationalist language, exaggerated claims, or unverifiable facts, which helps to distinguish real news from fake

Employee Reviews Dataset: this dataset consists of 1,000 employee company reviews sourced from a platform where current and former employees provide anonymous feedback on companies and their management.

Evaluation Metrics

Weighted F1-score
Used as the primary metric. It balances precision and recall while accounting for class proportions, crucial due to class imbalance in the Employee Reviews dataset.

For RoBERTa and ML models, k-fold cross-validation is used to mitigate overfitting due to small dataset sizes, and the mean weighted F1-score is reported.

Result Analytics

Model comparison

Prompting Method	Llama3 70B	Llama3 8B	Gemma2	Mistral	Mistral OO	Mistral OH	Zephyr	Xwin	Gpt4-turbo	RoBERTa	Naive Bayes	SVM
1) ZS	92.5	90.2	83.8	84.5	83.1	81.6	76.8	78.3	47.9	81.7		
2) ZS+CoT	91.1	79.4	80.4	80.9	81.6	84.6	84.6	77.8	76.0	81.7		
3) ZS+EM	92.9	90.2	83.8	84.1	84.1	79.1	75.9	79.4	45.9	83.2		
4) ZS+RP+NA	90.4	80.3	82.8	84.1	84.5	78.4	78.0	53.1	48.7	82.2		
5) ZS+CoT+EM	91.1	86.4	82.3	86.0	83.6	85.1	83.0	80.8	76.3	81.7	93.0	88.8
6) ZS+RP	92.9	81.4	83.3	83.2	83.5	77.4	75.4	47.8	46.9	82.2		
7) RP	91.1	81.2	81.2	81.2	81.2	81.2	75.4	47.8	46.9	82.2		
8) FS+RP+NA	92.5	90.0	82.8	83.5	79.7	80.7	71.2	41.1	50.8	83.7		
9) FS+CoT+RP+NA	93.6	92.1	84.2	86.0	81.2	81.2	79.4	50.1	50.1	82.2		
10) FS+CoT+RP+NA	93.6	93.2	81.7	87.8	81.9	81.4	81.1	46.1	69.1	82.8		

- Llama3 70B**: Achieved the highest F1-score of 94.4% using the ZS+RP+NA prompt, with consistently high scores above 91% across all other prompts.
- RoBERTa**: Achieved an excellent F1-score of 93.0% after 7 epochs of training, surpassing all other LLMs and ML models.

- ML Models (NB, SVM)**: Demonstrated competitive F1-scores of 90.0% and 88.8% respectively, outperforming all LLMs except Llama3 70B, 8B, and RoBERTa.

- Gpt4-turbo**: Underperformed with its best score being 83.7%, being surpassed by 5 other models.

- Zephyr, Xwin, Llama2**: Zephyr and Xwin showed mid-80s F1-scores, while Llama2 had the worst overall performance.

- Gpt4-turbo**: Achieved the highest F1-score of 87.6% using the FS+CoT+RP+NA prompting method.

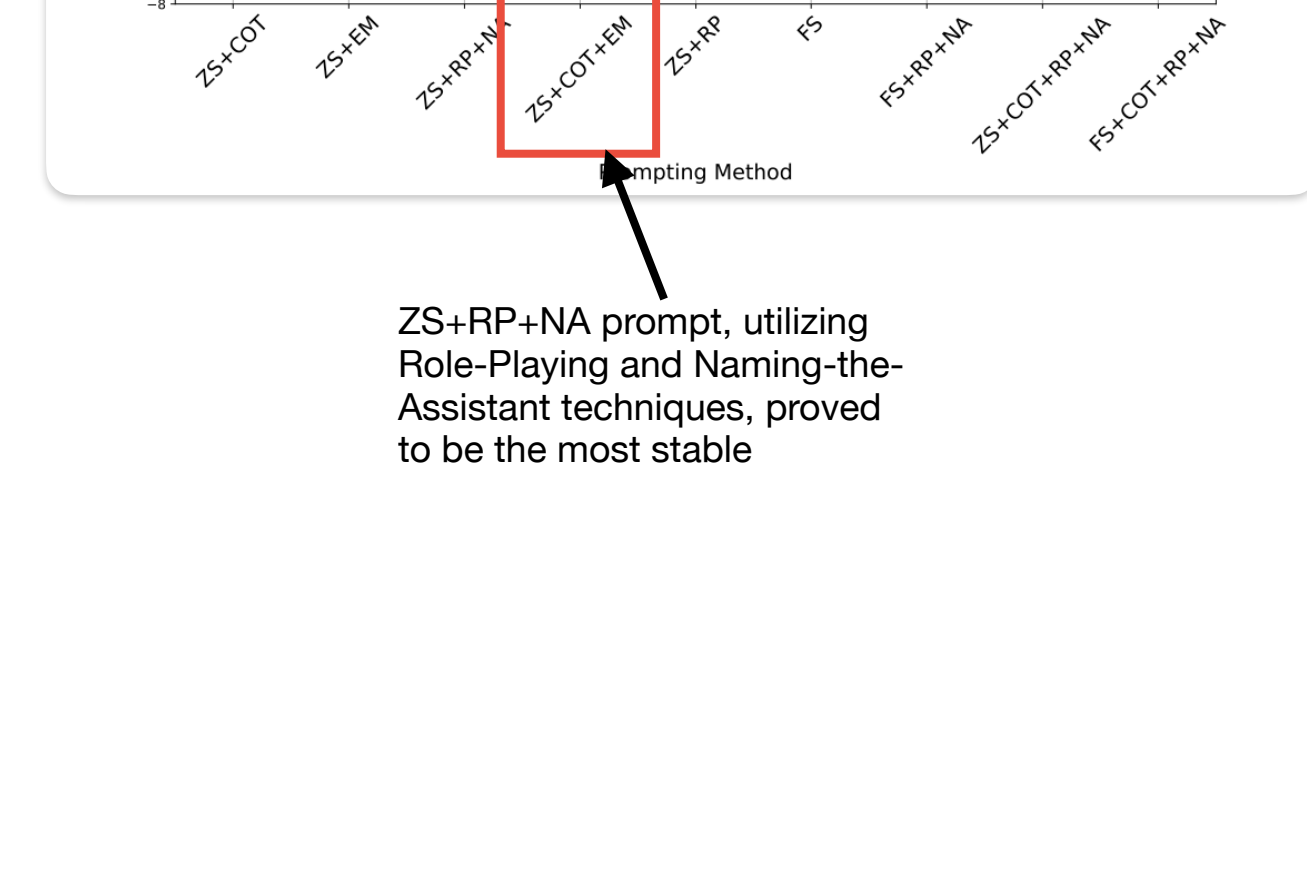
- Llama3 70B and Mistral OO**: Trained Gpt4-turbo by only 0.5% and 1.2% respectively.

- Llama3 8B**: Achieved the second highest score of 84.9% with the FS+CoT+RP+NA prompt.

- RoBERTa**: Maintained a strong performance with an F1-score of 83.8% after 5 epochs of training, though surpassed by 5 LLM models.

- Gemma2**: Showed improved performance compared to the FakeNewsNet task, with its F1-score reaching 84.9%.

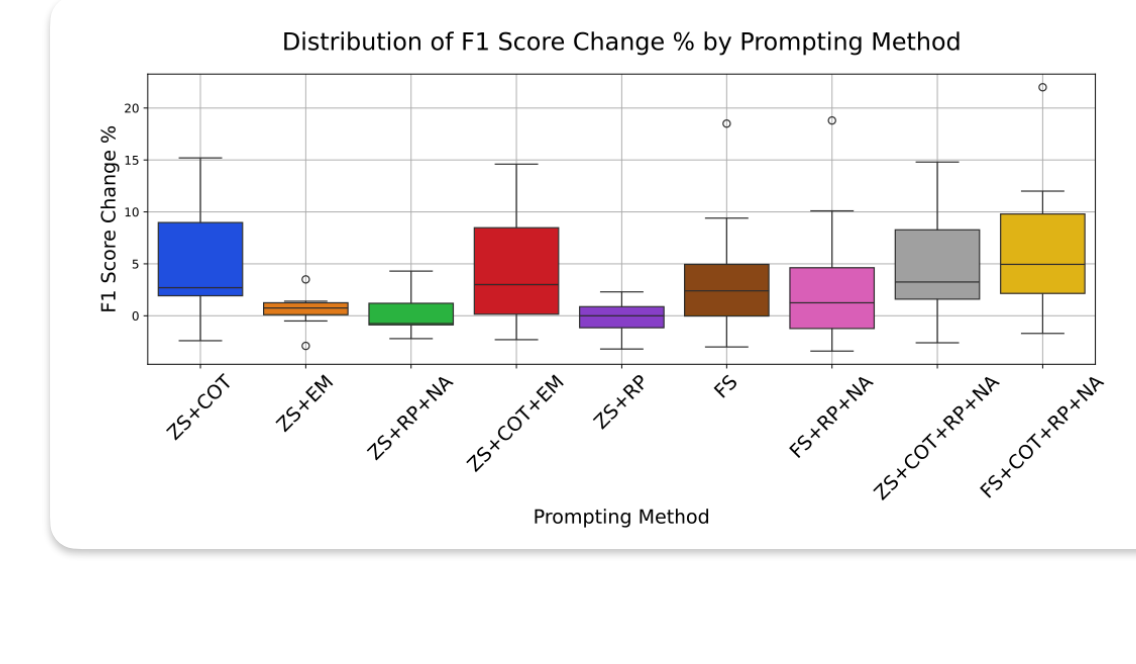
- Zephyr, Llama2, Xwin**: Zephyr showed variable performance, and Llama2 and Xwin particularly struggled on this task.



ZS+RP+NA prompt, utilizing Role-Playing and Naming-the-Assistant techniques, proved to be the most stable

This highlights the differences in model performance across various prompting techniques, indicating that prompts are not **universally effective**. While some models benefit from specific phrasing, others do not, emphasizing the lack of generalizability in prompt design across different LLMs.

Viability in model performance



variability in model performance = despite extensive training, likely haven't encountered every possible phrasing of a task.

The LLMs/Model therefore infer the user's intended task and answer by relying on semantic similarity with their training data.

critical importance of prompt engineering, as poorly constructed prompts can mislead models and degrade performance

correlation between a model's performance and its sensitivity to task phrasing: lower-performing models tend to show greater variability across prompts, while higher-performing models (being more confident in their answers) are less affected by phrasing

Model Scaling and Quantization

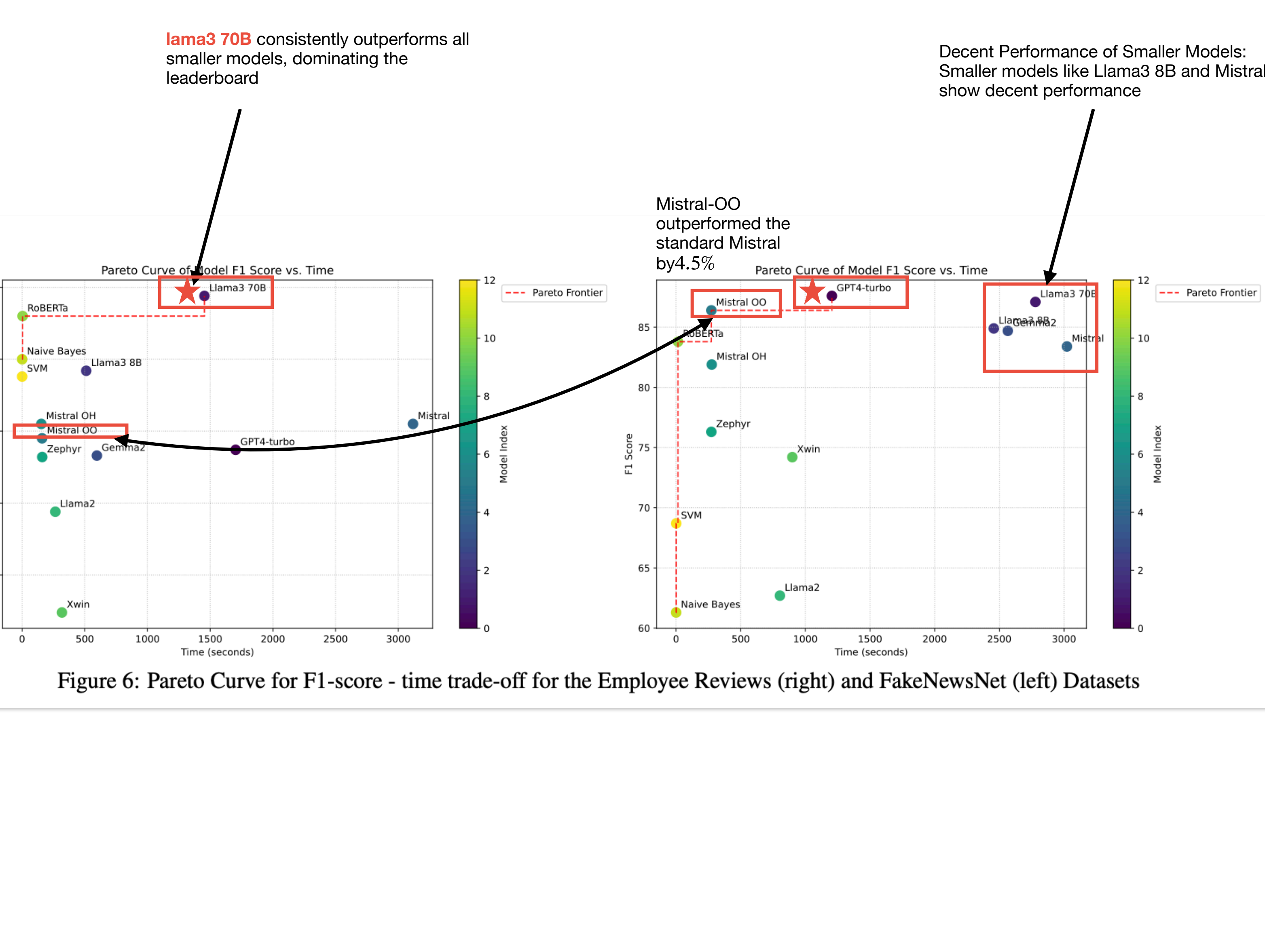


Figure 6: Pareto Curve for F1-score - time trade-off for the Employee Reviews (right) and FakeNewsNet (left) Datasets

My thoughts

- Complexity of Model Selection**: Can not guarantee top performance across all tasks. A model's performance can vary significantly depending on the task type
- Consider Trade-offs**: There's always a trade-off between F1-score (performance) and time/resources (cost -> **select the appropriate model based on your project's priorities (performance vs. speed/cost).**
- Importance of Prompt Engineering**: LLM performance is heavily influenced by prompt design. Specific prompts like "ZS+RP+NA" can lead to stable performance improvements on certain datasets, and techniques like "CoT (Chain-of-Thought)" can enhance performance.

- Importance of Prompt Engineering**: LLM performance is heavily influenced by prompt design. Specific prompts like "ZS+RP+NA" can lead to stable performance improvements on certain datasets, and techniques like "CoT (Chain-of-Thought)" can enhance performance.

Still Strong Contenders: Traditional deep learning models like RoBERTa can offer competitive performance comparable to or even surpassing state-of-the-art large LLMs in certain tasks (e.g., FakeNewsNet binary classification) with better efficiency (faster inference times)

Value of ML Models: Traditional machine learning models like SVM and Naive Bayes can provide decent performance at very high speed



Conclusion

must holistically consider the specific task requirements, available resources, and appropriate prompting strategies to find the optimal model.