# SCUNLP-1 at the NTCIR-18 FinArg-2 Task :
## Collaborative Large Language Models for Temporal Classification

Min-Chin Ho, Jheng-Long Wu
Department of Data Science, Soochow University, Taiwan

## Introduction

In recent years, understanding temporal information within financial texts has become essential for capturing the dynamics of market events and investor behavior. While prior research has extensively studied sentiment and argument mining, the role of temporal reasoning in financial argumentation remains relatively unexplored. This study addresses this gap by leveraging large language models (LLMs) to detect and analyze temporal references in financial arguments. Through novel collaboration mechanisms and tailored prompt designs, we aim to enhance the accuracy and interpretability of temporal predictions, providing deeper insights into the timing and duration of financial impacts.

## Purpose

• Identify temporal references in financial arguments and assess their impact on the expected duration of financial effects.

• Use large language models (LLMs) to perform temporal argument mining in financial texts.

• Apply negotiation-style prompting to guide the model in extracting and reasoning about time-related information.

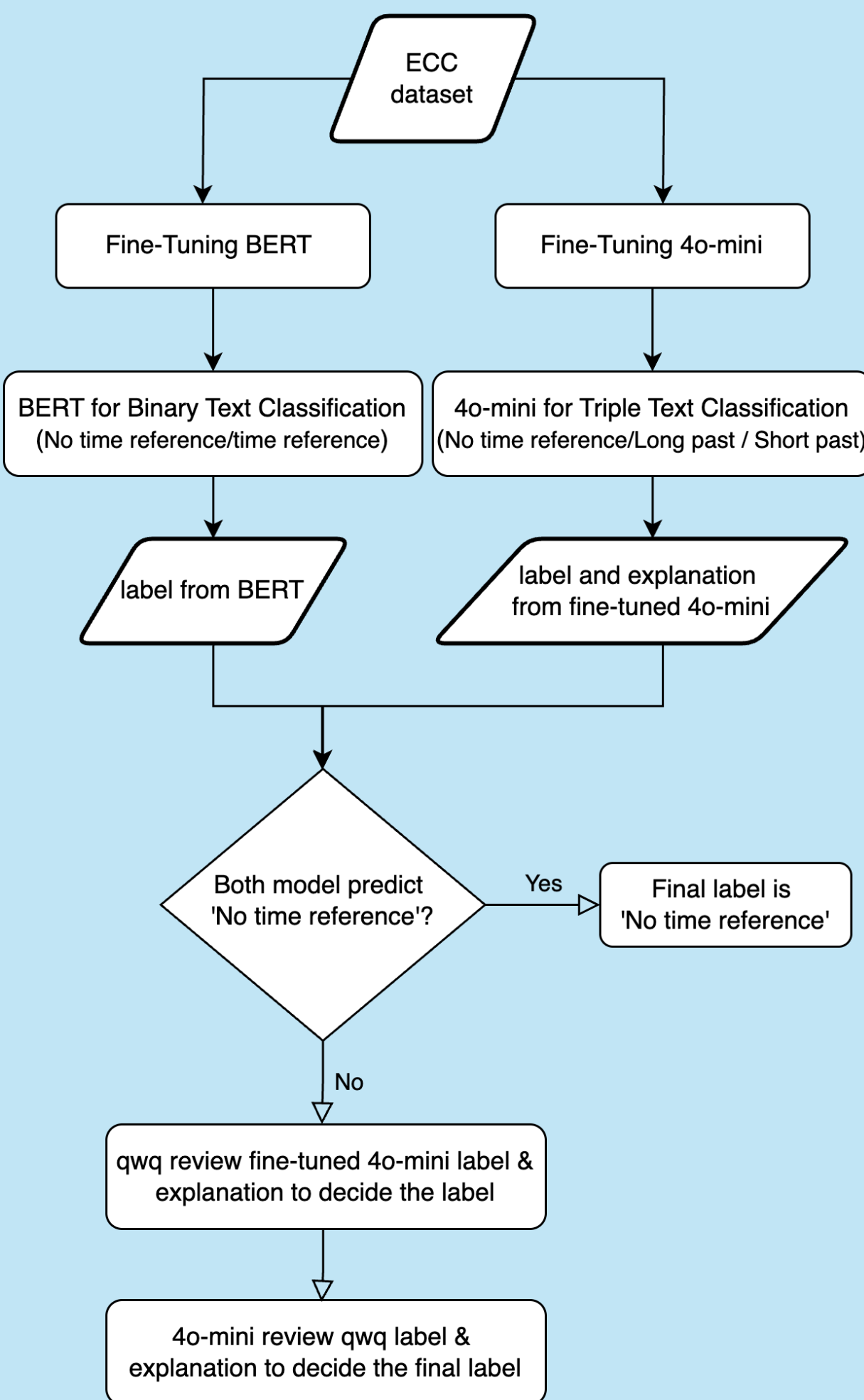• Simulate multi-agent or self-reflective dialogues to uncover and evaluate time-sensitive claims.

## Method

The flow of our study provides a comprehensive explanation of how large language models (LLMs) assign temporal tags to financial texts. First, we introduce the overall decision-making process. Then, we detail the design of custom prompts for each LLM, enabling them to analyze temporal elements and perform self-reflective reasoning.

Our approach adopts a two-stage framework for temporal argument classification. In the first stage, we use a fine-tuned BERT model to perform binary classification, distinguishing between "Time Reference" and "No Time Reference" content. A fine-tuned o4o-mini model is also used for verification. If both models predict "No Time Reference," the instance is finalized with that label. Otherwise, it proceeds to the second stage for further analysis.

In the second stage, we address BERT's limitations in assessing time duration by leveraging the advanced reasoning capabilities of LLMs. Through a negotiation-style framework, multiple LLMs evaluate temporal cues and collaboratively determine the most appropriate temporal label.

To support temporal reasoning, we design a prompt framework that guides LLMs through step-by-step reasoning using Chain of Thought. Each LLM is placed in a multi-turn, debate-style setting, where it reviews the prediction and explanation generated by the previous model. This allows the model to validate or challenge prior conclusions, encouraging deeper reasoning and producing more consistent and interpretable temporal classifications.
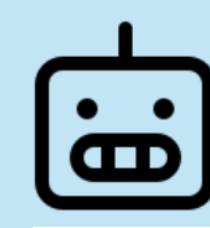


## Submission & Result



• **Official Results**
The official test results reveal that Run 3 achieved the highest overall performance, suggesting that our final prompt design—featuring both step-by-step reasoning and access to prior model outputs—was the most effective. This underscores the importance of prompt clarity and information flow in optimizing LLM-based classification.

| RUN | LLM 1 | LLM 2 | LLM 3 | Final Decision (LLMs + BERT) |
|---|---|---|---|---|
| RUN 1 | 71.43% | 61.9% | 65.48% | 66.67% |
| RUN 2 | 71.43% | 67.86% | 52.38% | 58.33% |
| RUN 3 | 71.43% | 66.67% | 63.1% | 67.86% |

• **Predictions of Each LLM and Final Decision**
The per-model performance shows that LLM1 consistently performed well across all runs, confirming its reliability as the first-stage reasoner. Variability in LLM2 and LLM3 outcomes reflects sensitivity to input structure and the quality of prior explanations. The final decision performance fluctuates depending on how effectively downstream models leverage earlier reasoning.

| RUN | Micro F1 | Macro F1 |
|---|---|---|
| SCUNLP-1_ECC_3 | 67.86% | 64.94% |
| SCUNLP-1_ECC_1 | 66.67% | 63.06% |
| SCUNLP-1_ECC_2 | 58.33% | 52.07% |

• **Prediction Results Across LLM Stages**
Stage-wise breakdown highlights the internal dynamics of the multi-turn reasoning process. While most samples are correctly classified by LLM1, about 28% require further deliberation. Notably, LLM3 occasionally overrides correct predictions from LLM2, indicating that later-stage models may prioritize their own logic over prior answers. This reinforces the need to better align reasoning consistency and trust across models.

| Stage | RUN1 | RUN2 | RUN3 |
|---|---|---|---|
| Correct in LLM1 | 71.43% | 71.43% | 71.43% |
| Transition to LLM2 | 28.57% | 28.57% | 28.57% |
| Transition to LLM3 | 38.10% | 32.14% | 33.33% |
| LLM2 Correct, LLM3 Rejected | 4.76% | 26.20% | 10.71% |

## Conclusion

This study proposes a two-stage temporal labeling framework combining BERT with multi-turn large language model (LLM) collaboration for classifying time-related arguments in financial texts. Our findings demonstrate that step-by-step reasoning and inter-model reference to prior predictions and explanations can enhance both classification accuracy and interpretability.

While there remains room for improvement in model design, reasoning flow, and computational efficiency, our results highlight the potential of collaborative LLMs in tackling complex temporal understanding tasks. Future work may further explore interaction strategies between models, prompt engineering, or variations in model scale to strengthen temporal reasoning performance.