

# MoC: Mixtures of Text Chunking Learners for Retrieval-Augmented Generation System

Anonymous ACL submission

## Abstract

Retrieval-Augmented Generation (RAG), while serving as a viable complement to large language models (LLMs), often overlooks the crucial aspect of text chunking within its pipeline. This paper initially introduces a dual-metric evaluation method, comprising Boundary Clarity and Chunk Stickiness, to enable the direct quantification of chunking quality. Leveraging this assessment method, we highlight the inherent limitations of traditional and semantic chunking in handling complex contextual nuances, thereby substantiating the necessity of integrating LLMs into chunking process. To address the inherent trade-off between computational efficiency and chunking precision in LLM-based approaches, we devise the granularity-aware Mixture-of-Chunkers (MoC) framework, which consists of a three-stage processing mechanism. Notably, our objective is to guide the chunker towards generating a structured list of chunking regular expressions, which are subsequently employed to extract chunks from the original text. Extensive experiments demonstrate that both our proposed metrics and the MoC framework effectively settle challenges of the chunking task, revealing the chunking kernel while enhancing the performance of the RAG system.

## 1 Introduction

Retrieval-augmented generation (RAG), as a cutting-edge technological paradigm, aims to address challenges faced by large language models (LLMs), such as data freshness (He et al., 2022), hallucinations (Bénédict et al., 2023; Chen et al., 2023; Zuccon et al., 2023; Liang et al., 2024), and the lack of domain-specific knowledge (Li et al., 2023; Shen et al., 2023). This is particularly relevant in knowledge-intensive tasks like open-domain question answering (QA) (Lazaridou et al., 2022). By integrating two key components: the retriever and the generator, this technology enables

more precise responses to input queries (Singh et al., 2021; Lin et al., 2023). While the feasibility of the retrieval-augmentation strategy has been widely demonstrated through practice, its effectiveness heavily relies on the relevance and accuracy of the retrieved documents (Li et al., 2022; Tan et al., 2022). The introduction of excessive redundant or incomplete information through retrieval not only fails to enhance the performance of the generation model but may also lead to a decline in answer quality (Shi et al., 2023; Yan et al., 2024).

In response to the aforementioned challenges, current research efforts mainly focus on two aspects: improving retrieval accuracy (Zhuang et al., 2024; Sidiropoulos and Kanoulas, 2022; Guo et al., 2023) and enhancing the robustness of LLMs against toxic information (Longpre et al.; Kim et al., 2024). However, in RAG systems, a commonly overlooked aspect is the chunked processing of textual content, which directly impacts the quality of dense retrieval for QA (Xu et al., 2023). This is due to the significant “weakest link” effect in the performance of RAG systems, where the quality of text chunking constrains the retrieved content, thereby influencing the accuracy of generated answers (Ru et al., 2024). Despite advancements in other algorithmic components, incremental flaws in the chunking strategy can still detract from the overall system performance to some extent.

Given the critical role of text chunking in RAG systems, optimizing this process has emerged as one of the key strategy to mitigate performance bottlenecks. Traditional text chunking methods, often based on rules or semantic similarity (Zhang et al., 2021; Langchain, 2023; Lyu et al., 2024), provide some structural segmentation but are inadequate in capturing subtle changes in logical relationships between sentences. The LumberChunker (Duarte et al., 2024) offers a novel solution by utilizing LLMs to receive a series of consecutive paragraphs and accurately identify where content

083	begins to diverge. However, it demands a high	subsequently rectifies the generated content.	135
084	level of instruction-following ability from LLMs,	The main contributions of this work are as fol-	136
085	which incurs significant resource and time costs.	lows:	137
086	Additionally, the effectiveness of current chunk-	• Breaking away from indirect evaluation	138
087	ing strategies is often evaluated indirectly through	paradigms, we introduce the dual metrics of	139
088	downstream tasks, such as the QA accuracy in RAG	Boundary Clarity and Chunk Stickiness to	140
089	systems, with a lack of independent metrics for	achieve direct quantification of chunking qual-	141
090	evaluating the inherent rationality of the chunking	ity. By deconstructing the failure mechanisms	142
091	process itself. These challenges give rise to two	of semantic chunking, we provide theoretic-	143
092	practical questions: This raises a practical question:	cal validation for the involvement of LLM in	144
093	How can we fully utilize the powerful reasoning	chunking tasks.	145
094	capabilities of LLMs while accomplishing the text	• We devise the MoC architecture, a hy-	146
095	chunking task at a lower cost? And how to devise	brid framework that dynamically orchestrates	147
096	evaluation metrics that directly quantify the validity	lightweight chunking experts via a multi-	148
097	of text chunking?	granularity-aware router. This architecture	149
098	Inspired by these observations, we innovatively	innovatively integrates: a regex-guided chunk-	150
099	propose two metrics, <b>Boundary Clarity</b> and	ing paradigm, a computation resource con-	151
100	<b>Chunk Stickiness</b> , to independently and effec-	straint mechanism based on sparse activation,	152
101	tively assess chunking quality. Concurrently, we	and a rectification algorithm driven by edit	153
102	leverage these metrics to delve into the reasons	distance.	154
103	behind the suboptimal performance of semantic	• To validate the effectiveness of our pro-	155
104	chunking in certain scenarios, thereby highlighting	posed metrics and chunking method, we con-	156
105	the necessity of LLM-based chunking. To miti-	duct multidimensional experiments using five	157
106	gate the resource overhead of chunking without	different language models across four QA	158
107	compromising the inference performance of LLMs,	datasets, accompanied by in-depth analysis.	159
108	we introduce the <b>Mixture-of-Chunkers (MoC)</b>		
109	framework. This framework primarily comprises	<b>2 Related Works</b>	160
110	a multi-granularity-aware router, specialized meta-	<b>Text Segmentation</b> It is a fundamental task in	161
111	chunkers, and a post-processing algorithm.	NLP, aimed at breaking down text content into its	162
112	This mechanism adopts a divide-and-conquer	constituent parts to lay the foundation for subse-	163
113	strategy, partitioning the continuous granularity	quent advanced tasks such as information retrieval	164
114	space into multiple adjacent subdomains, each cor-	(Li et al., 2020) and text summarization (Lukasik	165
115	responding to a lightweight, specialized chunker.	et al., 2020; Cho et al., 2022). By conducting	166
116	The router dynamically selects the most appropri-	topic modeling on documents, (Kherwa and Bansal,	167
117	ate chunker to perform chunking operation based	2020) and (Barde and Bainwad, 2017) demon-	168
118	on the current input text. This approach not only	strate the identification of primary and sub-topics within	169
119	effectively addresses the “granularity generaliza-	documents as a significant basis for text segmenta-	170
120	tion dilemma” faced by traditional single-model	tion. (Zhang et al., 2021) frames text segmentation	171
121	approaches but also maintains computational re-	as a sentence-level sequence labeling task, utiliz-	172
122	source consumption at the level of a single small	ing BERT to encode multiple sentences simulta-	173
123	language model (SLM) through sparse activation,	neously. It calculates sentence vectors after mod-	174
124	achieving an optimal balance between accuracy	eling longer contextual dependencies and finally	175
125	and efficiency for the chunking system. It is crucial	predicts whether to perform text segmentation after	176
126	to emphasize that our objective is not to require	each sentence. (Langchain, 2023) provides flexible	177
127	the meta-chunker to generate each text chunk in	and powerful support for various text processing	178
128	its entirety. Instead, we guide the model to gener-	scenarios by integrating multiple text segmenta-	179
129	ate a structured list of chunking regular expres-	tion methods, including character segmentation,	180
130	sions used to extract chunks from the original text.	delimiter-based text segmentation, specific docu-	181
131	To address potential hallucination phenomena of	ment segmentation, and recursive chunk segmen-	182
132	meta-chunker, we employ an edit distance recov-	tation. Although these methods better respect the	183
133	ery algorithm, which meticulously compares the		
134	generated chunking rules with the original text and		

structure of the document, they have limitations in deep contextual understanding. To address this issue, semantic-based segmentation (Kamradt, 2024) utilizes embeddings to aggregate semantically similar text chunks and identifies segmentation points by monitoring significant changes in embedding distances.

**Text Chunking in RAG** By expanding the input space of LLMs through introducing retrieved text chunks (Guu et al., 2020; Lewis et al., 2020), RAG significantly improves the performance of knowledge-intensive tasks (Ram et al., 2023). Text chunking allows information to be more concentrated, minimizing the interference of irrelevant information, enabling LLMs to focus more on the specific content of each text chunk and generate more precise responses (Yu et al., 2023; Besta et al., 2024; Su et al., 2024). LumberChunker (Duarte et al., 2024) iteratively harnesses LLMs to identify potential segmentation points within a continuous sequence of textual content, showing some potential for LLMs chunking. However, this method demands a profound capability of LLMs to follow instructions and entails substantial consumption when employing the Gemini model.

### 3 Methodology

#### 3.1 Deep Reflection on Chunking Strategies

As pointed out by Qu et al. (2024), semantic chunking has not shown a significant advantage in many experiments. This paper further explores this phenomenon and proposes two key metrics, "Boundary Clarity" and "Chunk Stickiness", to scientifically explain the limitations of semantic chunking and the effectiveness of LLM chunking. At the same time, it also provides independent evaluation indicators for the rationality of chunking itself.

##### 3.1.1 Boundary Clarity (BC)

Boundary clarity refers to the effectiveness of chunks in separating semantic units. Specifically, it focuses on whether the structure formed by chunking can create clear boundaries between text units at the semantic level. Blurred chunk boundaries may lead to a decrease in the accuracy of subsequent tasks. Specifically, boundary clarity is calculated utilizing the following formula:

$$BC(q, d) = \frac{ppl(q|d)}{ppl(q)} \quad (1)$$

where  $ppl(q)$  represents the perplexity of sentence sequence  $q$ , and  $ppl(q|d)$  denotes the *contrastive*

*perplexity* given the context  $d$ . Perplexity serves as a critical metric for evaluating the predictive accuracy of language models (LMs) on specific textual inputs, where lower perplexity values reveal superior model comprehension of the text, whereas higher values reflect greater uncertainty in semantic interpretation. When the semantic relationship between two text chunks is independent,  $ppl(q|d)$  tends to be closer to  $ppl(q)$ , resulting in the BC metric approaching 1. Conversely, strong semantic interdependence drives the BC metric toward zero. Therefore, higher boundary clarity implies that chunks can be effectively separated, whereas a lower boundary clarity indicates blurred boundaries between chunks, which may potentially lead to information confusion and comprehension difficulties.

##### 3.1.2 Chunk Stickiness (CS)

The objective of text chunking lies in achieving adaptive partitioning of documents to generate logically coherent independent chunks, ensuring that each segmented chunk encapsulates a complete and self-contained expression of ideas while preventing logical discontinuity during the segmentation process. Chunk stickiness specifically focuses on evaluating the tightness and sequential integrity of semantic relationships between text chunks. This is achieved by constructing a semantic association graph among text chunks, where structural entropy is introduced to quantify the network complexity. Within this graph, nodes represent individual text chunks, and edge weights are defined as follows:

$$Edge(q, d) = \frac{ppl(q) - ppl(q|d)}{ppl(q)} \quad (2)$$

where the theoretical range of the Edge value is defined as  $[0, 1]$ . Specifically, we initially compute the Edge value between any two text chunks within a long document. Values approaching 1 indicate that  $ppl(q|d)$  tends towards 0, signifying a high degree of inter-segment correlation. Conversely, an Edge value approaching 0 suggests that  $ppl(q|d)$  converges to  $ppl(q)$ , implying that text chunks are mutually independent. We establish a threshold parameter  $K \in (0, 1)$  to retain edges exceeding this value. Subsequently, the chunk stickiness is specifically calculated as:

$$CS(G) = - \sum_{i=1}^n \frac{d_i}{2m} \cdot \log_2 \frac{d_i}{2m} \quad (3)$$

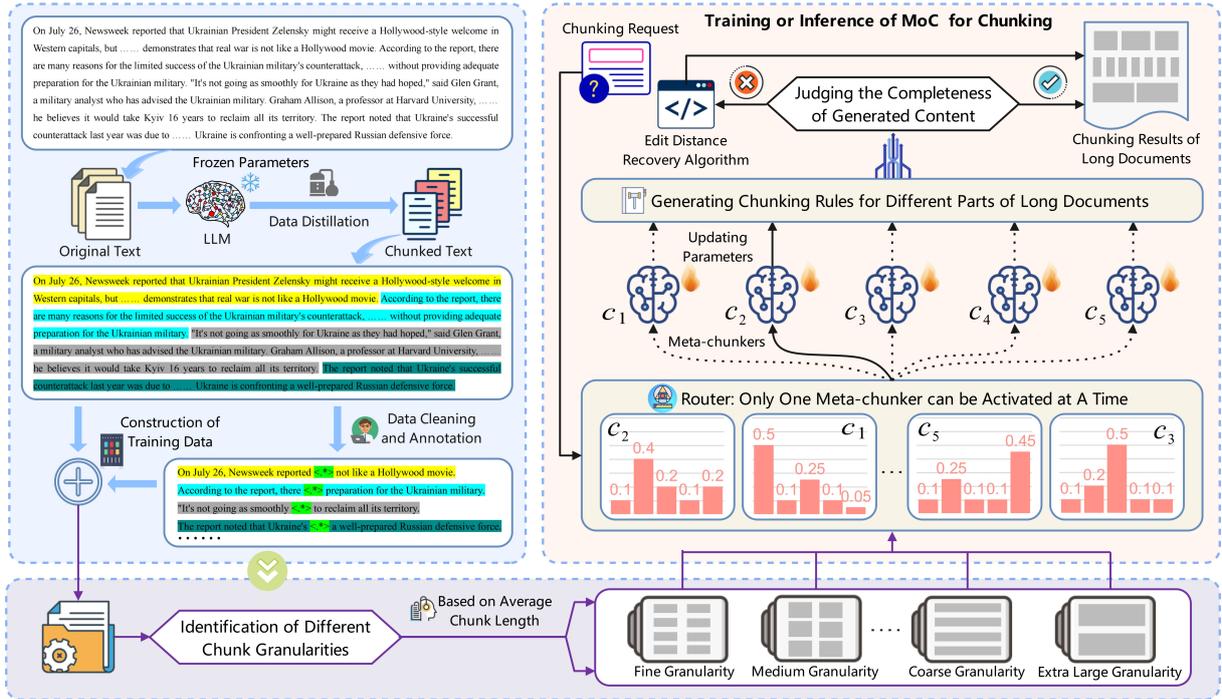


Figure 1: Overview of the entire process of granularity-aware MoC: Dataset construction, training of router and meta-chunkers, as well as chunking inference.

where  $G$  is the constructed semantic graph,  $d_i$  represents the degree of node  $i$ , and  $m$  denotes the total number of edges. This methodology constructs a complete graph, followed by redundancy reduction based on the inter-segment relationships.

On the other hand, to enhance computational efficiency, we construct a sequence-aware incomplete graph that preserves the original ordering of text chunks, which constitutes a graph construction strategy governed by sequential positional constraints. Specifically, given a long text partitioned into an ordered sequence of text chunks  $D = \{d_1, d_2, \dots, d_n\}$ , each node in the graph corresponds to a text chunk, while edge formation is subject to dual criteria: (1) Relevance Criterion: Edge weight  $\text{Edge}(d_i, d_j) > K$ , where  $K$  denotes a predefined threshold; (2) Sequential Constraint: Connections are permitted exclusively when  $j - i > \delta$ , with  $\delta$  representing the sliding window radius fixed at 0. This dual-constraint mechanism strategically incorporates positional relationships, thereby achieving a better equilibrium between semantic relevance and textual coherence.

The detailed design philosophy is elaborated in Appendix A.2. To more intuitively demonstrate the effectiveness of the two metrics, we construct a ‘‘Dissimilarity’’ metric based on the current mainstream semantic similarity, as detailed in Section

4.5. Stemming from the above analysis, we introduce a LM-based training and reasoning framework for text chunking, named granularity-aware MoC.

### 3.2 Granularity-Aware MoC

In response to the complex and variable granularity of large-scale text chunking in real-world scenarios, this paper proposes a multi-granularity chunking framework based on MoC. Our approach, whose overall architecture is illustrated in Figure 1, dynamically routes different granularity experts through a scheduling mechanism and optimizes the integrity of results with a post-processing algorithm.

#### 3.2.1 Dataset Construction

We instruct GPT-4o to generate text chunks from raw long-form texts according to the following criteria: (1) Segmentation: The given text should be segmented according to its logical and semantic structure, such that each resulting chunk maintains a complete and independent logical expression. (2) Fidelity: The segmentation outcome must remain faithful to the original text, preserving its vocabulary and content without introducing any fictitious elements. However, extracting such data from GPT-4o poses significant challenges, as the LLM does not always follow instructions, particularly when

dealing with long texts that contain numerous special characters. In preliminary experiments, we also observed that GPT-4o tends to alter the expressions used in the original text and, at times, generates fabricated content.

To address these challenges, we propose the following dataset distillation procedure. We enhance chunking precision in GPT-4o through structured instructions that enforce adherence to predefined rules. A sliding window algorithm, coupled with a chunk buffering mechanism, mitigates the impact of input text length on performance, ensuring seamless transitions between text subsequences. Furthermore, a rigorous data cleaning process, leveraging edit distance calculations and manual review, addresses potential hallucination, while strategic anchor point extraction and placeholder insertion facilitate efficient processing. Detailed implementation and technical specifics are provided in Appendix A.3.

### 3.2.2 Multi-granularity-aware Router

After the dataset construction is completed, the MoC architecture achieves efficient text processing through the training of the routing decision module and meta-chunkers. The router dynamically evaluates the compatibility of each chunk granularity level based on document features, thereby activating the optimal chunk expert. A major challenge in training the routing module lies in the implicit relationship between text features and chunk granularity, where the goal is to infer the potential granularity of the text without performing explicit chunking operations.

In view of this, we propose a specialized fine-tuning method for SLMs. Firstly, we truncate or concatenate long and short texts respectively, ensuring their lengths hover around 1024 characters. Both operations are performed on text chunks as the operational unit, preserving the semantic integrity of the training texts. By maintaining approximate text lengths, SLMs can better focus on learning features that affect chunk granularity, thus minimizing the impact of text length on route performance. Subsequently, leveraging the segmented data generated by GPT-4o, we assign granularity labels ranging from 0 to 3 to the text, corresponding to average chunk length intervals such as (0, 120], (120, 150], (150, 180], and (180, +∞). The

loss function is formulated as:

$$\mathcal{L}(\theta) = -\frac{1}{N} \sum_{i=1}^N y_i \log(p(y_i|X_i; \theta)) \quad (4)$$

where  $\theta$  represents the set of trainable parameters of the SLM,  $y_i$  denotes the ground-truth granularity label for the  $i$ -th sample,  $N$  signifies the total number of samples, and  $p(y_i|X_i; \theta)$  represents the probability of assigning granularity label  $y_i$ , given input  $X_i$  and current parameters  $\theta$ .

During inference, we implement marginal sampling over the probability distribution of the final token generated by the SLM in its contextual sequence, selecting the granularity category with the highest probability from the four available categories as the granularity for the corresponding text. Afterwards, the text to be chunked is routed to the corresponding chunking expert:

$$R(X_i) = \arg \max_k p(k|X_i; \theta) \quad (5)$$

where  $k$  represents the category of chunking granularity. Through this mechanism, the router enables dynamic expert selection without explicit chunking operations.

### 3.2.3 Meta-chunkers

Our objective is not to require meta-chunkers to generate each text chunk in its entirety, but rather to guide it in producing a structured list of segmented regular expressions. Each element in this list contains only the start  $S$  and end  $E$  of a text chunk  $C$ , with a special character  $r$  replacing the intervening content. The regular expression is represented as:

$$C_{\text{regex}} = S \oplus r \oplus E, \quad r \in \mathcal{R} \quad (6)$$

where  $\oplus$  denotes the string concatenation operation,  $\mathcal{R} = \{“ < omitted > ”, “ < ellipsis > ”, “[MASK]”, “[ELLIPSIS]”, “.*?”, “ < ... > ”, “ < .* > ”, “ < pad > ”\}$  is the set of eight special characters we have defined to represent the omitted parts in a text chunk. During the expert training phase, we employ a full fine-tuning strategy, utilizing datasets categorized by different segmentation granularities to optimize the model parameters. The loss function remains consistent with Equation 4. This design allows Meta-chunkers to comprehensively understand the composition of each chunk while significantly reducing the time cost of generation.

### 3.2.4 Edit Distance Recovery Algorithm

Let string  $A$  denote an element generated by a meta-chunker and string  $B$  represent a segment within the original text. The edit distance refers to the minimum number of operations required to transform  $A$  into  $B$ , where the permissible operations include the insertion, deletion, or substitution of a single character. We then define a two-dimensional array,  $ab[i][j]$ , which represents the minimum number of operations needed to convert the substring  $A[1 \dots i]$  into  $B[1 \dots j]$ . By recursively deriving the state transition formula, we can incrementally construct this array.

Initially, the conditions are as follows: (1) When  $i = 0$ ,  $A$  is an empty string, necessitating the insertion of  $j$  characters to match  $B$ , thus  $ab[0][j] = j$ ; (2) When  $j = 0$ ,  $B$  is an empty string, requiring the deletion of  $i$  characters, hence  $ab[i][0] = i$ ; (3) When  $i = j = 0$ , the edit distance between two empty strings is evidently  $ab[0][0] = 0$ . Subsequently, the entire  $ab$  array is populated using the following state transition formula:

$$ab[i][j] = \begin{cases} ab[i-1][j-1], & \text{if } A[i] = B[j] \\ 1 + \min(ab[i-1][j], \\ ab[i][j-1], \\ ab[i-1][j-1]), & \text{if } A[i] \neq B[j] \end{cases}$$

If the current characters are identical, no additional operation is required, and the problem reduces to a subproblem; if the characters differ, the operation with the minimal cost among insertion, deletion, or substitution is selected. Ultimately, by utilizing the minimum edit distance, we can accurately pinpoint the field in the original text that most closely matches the elements generated by the meta-chunker, thereby ensuring the precision of regular extraction.

## 4 Experiment

### 4.1 Datasets and Metrics

We conduct a comprehensive evaluation on three datasets, and covering multiple metrics. The CRUD benchmark (Lyu et al., 2024) contains single-hop and two-hop questions, evaluated using metrics including BLEU series and ROUGE-L. We utilize the DuReader dataset from LongBench benchmark (Bai et al., 2023), evaluated based on F1 metric. In addition, a dataset called WebCPM (Qin et al., 2023) specifically designed for long-text QA, is utilized to retrieve relevant facts and

generate detailed paragraph-style responses, with ROUGE-L as the metric.

### 4.2 Baselines

We primarily compare meta-chunker and MoC with two types of baselines, namely rule-based chunking and dynamic chunking, noting that the latter incorporates both semantic similarity models and LLMs. The original rule-based method simply divides long texts into fixed-length chunks, disregarding sentence boundaries. However, the Llama\_index method (Langchain, 2023) offers a more nuanced approach, balancing the maintenance of sentence boundaries while ensuring that token counts in each segment are close to a preset threshold. On the other hand, semantic chunking (Xiao et al., 2023) utilizes sentence embedding models to segment text based on semantic similarity. LumberChunker (Duarte et al., 2024) employs LLMs to predict optimal segmentation points within the text.

### 4.3 Experimental Settings

Without additional annotations, all LMs used in this paper adopt chat or instruction versions. When chunking, we primarily employ LMs with the following hyperparameter settings: temperature at 0.1 and top-p at 0.1. For evaluation, Qwen2-7B is applied with the following settings: top\_p = 0.9, top\_k = 5, temperature = 0.1, and max\_new\_tokens = 1280. When conducting QA, the system necessitates dense retrievals from the vector database, with top\_k set to 8 for CRUD, 5 for DuReader and WebCPM. To control variables, we maintain consistent chunk lengths for various chunking methods across each dataset. Detailed experimental setup information can be found in Appendix A.1.

### 4.4 Main Results

To comprehensively validate the effectiveness of the proposed meta-chunker and MoC architectures, we conduct experiments using three widely used QA datasets. During dataset preparation, we curate 20,000 chunked QA pairs through rigorous processing. Initially, we fine-tune the Qwen2.5-1.5B model using this data. As shown in Table 1, compared to traditional rule-based and semantic chunking methods, as well as the state-of-the-art LumberChunker approach based on Qwen2.5-14B, the Meta-chunker-1.5B exhibits both improved and more stable performance. Furthermore, we directly perform chunking employing Qwen2.5-14B and Qwen2.5-72B. The results demonstrate that these

Chunking Methods	CRUD (Single-hop)			CRUD (Two-hop)			DuReader	WebCPM
	BLEU-1	BLEU-Avg	ROUGE-L	BLEU-1	BLEU-Avg	ROUGE-L	F1	ROUGE-L
Original	0.3515	0.2548	0.4213	0.2322	0.1133	0.2613	0.2030	0.2642
Llama_index	0.3620	0.2682	0.4326	0.2315	0.1133	0.2585	0.2220	0.2630
Semantic Chunking	0.3382	0.2462	0.4131	0.2223	0.1075	0.2507	0.2157	0.2691
LumberChunker	0.3456	0.2542	0.4160	0.2204	0.1083	0.2521	0.2178	0.2730
Qwen2.5-14B	0.3650	0.2679	0.4351	0.2304	0.1129	0.2587	0.2271	0.2691
Qwen2.5-72B	0.3721	0.2743	0.4405	<b>0.2382</b>	<b>0.1185</b>	<b>0.2677</b>	0.2284	0.2693
Meta-chunker-1.5B	<b>0.3754</b>	<b>0.2760</b>	<b>0.4445</b>	<u>0.2354</u>	<u>0.1155</u>	<u>0.2641</u>	<b>0.2387</b>	<b>0.2745</b>

Table 1: Main experimental results are presented in four QA datasets. The best result is in bold, and the second best result is underlined.

Methods	BLEU-1	BLEU-2	BLEU-3	BLEU-4	ROUGE-L
<pad>	0.3683	0.2953	0.2490	0.2132	0.4391
<omitted>	0.3725	0.2985	0.2523	0.2165	0.4401
<ellipsis>	0.3761	0.3025	0.2554	0.2193	0.4452
[MASK]	0.3754	0.3012	0.2545	0.2188	0.4445
[ELLIPSIS]	0.3699	0.2966	0.2510	0.2159	0.4380
.*?	0.3745	0.3015	0.2553	0.2195	0.4437
<...>	0.3716	0.2988	0.2526	0.2167	0.4412
<.*>	0.3790	0.3054	0.2583	0.2221	0.4470
MoC	<b>0.3826</b>	<b>0.3077</b>	<b>0.2602</b>	<b>0.2234</b>	<b>0.4510</b>

Table 2: Performance impact of special characters and the effectiveness of granularity-aware MoC framework in text chunking.

LLMs, with their powerful context processing and reasoning abilities, also deliver outstanding performance in chunking tasks. However, Meta-chunker-1.5B slightly underperforms the 72B model only in the two-hop CRUD, while outperforming both LLMs in other scenarios.

Upon validating the effectiveness of our proposed chunking experts, we proceeded to investigate the impact of various special characters on performance, and extended chunking within the MoC framework. As illustrated in Table 2, we design eight distinct special characters, each inducing varying degrees of performance fluctuation in the meta-chunker. Notably, all character configurations demonstrate measurable performance enhancements compared to baseline approaches, with *[Mask]* and *<.\*>* exhibiting particularly remarkable efficacy. In our experiments, both the Meta-chunker-1.5B and the MoC framework employ *[Mask]* as an ellipsis to replace the middle sections of text chunks, while maintaining consis-

tent training data. The experimental results indicate that the chunking method based on the MoC architecture further enhances performance. Specifically, when handling complex long texts, MoC effectively differentiates the chunking granularity of various sections. Moreover, the time complexity of the MoC remains at the level of a single SLM, showcasing a commendable balance between computational efficiency and performance.

#### 4.5 Exploring Chunking Based on Boundary Clarity and Chunk Stickiness

To compare the effectiveness of the two metrics we designed, we introduce the "Dissimilarity" (DS) metric:  $DS = 1 - \text{sim}(q, d)$ , where  $\text{sim}(q, d)$  represents the semantic similarity score between the text chunks  $q$  and  $d$ , and this concept is further concretely illustrated in Appendix A.4. Figure 2(a) reveals the QA performance of RAG using different chunking strategies. It is important to note that, to ensure the validity of the evaluation, we maintained the same average text chunk length across all chunking methods.

**Why Does Semantic Chunking Underperform?** As illustrated in Figure 2(b), while semantic chunking scores are generally high, its performance in QA tasks is suboptimal. Moreover, there is no evident correlation between the scores of semantic dissimilarity and the efficacy of QA. This suggests that in the context of RAG, relying solely on semantic similarity between sentences is insufficient for accurately delineating the optimal boundaries of text chunks.

Furthermore, it can be observed from Table 3 that the clarity of semantic chunking boundaries is only marginally superior to fixed-length chunk-

Chunking Methods	Qwen2.5-1.5B			Qwen2.5-7B			Qwen2.5-14B			Internlm3-8B		
	BC	CS <sub>c</sub>	CS <sub>i</sub>									
Original	0.8210	2.397	1.800	0.8049	2.421	1.898	0.7704	2.297	1.459	0.8054	2.409	1.940
Llama_index	0.8590	2.185	1.379	0.8455	2.250	1.483	0.8117	2.081	1.088	0.8334	2.107	1.303
Semantic Chunking	0.8260	2.280	1.552	0.8140	2.325	1.650	0.7751	2.207	1.314	0.8027	2.255	1.546
Qwen2.5-14B	<b>0.8750</b>	<b>2.069</b>	<b>1.340</b>	<b>0.8641</b>	<b>2.125</b>	<b>1.438</b>	<b>0.8302</b>	<b>1.927</b>	<b>1.068</b>	<b>0.8444</b>	<b>1.889</b>	<b>1.181</b>

Table 3: Performance of different chunking methods under various LMs, directly calculated using two metrics we proposed: BC represents boundary clarity, which is preferable when higher; CS<sub>c</sub> denotes chunk stickiness utilizing a complete graph, and CS<sub>i</sub> indicates chunk stickiness employing an incomplete graph, both of which are favorable when lower.

ing. This implies that although semantic chunking attempts to account for the degree of association between sentences, its limited ability to distinguish logically connected sentences often results in incorrect segmentation of content that should remain coherent. Additionally, Table 3 reveals that semantic chunking also falls short in terms of capturing semantic relationships, leading to higher chunk stickiness and consequently affecting the independence of text chunks.

**Why Does LLM-Based Chunking Work?** As shown in Table 3, the text chunks generated by LLMs exhibit superior boundary clarity, indicating the heightened ability to accurately identify semantic shifts and topic transitions, thereby mitigating the erroneous segmentation of related sentences. Concurrently, the LLM-based chunking produces text chunks with reduced chunk stickiness, signifying that the internal semantics of chunks are more tightly bound, while a greater degree of independence is maintained between chunks. This combination of well-defined boundaries and diminished stickiness contributes to enhanced retrieval efficiency and generation quality within RAG systems, ultimately leading to superior overall performance.

#### 4.6 Hyper-parameter Sensitivity Analysis

In calculating the chunk stickiness, we rely on the  $K$  to filter out edges with weaker associations between text chunks in the knowledge graph. As presented in Table 4, an increase in the value of  $K$  leads to a gradual decrease in the metric. This occurs because a larger  $K$  value limits the number of retained edges, resulting in a sparser connectivity structure within the graph. Notably, regardless of the chosen  $K$  value, the LLM-based chunking method consistently maintains a low level of chunk stickiness. This indicates that it more accurately

identifies semantic transition points between sentences, effectively avoiding excessive cohesion between text chunks caused by interruptions within paragraphs. The sensitivity analysis for generative models in the MoC framework is presented in Appendix A.6.

Methods	Complete Graph			Incomplete Graph		
	0.7	0.8	0.9	0.7	0.8	0.9
Original	2.536	2.397	2.035	2.199	1.800	1.300
Llama_index	2.454	2.185	1.543	1.997	1.379	0.740
Semantic Chunking	2.455	2.280	1.733	2.039	1.552	0.835
Qwen2.5-14B	<b>2.364</b>	<b>2.069</b>	<b>1.381</b>	<b>1.972</b>	<b>1.340</b>	<b>0.623</b>

Table 4: Performance sensitivity of  $K$  in chunk stickiness.

## 5 Conclusion

Addressing the current void in the independent assessment of chunking quality, this paper introduces two novel evaluation metrics: boundary clarity and chunk stickiness. It systematically elucidates the inherent limitations of semantic chunking in long-text processing, which further leads to the necessity of LLM-based chunking. Amidst the drive for performance and efficiency optimization, we propose the MoC framework, which utilizes sparsely activated meta-chunkers through multi-granularity-aware router. It’s worth emphasizing that this study guides meta-chunkers to generate a highly structured list of chunking regular expressions, precisely extracting text chunks from the original text using only a few characters from the beginning and end. Our approach demonstrates superior performance compared to strong baselines.

## 6 Limitations

Despite the superior performance demonstrated by the proposed MoC framework for chunking tasks on various datasets, there are still some limitations that merit further exploration and improvement. Although we have implemented multiple quality control measures to ensure data quality and constructed a training set consisting of nearly 20,000 data entries, the current dataset size remains relatively limited compared to the massive scale and complex diversity of real-world text data. We have mobilized the power of the open-source community to further enrich our chunking dataset utilizing pre-training data from LMs. Additionally, while the dataset construction process is flexible and theoretically expandable to more scenarios, it has not yet undergone adequate multi-language adaptation and validation. We leave this aspect for future research.

## References

Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, et al. 2023. Longbench: A bilingual, multitask benchmark for long context understanding. *arXiv preprint arXiv:2308.14508*.

Bhagyashree Vyankatrao Barde and Anant Madhavrao Bainwad. 2017. An overview of topic modeling methods and tools. In *2017 International Conference on Intelligent Computing and Control Systems (ICICCS)*, pages 745–750. IEEE.

Garbiel Bénédicte, Ruqing Zhang, and Donald Metzler. 2023. Gen-ir@ sigir 2023: The first workshop on generative information retrieval. In *Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 3460–3463.

Maciej Besta, Ales Kubicek, Roman Niggli, Robert Gerstenberger, Lucas Weitzendorf, Mingyuan Chi, Patrick Iff, Joanna Gajda, Piotr Nyczyk, Jürgen Müller, et al. 2024. Multi-head rag: Solving multi-aspect problems with llms. *arXiv preprint arXiv:2406.05085*.

Yuyan Chen, Qiang Fu, Yichen Yuan, Zhihao Wen, Ge Fan, Dayiheng Liu, Dongmei Zhang, Zhixu Li, and Yanghua Xiao. 2023. Hallucination detection: Robustly discerning reliable answers in large language models. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management*, pages 245–255.

Sangwoo Cho, Kaiqiang Song, Xiaoyang Wang, Fei Liu, and Dong Yu. 2022. Toward unifying text segmentation and long document summarization. *arXiv preprint arXiv:2210.16422*.

André V Duarte, João Marques, Miguel Graça, Miguel Freire, Lei Li, and Arlindo L Oliveira. 2024. Lumberchunker: Long-form narrative document segmentation. *arXiv preprint arXiv:2406.17526*.

Zhicheng Guo, Sijie Cheng, Yile Wang, Peng Li, and Yang Liu. 2023. Prompt-guided retrieval augmentation for non-knowledge-intensive tasks. *arXiv preprint arXiv:2305.17653*.

Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938. PMLR.

Hangfeng He, Hongming Zhang, and Dan Roth. 2022. Rethinking with retrieval: Faithful large language model inference. *arXiv preprint arXiv:2301.00303*.

Greg Kamradt. 2024. Semantic chunking. <https://github.com/FullStackRetrieval-com/RetrievalTutorials>.

P Kherwa and P Bansal. 2020. Topic modeling: A comprehensive review. *eai endorsed transactions on scalable information systems*, 7 (24), 1–16.

Youna Kim, Hyuhng Joon Kim, Cheonbok Park, Choonghyun Park, Hyunsoo Cho, Junyeob Kim, Kang Min Yoo, Sang-goo Lee, and Taeuk Kim. 2024. Adaptive contrastive decoding in retrieval-augmented generation for handling noisy contexts. *arXiv preprint arXiv:2408.01084*.

Langchain. 2023. <https://github.com/langchain-ai/langchain>.

Angeliki Lazaridou, Elena Gribovskaya, Wojciech Stokowiec, and Nikolai Grigorev. 2022. Internet-augmented language models through few-shot prompting for open-domain question answering. *arXiv preprint arXiv:2203.05115*.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in Neural Information Processing Systems*, 33:9459–9474.

Huayang Li, Yixuan Su, Deng Cai, Yan Wang, and Lemao Liu. 2022. A survey on retrieval-augmented text generation. *arXiv preprint arXiv:2202.01110*.

Jing Li, Billy Chiu, Shuo Shang, and Ling Shao. 2020. Neural text segmentation and its application to sentiment analysis. *IEEE Transactions on Knowledge and Data Engineering*, 34(2):828–842.

Xianzhi Li, Samuel Chan, Xiaodan Zhu, Yulong Pei, Zhiqiang Ma, Xiaomo Liu, and Sameena Shah. 2023. Are chatgpt and gpt-4 general-purpose solvers for financial text analytics? a study on several typical tasks. *arXiv preprint arXiv:2305.05862*.

739	Xun Liang, Shichao Song, Zifan Zheng, Hanyu Wang, Qingchen Yu, Xunkai Li, Rong-Hua Li, Feiyu Xiong, and Zhiyu Li. 2024. Internal consistency and self-feedback in large language models: A survey. <i>arXiv preprint arXiv:2407.14507</i> .	794	Georgios Sidiropoulos and Evangelos Kanoulas. 2022. Analysing the robustness of dual encoders for dense retrieval against misspellings. In <i>Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 2132–2136.	795
740		796		797
741		798		799
742				
743				
744	Weizhe Lin, Rexhina Blloshmi, Bill Byrne, Adrià de Gispert, and Gonzalo Iglesias. 2023. Li-rage: Late interaction retrieval augmented generation with explicit signals for open-domain table question answering. In <i>Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 1557–1566.	800	Devendra Singh, Siva Reddy, Will Hamilton, Chris Dyer, and Dani Yogatama. 2021. End-to-end training of multi-document reader and retriever for open-domain question answering. <i>Advances in Neural Information Processing Systems</i> , 34:25968–25981.	801
745		802		803
746		804		805
747				
748				
749				
750				
751	S Longpre, G Yauney, E Reif, K Lee, A Roberts, B Zoph, D Zhou, J Wei, K Robinson, D Mimno, et al. A pre-trainer’s guide to training data: Measuring the effects of data age, domain coverage, quality, & toxicity, may 2023. URL <a href="http://arxiv.org/abs/2305.13169">http://arxiv.org/abs/2305.13169</a> .	806	Weihang Su, Yichen Tang, Qingyao Ai, Zhijing Wu, and Yiqun Liu. 2024. Dragin: Dynamic retrieval augmented generation based on the real-time information needs of large language models. <i>arXiv preprint arXiv:2403.10081</i> .	807
752		808		809
753				
754				
755				
756	Michal Lukasik, Boris Dadachev, Gonçalo Simoes, and Kishore Papineni. 2020. Text segmentation by cross segment attention. <i>arXiv preprint arXiv:2004.14535</i> .	810	Chao-Hong Tan, Jia-Chen Gu, Chongyang Tao, Zhen-Hua Ling, Can Xu, Huang Hu, Xiubo Geng, and Daxin Jiang. 2022. Tegtok: Augmenting text generation via task-specific and open-world knowledge. <i>arXiv preprint arXiv:2203.08517</i> .	811
757		812		813
758		814		815
759	Yuanjie Lyu, Zhiyu Li, Simin Niu, Feiyu Xiong, Bo Tang, Wenjin Wang, Hao Wu, Huanyong Liu, Tong Xu, and Enhong Chen. 2024. Crud-rag: A comprehensive chinese benchmark for retrieval-augmented generation of large language models. <i>arXiv preprint arXiv:2401.17043</i> .	816	Shitao Xiao, Zheng Liu, Peitian Zhang, and N Muenighof. 2023. C-pack: packaged resources to advance general chinese embedding. 2023. <i>arXiv preprint arXiv:2309.07597</i> .	817
760		818		819
761				
762				
763				
764				
765	Yujia Qin, Zihan Cai, Dian Jin, Lan Yan, Shihao Liang, Kunlun Zhu, Yankai Lin, Xu Han, Ning Ding, Huadong Wang, et al. 2023. Webcpm: Interactive web search for chinese long-form question answering. <i>arXiv preprint arXiv:2305.06849</i> .	820	Shicheng Xu, Liang Pang, Huawei Shen, and Xueqi Cheng. 2023. Berm: Training the balanced and extractable representation for matching to improve generalization ability of dense retrieval. <i>arXiv preprint arXiv:2305.11052</i> .	821
766		822		823
767				
768				
769				
770	Renyi Qu, Ruixuan Tu, and Forrest Bao. 2024. Is semantic chunking worth the computational cost? <i>arXiv preprint arXiv:2410.13070</i> .	824	Shi-Qi Yan, Jia-Chen Gu, Yun Zhu, and Zhen-Hua Ling. 2024. Corrective retrieval augmented generation. <i>arXiv preprint arXiv:2401.15884</i> .	825
771		826		
772				
773	Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay, Amnon Shashua, Kevin Leyton-Brown, and Yoav Shoham. 2023. In-context retrieval-augmented language models. <i>Transactions of the Association for Computational Linguistics</i> , 11:1316–1331.	827	Wenhao Yu, Hongming Zhang, Xiaoman Pan, Kaixin Ma, Hongwei Wang, and Dong Yu. 2023. Chain-of-note: Enhancing robustness in retrieval-augmented language models. <i>arXiv preprint arXiv:2311.09210</i> .	828
774		829		830
775				
776				
777				
778	Dongyu Ru, Lin Qiu, Xiangkun Hu, Tianhang Zhang, Peng Shi, Shuaichen Chang, Cheng Jiayang, Cunxiang Wang, Shichao Sun, Huanyu Li, et al. 2024. Ragchecker: A fine-grained framework for diagnosing retrieval-augmented generation. <i>arXiv preprint arXiv:2408.08067</i> .	831	Qinglin Zhang, Qian Chen, Yali Li, Jiaqing Liu, and Wen Wang. 2021. Sequence model with self-adaptive sliding window for efficient spoken document segmentation. In <i>2021 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)</i> , pages 411–418. IEEE.	832
779		833		834
780		835		836
781				
782				
783				
784	Xinyue Shen, Zeyuan Chen, Michael Backes, and Yang Zhang. 2023. In chatgpt we trust? measuring and characterizing the reliability of chatgpt. <i>arXiv preprint arXiv:2304.08979</i> .	837	Ziyuan Zhuang, Zhiyang Zhang, Sitao Cheng, Fangkai Yang, Jia Liu, Shujian Huang, Qingwei Lin, Saravan Rajmohan, Dongmei Zhang, and Qi Zhang. 2024. Efficientrag: Efficient retriever for multi-hop question answering. <i>arXiv preprint arXiv:2408.04259</i> .	838
785		839		840
786		841		
787				
788	Freda Shi, Xinyun Chen, Kanishka Misra, Nathan Scales, David Dohan, Ed H Chi, Nathanael Schärli, and Denny Zhou. 2023. Large language models can be easily distracted by irrelevant context. In <i>International Conference on Machine Learning</i> , pages 31210–31227. PMLR.	842	Guido Zuccon, Bevan Koopman, and Razia Shaik. 2023. Chatgpt hallucinates when attributing answers. In <i>Proceedings of the Annual International ACM SIGIR Conference on Research and Development in Information Retrieval in the Asia Pacific Region</i> , pages 46–51.	843
789		844		845
790		846		847
791				
792				
793				

## A Appendix

### A.1 Main Experimental Details

All language models utilized in this paper employ the chat or instruct versions where multiple versions exist, and are loaded in full precision (Float32). The vector database is constructed using Milvus, where the embedding model is bge-large-zh-v1.5. In experiments, we utilized a total of four benchmarks, and their specific configurations are detailed as follows:

#### (a) Rule-based Chunking Methods

- **Original:** This method divides long texts into segments of a fixed length, such as two hundred Chinese characters or words, without considering sentence boundaries.
- **Llama\_index (Langchain, 2023):** This method considers both sentence completeness and token counts during segmentation. It prioritizes maintaining sentence boundaries while ensuring that the number of tokens in each chunk are close to a preset threshold. We use the SimpleNodeParser function from Llama\_index, adjusting the chunk\_size parameter to control segment length. Overlaps are handled by dynamically overlapping segments using the chunk\_overlap parameter, ensuring sentence completeness during segmentation and overlapping.

#### (b) Dynamic Chunking Methods

- **Semantic Chunking (Xiao et al., 2023):** Utilizes pre-trained sentence embedding models to calculate the cosine similarity between sentences. By setting a similarity threshold, sentences with lower similarity are selected as segmentation points, ensuring that sentences within each chunk are highly semantically related. This method employs the SemanticSplitterNodeParser from Llama\_index, exploiting the bge-base-zh-v1.5 model. The size of the text chunks is controlled by adjusting the similarity threshold.
- **LumberChunker (Duarte et al., 2024):** Leverages the reasoning capabilities of LLMs to predict suitable segmentation

points within the text. We utilize Qwen2.5 models with 14B parameters, set to full precision.

### A.2 Design Philosophy of Chunk Stickiness

In the context of network architecture, high structural entropy tends to exhibit greater challenges in predictability and controllability due to its inherent randomness and complexity. Our chunking strategy aims to maximize semantic independence between text chunks while maintaining a coherent semantic expression. Consequently, a higher chunk stickiness implies greater interconnectedness among these chunks, resulting in a more intricate and less ordered semantic network. Furthermore, to ensure a robust comparison between different chunking methods, we enforce a uniform average chunking length. This standardization provides a fair basis for evaluation, mitigating potential biases arising from discrepancies in chunking size. Ultimately, a lower CS score signifies that the chunking method is more accurate in identifying semantic transition points between sentences, thereby avoiding the fragmentation of coherent passages and the consequent excessive stickiness between resulting chunks.

To more intuitively demonstrate the effectiveness of the two metrics we designed, we construct a “Dissimilarity” metric based on the current mainstream semantic similarity, as detailed in Section 4.5. Furthermore, employing several chunking techniques and LLMs, we conduct an in-depth investigation of boundary clarity and chunk stickiness, conducting comparative experiments with the dissimilarity metric. The experimental results clearly show that the two proposed metrics exhibit a consistent trend with RAG performance when evaluating the quality of text chunking. In contrast, the dissimilarity metric fail to display a similar variation. This suggests that, even without relying on QA accuracy, the two proposed metrics can independently and effectively assess chunking quality.

### A.3 Dataset Construction Process

**Structured Instruction Design** By explicitly enumerating rules, GPT-4o is compelled to adhere to predefined chunking regulations, such as ensuring semantic unit integrity, enforcing punctuation boundaries, and prohibiting content rewriting.

**Sliding Window and Chunk Buffering Mechanism** Drawing from the research conducted by Duarte et al. (2024) and practical experience, we

observe that the length of the original text significantly influences the chunking performance of LLMs. To address this problem, we initially apply a sliding window algorithm to segment the input text into subsequences, each below a threshold of 1024 tokens. Segmentation points are prioritized at paragraph boundaries or sentence-ending positions. These subsequences are then processed sequentially by GPT-4o. To maintain continuity between two consecutive subsequences, we implement a chunk buffer mechanism by removing the last generated text chunk of the preceding sequence and using it as the prefix for the subsequent sequence, thereby ensuring smooth information flow.

**Data Cleaning and Annotation** To identify and eliminate hallucinated content during the generation process, we calculate the difference between each chunk and the paragraphs in the original text through the edit distance, as outlined in Section 3.2.4. If the minimum edit distance exceeds 10% of the chunk length, we manually review the location of the chunk error and make corrections accordingly. Additionally, for a long text, we extract several characters at the beginning and end of each text chunk as anchor points, while replacing the intermediate content with eight preset special placeholders, as demonstrated in Sections 3.2.2 and 3.2.3.

#### A.4 Design Philosophy of Dissimilarity

To compare the effectiveness of the two metrics we designed, we introduce the "Dissimilarity" (DS) metric:

$$DS = 1 - \text{sim}(q, d)$$

where  $\text{sim}(q, d)$  represents the semantic similarity score between the text chunks  $q$  and  $d$ . With this definition, the DS metric ranges from  $[0, 1]$ , where 0 indicates perfect similarity and 1 indicates complete dissimilarity. The design of the DS metric is based on the following considerations: first, semantic similarity measures are typically employed to assess the degree of semantic proximity between two text segments. By converting this to the dissimilarity measure, we can more directly observe the semantic differences between chunks. Second, the linear transformation of DS preserves the monotonicity of the original similarity measure without losing any information.

Methods	Qwen-1.5B	Qwen-7B	Qwen-14B	Internlm-8B
Original	2.206	2.650	2.560	1.636
Llama_index	1.964	2.412	2.353	1.486
Semantic Chunking	1.865	2.331	2.238	1.411
LumberChunker	2.184	2.593	2.589	1.652
Qwen2.5-14B	1.841	2.313	2.209	1.373
Meta-chunker-1.5B	<b>1.835</b>	<b>2.267</b>	<b>2.199</b>	<b>1.367</b>

Table 5: Information-based performance evaluation for the RAG system.

#### A.5 Another Perspective on Chunking Performance Comparison

The performance evaluation of RAG systems primarily focuses on the similarity between generated answers and reference answers. However, this evaluation method introduces additional noise during the decoding strategy in the generation stage, making it difficult to distinguish whether the performance defects originate from the retrieved chunk or the generation module. To address this constraint, we propose an evaluation approach based on information support, which centers on quantifying the supporting capability of retrieved text chunks for the target answer through conditional probability modeling.

Given a set of retrieved chunks  $C = \{c_1, c_2, \dots, c_n\}$  and the reference answer  $A = \{a_1, a_2, \dots, a_m\}$ , we employ a LLM to compute the average conditional probability (CP) of the target answer:

$$CP = -\frac{1}{M} \sum_{i=1}^M \log P(a_i | c_1, c_2, \dots, c_n) \quad (7)$$

A smaller CP value indicates a higher likelihood of the correct answer being inferred from the retrieved text chunks, signifying stronger support. The results presented in Table 5 show that, even when evaluated with different LMs, our chunking method consistently exhibits high support. This suggests that our chunking strategy, by optimizing the semantic integrity and independence of text chunks, enhances the relevance of the retrieved text to the question, thereby reducing the difficulty of generating the correct answer.

#### A.6 Hyper-parameter Sensitivity Analysis of Meta-chunker

We conducted experiments on the decoding sampling hyperparameters of the meta-chunker within the MoC framework, with specific results presented

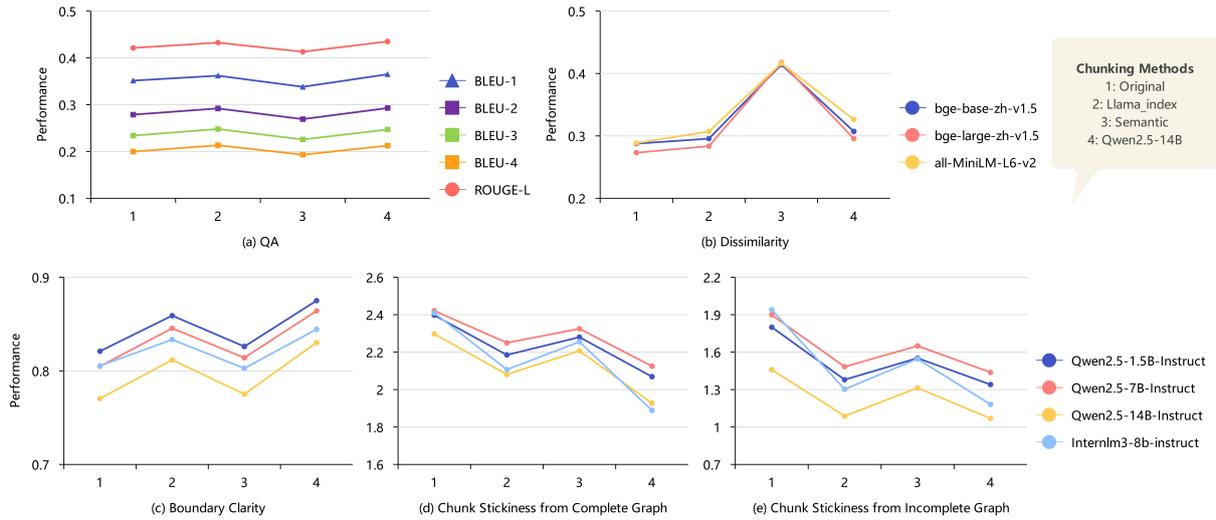


Figure 2: Trends in evaluating chunking performance using different metrics.

1030 in Table 3. Experimental data demonstrates that  
 1031 higher values of temperature and top-k sampling  
 1032 strategies introduce increased randomness, thereby  
 1033 exerting a certain impact on the chunking effect.  
 1034 Conversely, when these two hyperparameters are  
 1035 set to lower values, the model typically provides  
 1036 more stable and precise chunking, leading to a  
 1037 more significant performance improvement.

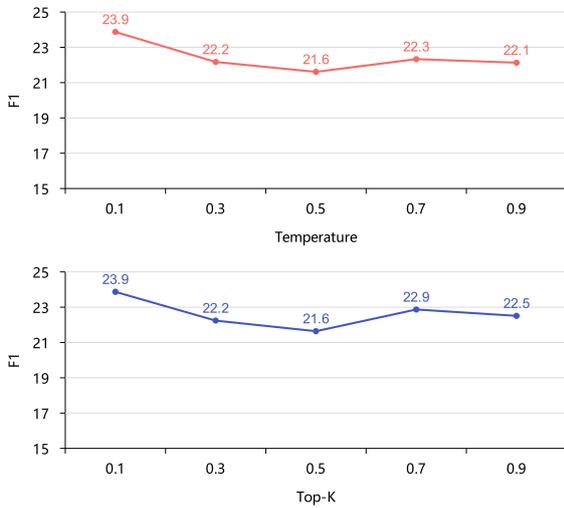


Figure 3: Performance sensitivity to temperature and top-k.

### 1038 A.7 Prompt utilized in Chunking

1039 When preparing datasets using GPT-4o and generat-  
 1040 ing chunking rules with MoC, prompts are neces-  
 1041 sary, as illustrated in Tables 6 and 7. The design  
 1042 and implementation of these prompts are crucial,  
 1043 as they directly influence the quality and charac-  
 1044 teristics of the resulting datasets and chunking rules.

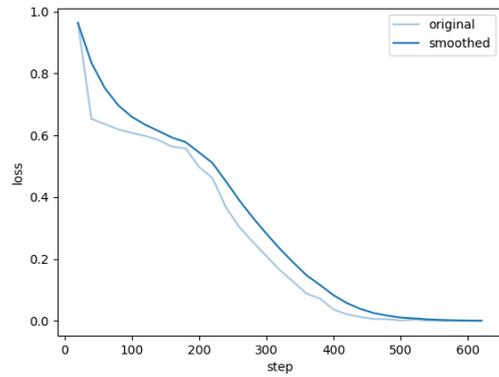


Figure 4: Trend of loss change during router training.

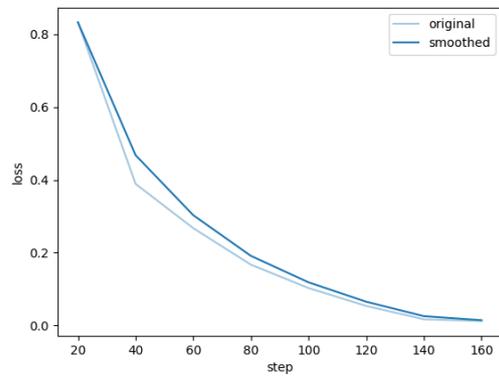


Figure 5: Trend of loss change during meta-chunker training with granularity range [0,120].

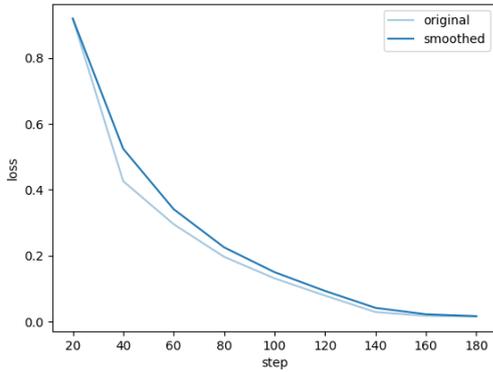


Figure 6: Trend of loss change during meta-chunker training with granularity range (120,150].

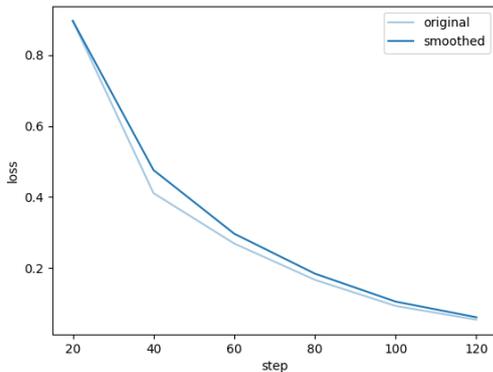


Figure 7: Trend of loss change during meta-chunker training with granularity range (150,180].

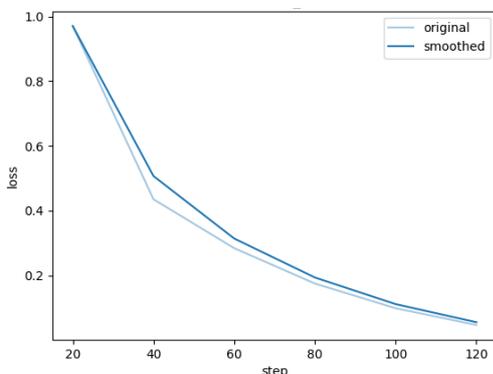


Figure 8: Trend of loss change during meta-chunker training with granularity range (180, +infinity).

---

### Chunking Prompt

---

This is a text chunking task, and you are an expert in text segmentation, responsible for dividing the given text into text chunks. You must adhere to the following four conditions:

1. Segment the text based on its logical and semantic structure, ensuring each text chunk expresses a complete logical thought.
2. Avoid making the text chunks too short, balancing the recognition of content transitions with appropriate chunk length.
3. Do not alter the original vocabulary or content of the text.
4. Do not add any new words or symbols.

If you understand, please segment the following text into text chunks, with each chunk separated by "\n—\n". Output the complete set of segmented chunks without omissions.

Document content: [Text to be segmented]

The segmented text chunks are:

---

Table 6: Prompt for direct chunking of GPT-4o.

---

### **Chunking Prompt**

---

This is a text chunking task. As an expert in text segmentation, you are responsible for segmenting the given text into text chunks. You must adhere to the following four conditions:

1. Combine several consecutive sentences with related content into text chunks, ensuring that each text chunk has a complete logical expression.
2. Avoid making the text chunks too short, and strike a good balance between recognizing content transitions and chunk length.
3. The output of the chunking result should be in a list format, where each element represents a text chunk in the document.
4. Each text chunk in the output should consist of the first few characters of the text chunk, followed by "[MASK]" to replace the intermediate content, and end with the last few characters of the text chunk. The output format is as follows:

```
[  
  "First few characters of text chunk [MASK] Last few characters of text chunk",  
  ...  
]
```

If you understand, please segment the following text into text chunks and output them in the required list format.

Document content: [Text to be segmented]

---

Table 7: Prompt for chunking of MoC.