

Auto-Prompting with Retrieval Guidance for Frame Detection in Logistics

Do Minh Duc¹[0009–0008–4116–4067], **Quan Xuan Truong**¹[0009–0001–4911–3394],
Nguyen Tat Dat¹[0009–0005–7422–2663], **Nguyen Van Vinh**^{1,*}[0009–0004–0912–6328]
 dominhduc@vnu.edu.vn, xuantruongvnuet@gmail.com, nguyentatdat2811@gmail.com,
 vinhnv@vnu.edu.vn

¹Faculty of Information Technology, VNU University of Engineering and Technology, Hanoi, Vietnam

Abstract. Prompt engineering plays a critical role in adapting large language models (LLMs) to complex reasoning and labeling tasks without the need for extensive fine-tuning. In this paper, we propose a novel prompt optimization pipeline for frame detection in logistics texts, combining retrieval-augmented generation (RAG), few-shot prompting, chain-of-thought (CoT) reasoning, and automatic CoT synthesis (Auto-CoT) to generate highly effective task-specific prompts. Central to our approach is an LLM-based prompt optimizer agent that iteratively refines the prompts using retrieved examples, performance feedback, and internal self-evaluation. Our framework is evaluated on a real-world logistics text annotation task, where reasoning accuracy and labeling efficiency are critical. Experimental results show that the optimized prompts—particularly those enhanced via Auto-CoT and RAG—improve real-world inference accuracy by up to 15% compared to baseline zero-shot or static prompts. The system demonstrates consistent improvements across multiple LLMs, including GPT-4o, Qwen 2.5 (72B), and LLaMA 3.1 (70B), validating its generalizability and practical value. These findings suggest that structured prompt optimization is a viable alternative to full fine-tuning, offering scalable solutions for deploying LLMs in domain-specific NLP applications such as logistics.

Keywords: Prompt Optimization, Prompt Engineering, LLMs, Frame Detection, Logistics Texts

1 Introduction

The recent breakthroughs in large language models (LLMs), such as ChatGPT [1], Gemini [2], and other foundation models, have significantly reshaped how businesses approach automation. These models demonstrate remarkable abilities in understanding natural language, reasoning over complex instructions, and adapting across diverse use cases without extensive task-specific fine-tuning [3]. As a result, AI-driven automation is being increasingly deployed across industries to streamline workflows, reduce operational costs, and augment human decision-making [4].

¹ * Corresponding author

In the logistics sector, where large volumes of informal text such as shipment updates, incident reports, or customer inquiries must be processed daily, LLMs have the potential to automate previously manual workflows [5]. A key task in this context is **frame detection**, which involves identifying structured semantic events and their roles (e.g., “delayed delivery,” “handover responsibility”) from unstructured short messages. This can be modeled as a multi-level classification task, where each message may carry nested frame labels representing process stage, location, and agent responsibility. However, real-world deployments face two major challenges: (1) extremely low-resource conditions, with only 1–2 labeled examples per frame type [6], and (2) domain-specific terminology that LLMs are not exposed to during pretraining [7].

Building on the prompting paradigm introduced by McCann et al. [8], LLMs are now widely used via prompt-based interfaces. Prompting has been enhanced through few-shot learning [9], instruction tuning [10], and zero-shot reasoning [11]. However, prompting remains sensitive: small variations in wording, ordering, or formatting can result in large differences in output quality [12], making it difficult for non-expert users to construct reliable prompts.

To overcome these limitations, **Automatic Prompt Optimization (APO)** methods have been proposed to refine prompts in a black-box setting, without requiring model fine-tuning. Although effective in general domains, APO remains underexplored in domain-specific, low-resource contexts such as logistics.

In this work, we introduce **Auto-Prompting with Retrieval Guidance**, a retrieval-enhanced APO approach tailored to frame detection in logistics. Our method:

- Maintains a library of prompt components, including structured examples and instruction templates;
- Uses semantic similarity retrieval to dynamically select relevant components per input;
- Composes optimized prompts without human intervention;
- Achieves more than 80% accuracy on a proprietary logistics dataset, outperforming zero-shot and manually engineered prompts.

In summary, by bridging APO with retrieval-based contextualization [13], our method reduces prompt engineering burden, improves robustness, and facilitates practical deployment of LLMs for structured frame detection in specialized domains. Conceptually, the framework can be viewed as an iterative search in the prompt space, where retrieval narrows contextual uncertainty and self-evaluation acts as a feedback signal approximating gradient-free optimization. This provides a theoretical justification for its effectiveness in black-box settings.

2 Related Work

2.1 LLMs in Domain-Specific and Low-Resource Scenarios

The unprecedented capabilities of LLMs like GPT-4 [1] have spurred their adoption in various specialized, high-stakes domains that often suffer from data scarcity. In medicine, models are being developed to assist with clinical note summarization and diagnostic

reasoning [14]. Similarly, in the legal sector, LLMs are being explored for contract analysis and legal question-answering, where domain-specific terminology and nuanced reasoning are critical. The financial domain has also seen the emergence of models like FinGPT, trained on financial data to handle tasks such as sentiment analysis and report generation [15].

The logistics sector presents similar challenges: a reliance on proprietary data formats, a specialized vocabulary, and the need to process vast streams of unstructured text. Recent surveys suggest that LLMs hold significant potential for automating logistics workflows, such as shipment tracking and incident reporting[12]. Our work addresses a core task in this domain: frame detection and many labels we only have 1 or 2 samples.

2.2 Automatic Prompt Optimization

The effectiveness of LLMs is heavily dependent on the quality of the input prompt, a phenomenon that has given rise to the field of "prompt engineering." However, manual prompt crafting is often brittle and requires significant expertise [16,17]. To address this, Automatic Prompt Optimization (APO) has emerged as a key research direction.

APO techniques can be broadly categorized. Gradient-based methods, such as prompt tuning [18], learn continuous "soft prompts" that are prepended to the input embedding. However, these methods typically require a moderate amount of labeled data for effective training. In our case, the dataset is extremely low-resource, with only 1–2 labeled examples per label, making conventional fine-tuning or even prompt tuning impractical. Given the large number of distinct labels and the limited supervision available, we instead leverage the generative capacity of LLMs themselves to optimize prompts in a data-efficient manner. A seminal work in this area, Automatic Prompt Engineer (APE) [19], uses an LLM to generate and search for instruction candidates. More recent approaches have framed this as a reinforcement learning problem. For instance, Pryzant et al [20] employ "gradient-free" exploration, using an editor model to propose prompt edits and a reward model to evaluate their efficacy. These methods are particularly appealing for low-resource settings, as they minimize dependence on labeled examples while still enabling prompt adaptation to task-specific semantics.

A critical component of prompting, especially for complex reasoning tasks, is the Chain-of-Thought (CoT) [21]. CoT prompting, which encourages the model to generate intermediate reasoning steps, has been shown to significantly improve performance. Research has also focused on automating the generation of these reasoning chains. Auto-CoT [22] clusters questions and generates reasoning chains for each cluster, creating a diverse set of exemplars for few-shot prompting. Shum proposed generating multiple reasoning paths and then selecting the most consistent one to improve robustness[23]. Our approach incorporates reward-based learning to optimize prompts that may include CoT-style instructions, but it uniquely integrates a retrieval mechanism to dynamically select the few-shot examples that form the prompt's core.

2.3 Retrieval-Augmented Language Models

LLMs possess vast parametric knowledge but it is static and may lack domain-specific or real-time information. Retrieval-Augmented Generation (RAG) was introduced to

mitigate this by augmenting prompts with relevant information retrieved from an external knowledge corpus [13]. This paradigm has proven effective at reducing factual hallucinations and improving performance on knowledge-intensive tasks.

The synergy between retrieval and in-context learning is particularly relevant to our work. Instead of retrieving factual passages, recent methods retrieve entire exemplars (i.e., input-output pairs) to populate the few-shot context of a prompt. For instance, Ram demonstrated that retrieving examples that are semantically similar to the test input significantly improves few-shot performance[24]. This dynamic selection of examples is more effective than using a fixed, randomly chosen set. These systems typically employ dense retrieval models, such as Sentence-BERT [25], to embed the corpus and query into a shared vector space for similarity search.

Our method, Auto-Prompting with Retrieval Guidance, builds directly on this line of research. However, we extend the paradigm in a crucial way: retrieval is not just a pre-processing step to find examples, but an integral part of the automatic prompt construction loop. We maintain a structured library of prompt components—including instruction templates and domain-specific few-shot examples—and use dense retrieval to dynamically select the most relevant components for each input instance. This allows our system to compose highly specialized prompts tailored to the nuances of hierarchical frame detection in logistics, effectively bridging the gap between automatic prompt optimization and retrieval-augmented in-context learning in a challenging, low-resource industrial setting.

3 Dataset

We construct a domain-specific dataset comprising 1,500 text messages collected from a Vietnamese logistics platform. Each instance is a short sentence extracted from real-world customer or internal staff communications. The messages are informal, highly domain-specific, and often include abbreviations, shorthand expressions, and non-standard language structures—reflecting the linguistic challenges typically encountered in operational logistics settings.

Each message is annotated with a three-level hierarchical frame label, capturing the semantic structure of the logistics issue being described:

- Level 1 (actor): the primary agent or entity responsible for the issue (e.g., Customer, Shop, Delivery Service, External Factors),
- Level 2 (reason): the underlying reason or category of the event (e.g., Incorrect Information, Unavailable, Changed Address, Refused Delivery),
- Level 3 (fine-grained cause): a specific instantiation of the reason (e.g., Wrong Product, Out of Stock, Customer on Vacation).

The annotation schema includes 73 unique frame labels in total, representing the full cross-product of these three hierarchical levels. The full label space is defined in a structured format such as Shop – Thay đổi thông tin – Thời gian lấy hàng, which translates to (Actor=Shop, Reason=Change Info, Detail=Pickup Time).

Annotation was performed manually by domain experts within the logistics company who possess extensive knowledge of business operations and customer communication

patterns. All annotations were reviewed and agreed upon through internal consensus to ensure label consistency and semantic fidelity. Notably, the dataset reflects real-world class imbalance: while some frame labels are associated with dozens of instances, others appear only once or twice, creating a naturally low-resource setting for many categories.

This dataset serves as a valuable resource for developing and evaluating frame detection methods under conditions that are both semantically complex and data-scarce, highlighting the practical challenges of applying large language models in real-world industrial domains.

3.1 Dataset Statistics and Characteristics

To better understand the linguistic characteristics of the dataset, we conducted descriptive statistical analysis on the first 1,500 annotated messages. The results are summarized in Table 1, while Figure 1 and Figure 2 provide visual insights.

Table 1. Sentence-level statistics of the annotated dataset

| Statistic | Value |
|--------------------------------|------------|
| Total number of sentences | 1,500 |
| Average sentence length | 8.48 words |
| Maximum sentence length | 49 words |
| Minimum sentence length | 1 word |
| Sentences longer than 10 words | 393 |
| Training set (70%) | 1,050 |
| Validation set (15%) | 225 |
| Test set (15%) | 225 |



Fig. 1. Word cloud of the most frequent tokens in the dataset

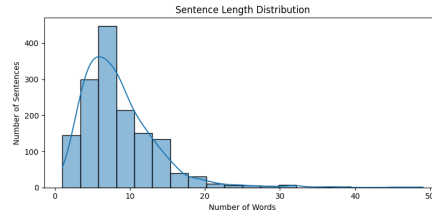


Fig. 2. Distribution of sentence lengths (in number of words)

From the analysis, we observe that the dataset contains mostly short utterances, with a median around 7–9 words. The vocabulary includes many domain-specific terms such as *shop*, *khách*, *giao*, and informal tokens like *ck*, *kt*, or *sđt*, which reflect both abbreviation and inconsistency in natural language usage. These properties further

emphasize the challenges in applying structured learning or fine-tuning methods to the data, and highlight the importance of robust prompting and retrieval mechanisms.

4 Method

The loop for generating optimized prompts is repeated three times. Each times we have three prompts using for compare and choose the best prompt.

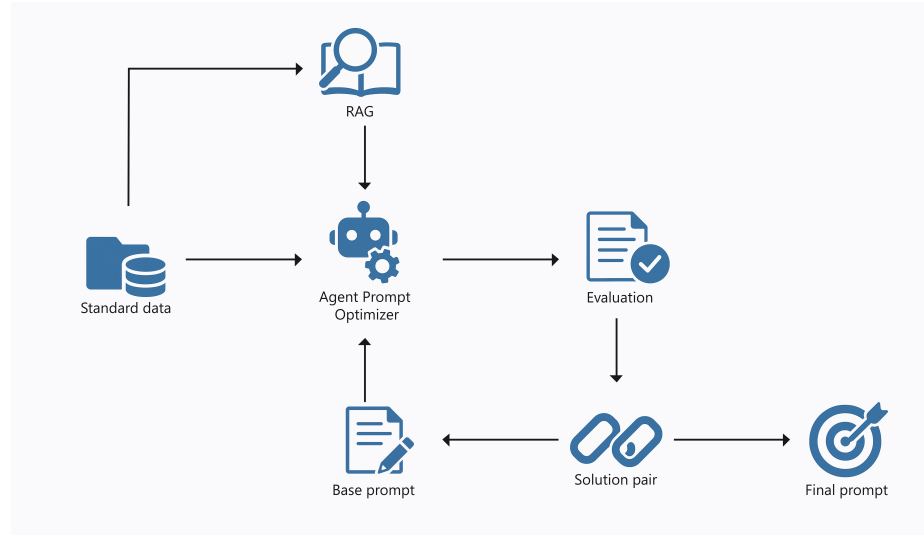


Fig. 3. Overview of the proposed prompt optimization pipeline.

To enhance the performance of large language models (LLMs) on text annotation tasks, we propose a prompt optimization pipeline that integrates retrieval-augmented generation (RAG) with an LLM-based prompt optimizer. Our method aims to generate effective task-specific prompts that generalize well during inference. Figure 3 illustrates the architecture of our system.

4.1 Retrieval-Augmented Generation (RAG)

To enrich the prompt context with relevant examples, we utilize a dense retrieval mechanism. Each input query is used to retrieve semantically similar examples from a pre-encoded vector database. This process enables the system to dynamically augment the prompt with high-quality, context-relevant exemplars, improving the model’s reasoning ability and robustness.

4.2 Prompt Optimizer Agent

At the core of our system is the *Prompt Optimizer*, an autonomous agent designed to iteratively improve prompts using a suite of large language models (LLMs). Specifically, we employ a combination of state-of-the-art models, including **LLaMA 3.1 70B**, **Qwen 2.5 72B-chat**, **GPT-4o mini**, and **GPT-4o**. These models serve both as prompt generators and evaluators, enabling a powerful closed-loop prompt engineering framework.

The optimizer receives as input a *base prompt* and a set of retrieved examples obtained via the RAG module. It applies multiple prompt engineering techniques to generate candidate prompts optimized for the downstream task. Among the key strategies used are:

- **Few-shot prompting:** Retrieved examples are appended to the prompt as demonstrations, allowing the model to learn from specific task instances and better generalize to unseen inputs.
- **Chain-of-Thought (CoT) prompting:** Prompts are structured to elicit intermediate reasoning steps before final predictions. This is particularly effective for tasks involving multi-step logic or weak supervision.
- **Auto-CoT:** Building upon CoT, we apply Auto-CoT methods to automatically generate rationales from unlabeled data using a two-phase process: (1) reasoning generation with self-consistency sampling, and (2) rationale-filtering based on agreement with ground-truth or simulated outputs. This allows us to bootstrap high-quality CoT-style demonstrations without requiring costly manual annotation.

The prompt refinement process operates in a feedback loop. For each generated candidate prompt, we simulate outputs using one or more of the aforementioned LLMs and assess performance using a mix of automated metrics and internal scoring functions. These functions may include comparisons with expert annotations or meta-evaluation via a scoring model (e.g., GPT-4o). Poorly performing prompts are discarded, while promising ones are further refined through transformations such as reordering examples, adjusting instruction specificity, or modifying CoT strategy. Through this iterative optimization process, the Prompt Optimizer acts as a meta-level agent capable of self-evaluating and improving its prompt generation policy. The final prompt selected from this pool is the one that consistently yields the best annotation accuracy and reasoning quality across multiple validation settings. This prompt is then deployed during inference, forming the backbone of the downstream system.

4.3 Evaluation and Selection

Each candidate prompt generated by the optimizer is evaluated on a held-out validation set. Evaluation metrics include both automatic measures (e.g., accuracy) and human-in-the-loop assessments performed by domain experts. Prompts are compared in terms of their effectiveness in producing correct annotations.

4.4 Solution Pairing and Final Selection

The evaluation phase yields a set of candidate *solution pairs*, where each pair consists of a prompt and its associated performance score. These pairs are ranked and analyzed

to determine the best-performing prompt. The selected *Final Prompt* is then used for inference in downstream applications.

The final prompt is deployed for real-world inference scenarios, enabling the LLM to generate accurate and contextually appropriate annotations for new inputs. This prompt encapsulates both task-specific requirements and contextual examples retrieved dynamically during the optimization process.

5 Experiment

The main configuration parameters used across all evaluated models are summarized in Table 2. To ensure consistent reasoning behavior and fair comparison, we fixed the sampling and decoding parameters during all experiments. Specifically, the temperature was set to 0.3 to reduce randomness in generation and encourage deterministic reasoning chains. The top-p (nucleus sampling) and top-k values were set to 0.95 and 70, respectively, balancing diversity and precision in token selection.

Table 2. Model configurations and reasoning capabilities

| Parameter | LLaMA 3.1 70B | Qwen 2.5 72B | GPT-4o mini | GPT-4o |
|--------------------|---------------|--------------|-------------|--------|
| top_p | 0.95 | 0.95 | 0.9 | 0.9 |
| top_k | 70 | 70 | 50 | 50 |
| temperature | 0.3 | 0.3 | 0.3 | 0.3 |
| max tokens | 1024 | 1024 | 1024 | 1024 |
| repetition penalty | 0 | 0 | 0 | 0 |
| presence penalty | 0 | 0 | 0 | 0 |
| reasoning | Yes | Yes | Yes | Yes |

We categorize prompting strategies as follows:

- **Manual Prompt (6-shot)**: handcrafted prompt and 6 manually selected examples.
- **Auto Prompt (0-shot)**: prompt automatically generated by the optimizer, no examples.
- **Auto Prompt + RAG (k-shot)**: optimizer-generated prompt with k retrieved in-context examples via semantic retrieval.

Table 4 summarizes performance across several LLMs under varying prompting strategies. We evaluate zero-shot, manual few-shot, and our proposed auto-prompting with retrieval-augmented few-shot examples. The goal is to assess how models adapt with optimized prompts selected dynamically based on input similarity.

Auto Prompt + RAG is Most Effective. Across all models, the best results are consistently obtained using our proposed 6-shot Auto Prompt + RAG strategy. For example, **GPT-4o** reaches **90%** test and **92%** real-world accuracy, outperforming both manual and zero-shot prompts. This supports the effectiveness of combining LLM-based prompt optimization with retrieval-guided example selection.

Table 3. Comparison of fine-tuning strategies (standard vs. soft prompt) on the test set.

| Model | Fine-tuning Strategy | Test Accuracy |
|----------------|----------------------|---------------|
| DeBERTa-large | Standard | 40% |
| ViDeBERTa-base | Standard | 45% |
| mDeBERTa | Standard | 61% |
| DeBERTa-large | Soft Prompt | 37% |
| ViDeBERTa-base | Soft Prompt | 40% |
| mDeBERTa | Soft Prompt | 57% |

Table 4. Performance of different prompting strategies across models.

| Model | Prompting Strategy | Test Accuracy | Real-World Accuracy |
|----------------------|-----------------------------------|---------------|---------------------|
| GPT-4o mini | Manual Prompt (6-shot) | 83% | 83% |
| GPT-4o mini | Auto Prompt (0-shot) | 82% | 80% |
| GPT-4o mini | Auto Prompt + RAG (3-shot) | 85% | 84% |
| GPT-4o mini | Auto Prompt + RAG (6-shot) | 88% | 88% |
| GPT-4o | Manual Prompt (6-shot) | 84% | 84% |
| GPT-4o | Auto Prompt (0-shot) | 84% | 84% |
| GPT-4o | Auto Prompt + RAG (3-shot) | 86% | 88% |
| GPT-4o | Auto Prompt + RAG (6-shot) | 90% | 92% |
| LLaMA 3.1 70B | Manual Prompt (6-shot) | 72% | 70% |
| LLaMA 3.1 70B | Auto Prompt (0-shot) | 74% | 80% |
| LLaMA 3.1 70B | Auto Prompt + RAG (3-shot) | 79% | 81% |
| LLaMA 3.1 70B | Auto Prompt + RAG (6-shot) | 87% | 87% |
| Qwen 2.5 72B | Manual Prompt (6-shot) | 73% | 71% |
| Qwen 2.5 72B | Auto Prompt (0-shot) | 76% | 80% |
| Qwen 2.5 72B | Auto Prompt + RAG (3-shot) | 79% | 81% |
| Qwen 2.5 72B | Auto Prompt + RAG (6-shot) | 87% | 87% |

Manual vs. Retrieval-based Prompting. Manual 6-shot prompts generally underperform. For instance, GPT-4o mini improves from 83% (manual) to 88% (RAG-6-shot), and similar gains are seen with Qwen and LLaMA. This shows that automatic retrieval often selects more relevant demonstrations than handcrafted examples.

Model Trends. GPT-4o and its mini variant show strong, consistent performance. In contrast, Qwen and LLaMA start lower in zero-shot (74–76%) but improve markedly with retrieval-augmented prompts (up to 87%). This highlights the value of prompt quality over raw model size.

Real-World Generalization. High test accuracy correlates well with real-world performance. Models like GPT-4o and Qwen retain or improve accuracy post-deployment. Slight drops in LLaMA suggest generalization may depend on alignment and robustness to input variation.

Fine-tuning Baselines. Table 3 compares standard and soft prompt fine-tuning on encoder-based models. Even the best result (57% for mDeBERTa with soft prompt) lags

far behind prompting-based LLMs. This reinforces the strength of in-context learning with optimized prompting pipelines over traditional fine-tuning approaches.

Conclusion. Our findings show that combining automatic prompt generation with retrieval-augmented few-shot examples offers a scalable and effective alternative to fine-tuning. It improves both performance and generalization, while being model-agnostic and parameter-efficient.

6 Limitations

While our method achieves strong performance in both controlled and real-world settings, it has several notable limitations:

- **Domain specificity:** Evaluation is confined to the **logistics domain**, which, while enabling precise alignment, limits **generalizability** to other domains or languages.
- **Evaluation subjectivity:** Real-world performance is assessed by **domain experts**, not NLP researchers, introducing **subjectivity** and lacking formal metrics for reasoning depth or faithfulness.
- **RAG sensitivity to informal input:** The retrieval component struggles with **slang, shorthand, or non-standard phrasing**, leading to **suboptimal few-shot selection**.
- **Fixed iteration prompting:** The prompt optimizer performs a **fixed number of refinement steps**, which may **fail to converge** and adds **latency**, limiting **real-time or large-scale applicability**.
- **Compute requirements:** High-performing models (e.g., **GPT-4o, Qwen 72B**) require substantial **compute resources**, which is impractical for **low-resource or edge deployment** scenarios.
- **Limited label diversity:** Training data often includes only **1–2 examples per class**, which reduces the system’s ability to select **representative demonstrations**, limiting **few-shot generalization** even with RAG support.

References

1. O. J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman *et al.*, “Gpt-4 technical report,” 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:257532815>
2. G. Team, P. Georgiev, V. I. Lei, R. Burnell, L. Bai, A. Gulati *et al.*, “Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context,” *ArXiv*, vol. abs/2403.05530, 2024. [Online]. Available: <https://api.semanticscholar.org/CorpusID:268297180>
3. E. M. Bender, T. Gebru, A. McMillan-Major, and S. Shmitchell, “On the dangers of stochastic parrots: Can language models be too big?” *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, 2021. [Online]. Available: <https://api.semanticscholar.org/CorpusID:262580630>
4. H. Wang, J. Jiang, L. J. Hong, and G. Jiang, “Llms for supply chain management,” *ArXiv*, vol. abs/2505.18597, 2025. [Online]. Available: <https://api.semanticscholar.org/CorpusID:278904366>

5. S. S. Jamee, M. R. Hossain, M. Hasan, M. K. Sharif, M. S. Khan, M. I. Islam *et al.*, “Enhancing supply chain decision-making with large language models: A comparative study of ai-driven optimization,” *International Journal of Economics Finance & Management Science*, 2025. [Online]. Available: <https://api.semanticscholar.org/CorpusID:277639822>
6. Z. Jiang, M. Y. R. Yang, M. Tsirlin, R. Tang, Y. Dai, and J. J. Lin, ““low-resource” text classification: A parameter-free classification method with compressors,” in *Annual Meeting of the Association for Computational Linguistics*, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:259858947>
7. Y. Gu, R. Tinn, H. Cheng, M. R. Lucas, N. Usuyama, X. Liu *et al.*, “Domain-specific language model pretraining for biomedical natural language processing,” *ACM Transactions on Computing for Healthcare (HEALTH)*, vol. 3, pp. 1 – 23, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:220919723>
8. B. McCann, N. S. Keskar, C. Xiong, and R. Socher, “The natural language decathlon: Multitask learning as question answering,” *ArXiv*, vol. abs/1806.08730, 2018. [Online]. Available: <https://api.semanticscholar.org/CorpusID:49393754>
9. T. B. Brown, B. Mann, N. Ryder, M. Subbiah, J. Kaplan, P. Dhariwal *et al.*, “Language models are few-shot learners,” *ArXiv*, vol. abs/2005.14165, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:218971783>
10. L. Ouyang, J. Wu, X. Jiang, D. Almeida, C. L. Wainwright, P. Mishkin *et al.*, “Training language models to follow instructions with human feedback,” *ArXiv*, vol. abs/2203.02155, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:246426909>
11. T. Kojima, S. S. Gu, M. Reid, Y. Matsuo, and Y. Iwasawa, “Large language models are zero-shot reasoners,” *ArXiv*, vol. abs/2205.11916, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:249017743>
12. L. Wang, C. Ma, X. Feng, Z. Zhang, H. ran Yang, J. Zhang *et al.*, “A survey on large language model based autonomous agents,” *Frontiers Comput. Sci.*, vol. 18, p. 186345, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:261064713>
13. P. Lewis, E. Perez, A. Piktus, F. Petroni, V. Karpukhin, N. Goyal *et al.*, “Retrieval-augmented generation for knowledge-intensive nlp tasks,” *ArXiv*, vol. abs/2005.11401, 2020. [Online]. Available: <https://api.semanticscholar.org/CorpusID:218869575>
14. K. Singhal, S. Azizi, T. Tu, S. Mahdavi, J. Wei, H. W. Chung *et al.*, “Large language models encode clinical knowledge,” *Nature*, vol. 620, pp. 172 – 180, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:255124952>
15. H. Yang, X.-Y. Liu, and C. Wang, “Fingpt: Open-source financial large language models,” *ArXiv*, vol. abs/2306.06031, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:259129734>
16. Z. Jiang, F. F. Xu, J. Araki, and G. Neubig, “How can we know what language models know?” *Transactions of the Association for Computational Linguistics*, vol. 8, pp. 423–438, 2019. [Online]. Available: <https://api.semanticscholar.org/CorpusID:208513249>
17. T. Zhao, E. Wallace, S. Feng, D. Klein, and S. Singh, “Calibrate before use: Improving few-shot performance of language models,” in *International Conference on Machine Learning*, 2021. [Online]. Available: <https://api.semanticscholar.org/CorpusID:231979430>
18. B. Lester, R. Al-Rfou, and N. Constant, “The power of scale for parameter-efficient prompt tuning,” in *Conference on Empirical Methods in Natural Language Processing*, 2021. [Online]. Available: <https://api.semanticscholar.org/CorpusID:233296808>
19. Y. Zhou, A. I. Muresanu, Z. Han, K. Paster, S. Pitis, H. Chan *et al.*, “Large language models are human-level prompt engineers,” *ArXiv*, vol. abs/2211.01910, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:253265328>
20. R. Pryzant, D. Iter, J. Li, Y. T. Lee, C. Zhu, and M. Zeng, “Automatic prompt optimization with “gradient descent” and beam search,” in *Conference on Empirical Methods in Natural*

- Language Processing*, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:258546785>
21. J. Wei, X. Wang, D. Schuurmans, M. Bosma, E. H. Chi, F. Xia *et al.*, “Chain of thought prompting elicits reasoning in large language models,” *ArXiv*, vol. abs/2201.11903, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:246411621>
 22. Z. Zhang, A. Zhang, M. Li, and A. J. Smola, “Automatic chain of thought prompting in large language models,” *ArXiv*, vol. abs/2210.03493, 2022. [Online]. Available: <https://api.semanticscholar.org/CorpusID:252762275>
 23. K. Shum, S. Diao, and T. Zhang, “Automatic prompt augmentation and selection with chain-of-thought from labeled data,” in *Conference on Empirical Methods in Natural Language Processing*, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:257205763>
 24. O. Ram, Y. Levine, I. Dalmedigos, D. Muhlgay, A. Shashua, K. Leyton-Brown *et al.*, “In-context retrieval-augmented language models,” *Transactions of the Association for Computational Linguistics*, vol. 11, pp. 1316–1331, 2023. [Online]. Available: <https://api.semanticscholar.org/CorpusID:256459451>
 25. N. Reimers and I. Gurevych, “Sentence-bert: Sentence embeddings using siamese bert-networks,” in *Conference on Empirical Methods in Natural Language Processing*, 2019. [Online]. Available: <https://api.semanticscholar.org/CorpusID:201646309>

A Example Prompt

In this appendix, we present the prompt templates that were employed in the logistics frame detection task. These prompts were designed to guide large language models (LLMs) such as GPT-4o, Qwen 2.5, and LLaMA 3.1 in generating preliminary annotations. We experimented with zero-shot, few-shot, chain-of-thought, automatic chain-of-thought (Auto-CoT), and retrieval-augmented prompts.

A.1 Zero-shot Prompt

Zero-shot Frame Detection Prompt

User: Identify the frame structure in the following logistics message. The frame label has three levels:

- Actor (Customer, Shop, Delivery Service, External Factor)
- Reason (Incorrect Information, Unavailable, Changed Address, Refused Delivery, ...)
- Fine-grained Cause (e.g., Wrong Product, Out of Stock, Customer on Vacation)

Text: “<input logistics message>”

Assistant: { "actor": ..., "reason": ..., "cause": ... }

A.2 Few-shot Prompt

Few-shot Frame Detection Prompt

User: Here are examples of annotated logistics messages. Please follow the same format for the new message.

Example 1: Text: “Khách không nhận vì giao trễ.” ” Output:
("Customer", "Refused Delivery", "Late Delivery")

Example 2: Text: “Shop báo sai địa chỉ nhận.” ” Output:
("Shop", "Incorrect Information", "Wrong Address")

Now annotate the following text: “<input logistics message>”

Assistant: ("...", "...", "...")

A.3 Chain-of-Thought Prompt**Chain-of-Thought Frame Detection Prompt**

User: Analyze the logistics message step by step.

1. Identify the main actor.
2. Identify the reason category.
3. Determine the fine-grained cause.

Finally, return the structured frame label.

Text: “<input logistics message>”

Assistant: Step 1: Actor = ... Step 2: Reason = ... Step 3: Cause = ... Final
Output: { "actor": ..., "reason": ..., "cause": ... }

A.4 Auto-CoT Prompt**Automatic Chain-of-Thought Prompt**

User: “<input logistics message>”

Assistant: Reasoning Path A: ... → Candidate frame label. Reasoning Path B:
... → Candidate frame label. Reasoning Path C: ... → Candidate frame label.

User: Select the most consistent reasoning path and output the final annotation.

Assistant: { "actor": ..., "reason": ..., "cause": ... }

A.5 Retrieval-Augmented Prompt**RAG-Enhanced Few-shot Prompt**

User: Given the following retrieved examples from similar logistics cases, use them to annotate the new message.

Retrieved Example 1: Text: “Khách đi vắng nên không nhận hàng.” ” Output: ("Customer", "Unavailable", "On Vacation")

Retrieved Example 2: Text: “Giao hàng thất bại do thiếu số điện thoại.” ” Output: ("Delivery Service", "Incorrect Information", "Missing Contact Info")

New text: "<input logistics message>"
Assistant: ("...", "...", "...")

A.6 Evaluation-Refinement Prompt

Evaluation and Refinement Prompt

User: You previously used the following prompt for frame detection: <previous prompt text>

Here are some annotated cases where the model prediction was incorrect:

Case 1: ...

Case 2: ...

Assistant: Step 1: Analyze the error in each case and explain why the prediction was wrong. Step 2: Identify missing or ambiguous instructions in the current prompt. Step 3: Propose modifications to improve the prompt. Step 4: Output a refined version of the prompt.

Assistant (Example Reasoning): - Case 1: The model confused "Out of Stock" with "Wrong Product" because the prompt did not emphasize distinguishing product mismatch from stock issues. - Case 2: The model misclassified the actor. The prompt should explicitly instruct to check who is responsible (Customer vs. Delivery Service).

Refined Prompt: "When detecting frames, always check

- (1) who is responsible for the issue (actor),
- (2) whether the problem comes from product mismatch vs. stock shortage
- (3) if the delivery service failed due to missing information. Return output as actor, reason, cause."

A.7 Multi-Agent Debate and Self-Improvement Prompt

Multi-Agent Debate Prompt

User: You are Agent A. Given the following logistics message, provide your frame detection output: Text: ... **Assistant (Agent A):** ("Customer", "Refused Delivery", "Late Delivery") **User:** Now here are the outputs from other agents: - Agent B: ("Delivery Service", "Unavailable", "Late Delivery") - Agent C: ("Customer", "Changed Address", "Wrong Time") Using these alternative answers, explain step by step which reasoning is more consistent, identify weaknesses in your own answer, and propose an improved solution.

Assistant (Agent A Revised): Step 1: Compare outputs: Agent B incorrectly assigns the actor; Agent C mislabels the reason. Step 2: My original output is mostly correct, but I need to emphasize actor detection rules. Step 3: Improved answer: ("Customer", "Refused Delivery", "Late Delivery")

Self-Improvement Debate Prompt

User: Multiple agents have generated candidate prompts for frame detection. Here are three prompt variants:

Prompt A: “Identify actor, reason, and cause. Return as JSON.” Prompt B: “Step by step, decide actor first, then reason, then cause.” Prompt C: “Focus on who is responsible, why it happened, and the fine-grained detail.”

Debate with each other to decide which prompt is most effective. Consider coverage, clarity, and ability to prevent errors seen in past cases. After the debate, propose a new improved prompt that combines the strengths of A, B, and C.

Assistant (Consensus): - Agent 1: Prompt A is too general, may cause confusion. - Agent 2: Prompt B enforces reasoning steps, but might be verbose. - Agent 3: Prompt C adds clarity on actor responsibility.

Final Improved Prompt: “When analyzing a logistics message: 1. Identify who is responsible (actor). 2. Identify the general reason category. 3. Identify the fine-grained cause. Always distinguish product mismatch vs. stock issues, and customer refusal vs. delivery service error. Return output in JSON format as "actor": ..., "reason": ..., "cause":”