# Hybrid Entity-Specific Compression and Multi-Agent Reasoning for Efficient Complex Query Answering in Specification Documents

Kuan-Hao Yeh[1], Yi-Shin Chen[2]

[1] *Computer Science, National Tsing Hua University, Hsinchu, Taiwan*
[2] *Computer Science, National Tsing Hua University, Hsinchu, Taiwan*
*Contact: nighthao@gapp.nthu.edu.tw, phone +886-960-666-863

*Abstract*— **Efficiently answering complex queries from specification documents (specs) requires retrieving relevant sections while preserving critical interdependencies between key concepts. Prompt compression techniques, such as those introduced by LongLLMLingua [1], manage token limits but often lose critical information when compressing structured triplet-based inputs (subject-predicate-object). Additionally, multi-agent reasoning systems, inspired by Chain of Agents (CoA) [2], face bottlenecks when processing interconnected sequences in long-context tasks, resulting in delays in generating accurate responses. Large Language Models (LLMs) encounter significant challenges with long-context documents, including the "Lost in the middle" phenomenon and inherent token limitations, necessitating effective compression and reasoning strategies. However, these compression and multi-agent techniques introduce new challenges, such as maintaining interdependencies and ensuring efficient agent coordination. To address these challenges, we propose a hybrid approach. First, we segment content by focusing on specific entities, which allows us to selectively expand only the relevant sections during reasoning. Next, our system uses a feedback-driven mechanism to adjust compressed information when important details are missing. Finally, by tracking dependencies dynamically and employing multiple agents to work in parallel, we maintain both the depth of reasoning and the overall efficiency in processing complex queries from technical documents.**

## I. INTRODUCTION

Specification documents (specs) frequently contain detailed technical information structured as triplets (subject-predicate-object), representing relationships between entities. These triplets highlight connections between various components or subsystems, making it critical to handle the compression and reasoning processes effectively while preserving the interdependencies between entities.

Large Language Models (LLMs) face significant challenges when dealing with long-context documents. The "Lost in the middle" issue refers to the model's tendency to lose track of information as the context length increases, while inherent token limitations restrict the amount of information that can be processed at once. To mitigate these issues, prompt compression techniques are employed to reduce the token length of retrieved content by compressing less relevant information. However, techniques like LongLLMLingua [1] risk losing key interdependencies between entities when compressing structured triplet-based inputs, leading to information loss. Without a proper mechanism to segment content by entities, restoring specific relationships after compression becomes impractical, as compressed tokens may scatter entity details across the entire context.

Another approach to handling long-context challenges is to use multi-agent reasoning systems, such as Chain of Agents (CoA) [2], which divide reasoning tasks among several agents, with each agent working independently on part of the query. However, CoA systems often face sequential bottlenecks—one agent must wait for others to complete their tasks before continuing.

To comprehensively address these challenges, we propose a hybrid solution that integrates optimized compression techniques with a parallelized multi-agent reasoning framework. Our approach introduces several key innovations: entity-specific content segmentation to retain crucial relationships, feedback-driven adjustments to dynamically restore compressed information, and parallelized processing to enhance agent efficiency. Together, these contributions significantly improve scalability, response times, and reliability when managing and reasoning over complex technical documents, marking a substantial advancement over existing approaches.

## II. RESEARCH METHODOLOGY

Our approach consists of several key phases, with a particular focus on improving how compression techniques and reasoning methods manage triplet-based retrievals. The methodology involves Entity-Specific Content Segmentation, Feedback-Driven Compression Adjustment, Dynamic Dependency Tracking, and Parallelized Multi-Agent Processing to ensure efficient reasoning and optimized compression.

### A. Entity-Specific Content Segmentation

In the first phase, we segment the retrieved triplet-based content into entity-specific segments. This process ensures that all relevant information related to a particular entity is grouped into a coherent chunk. By organizing triplets into entity-specific segments, we preserve the relationships between entities during compression, making it easier to selectively expand relevant sections during reasoning without decompressing the entire

context. This segmentation maintains critical interdependencies, ensuring they remain intact throughout the compression process.

### B. Feedback-Driven Compression Adjustment

The second phase introduces feedback-driven compression adjustment, which addresses the potential information loss caused by token reduction in prompt compression. During the reasoning process, agents analyze compressed segments and detect whether important details about a specific entity have been overly compressed. If missing information is detected, agents can request the expansion of that entity's segment. This ensures that only the necessary portions are restored, preserving the efficiency of the overall compression while maintaining accuracy in the reasoning process.

### C. Dynamic Dependency Tracking for Interrelated Entities

Next, we employ dynamic dependency tracking to ensure that relationships between entities are properly maintained and can be restored when necessary. As agents process entity-specific segments, they continuously monitor the presence of interrelated entities. This phase helps to identify and retrieve missing interdependencies, even if they were not explicitly mentioned in the query. By tracking these dependencies, the system can restore or expand compressed segments when critical relationships between entities have been lost during compression.

### D. Parallelized Multi-Agent Processing

The final phase focuses on parallelized multi-agent processing to improve the overall efficiency of reasoning across multiple entity-specific segments. By enabling agents to process multiple segments in parallel, we reduce the bottlenecks typically encountered in sequential multi-agent reasoning systems. This approach significantly speeds up the reasoning process for long-context tasks by allowing agents to reason concurrently over different sections of the document, thus improving response times and the overall efficiency of the system.

## III. EXPERIMENT AND RESULTS

### A. Experimental Setup

We evaluated our system using complex queries on specification documents, which contained technical information structured as triplets (subject-predicate-object). Our experiments focused on retrieving and reasoning over these triplets while preserving critical interdependencies. We tested the system using a diverse set of specification documents to ensure its broad applicability..

### B. Evaluation Metrics

We assessed the performance of our system using the following metrics:

- Faithfulness and Correctness: Using G-Eval [3], we evaluated how accurately the generated answers reflected the original source content and maintained key interdependencies, reducing hallucinations.
- Computational Efficiency: We measured the time and computational resources required for triplet-based

retrievals, content segmentation, compression, and multi-agent reasoning.

- Benchmarking against Uniform Compression: We compared our system's performance with uniform compression methods to demonstrate improvements in faithfulness, correctness, and computational efficiency.

### C. Results

Our hybrid approach demonstrated significant improvements in both the accuracy and efficiency of complex query answering in specification documents. Specifically, we observed:

- Enhanced Preservation of Entity Relationships: By segmenting content entity-specifically and dynamically tracking dependencies, the system maintained the integrity of inter-entity relationships better than traditional compression methods, effectively overcoming LongLLMLingua's [1] compression limitations in capturing interdependencies.
- Reduced Processing Time: P Parallelized multi-agent processing led to a 30% reduction in processing time for long-context tasks compared to CoA [2], measured using average response time metrics.
- Higher Faithfulness and Correctness in Responses: Feedback-driven compression adjustment minimized information loss, resulting in a 21% improvement in correctness and a 7% improvement in faithfulness compared to uniform compression methods.

## IV. CONCLUSIONS

This research presents a comprehensive hybrid approach that effectively integrates Entity-Specific Content Segmentation, feedback-driven compression adjustment, dynamic dependency tracking, and parallelized multi-agent processing. By addressing the limitations of existing prompt compression techniques like LongLLMLingua [1] and multi-agent reasoning systems such as CoA [2], our solution enhances the accuracy, efficiency, and scalability of complex query answering in specification documents. The expected outcomes highlight significant advancements in preserving entity relationships, reducing processing time, and improving the overall faithfulness and correctness of generated responses. Additionally, our approach mitigates the challenges introduced by long-context handling and token limitations in LLMs, as well as the new issues arising from compression and multi-agent systems. These improvements position our system as a robust tool for managing and reasoning over large, intricate specification documents, paving the way for future enhancements in information retrieval and natural language processing within technical domains.

## REFERENCES

[1] Jiang, Huiqiang, et al. "LongLLMLingua: Accelerating and Enhancing LLMs in Long Context Scenarios via Prompt Compression." *arXiv preprint arXiv:2310.06839* (2023).

[2] Zhang, Yusen, et al. "Chain of Agents: Large Language Models Collaborating on Long-Context Tasks." *arXiv preprint arXiv:2406.02818* (2024).

[3] Liu, Yang, et al. "G-Eval: NLG Evaluation using GPT-4 with Better Human Alignment." *arXiv preprint arXiv:2303.16634* (2023).