

Beyond Chunking: Discourse-Aware Hierarchical Retrieval for Long Document Question Answering

Huiyao Chen^{1,2}, Yi Yang¹, Yinghui Li³, Meishan Zhang^{1,*}, Baotian Hu^{1,2}, Min Zhang^{1,2}

¹Institute of Computing and Intelligence, Harbin Institute of Technology (Shenzhen), China

² Shenzhen Loop Area Institute (SLAI) ³ Tsinghua University

{chenhy1018, yangi10010730, mason.zms}@gmail.com

Abstract

Existing long-document question answering systems typically process texts as flat sequences or use heuristic chunking, which overlook the discourse structures that naturally guide human comprehension. We present a discourse-aware hierarchical framework that leverages rhetorical structure theory (RST) for long document question answering. Our approach converts discourse trees into sentence-level representations and employs LLM-enhanced node representations to bridge structural and semantic information. The framework involves three key innovations: language-universal discourse parsing for lengthy documents, LLM-based enhancement of discourse relation nodes, and structure-guided hierarchical retrieval. Extensive experiments on four datasets demonstrate consistent improvements over existing approaches through the incorporation of discourse structure, across multiple genres and languages. Moreover, the proposed framework exhibits strong robustness across diverse document types and linguistic settings.

1 Introduction

Document question answering represents a fundamental challenge in natural language processing, with applications spanning from academic research to enterprise knowledge management (Chen et al., 2017; Karpukhin et al., 2020). To date, large language models (LLMs) have achieved remarkable success on the task, particularly for short documents such as SQuAD (Rajpurkar et al., 2016), which has an average context length of 117 words (Reddy et al., 2019), reaching human-level performance with F1 scores exceeding 85% (Chowdhery et al., 2023; Malladi et al., 2023). However, as document length increases, their performance degrades significantly. For example, on challenging long document datasets such as QASPER, state-of-the-art LLM models only achieve performance

* Corresponding author: Meishan Zhang.

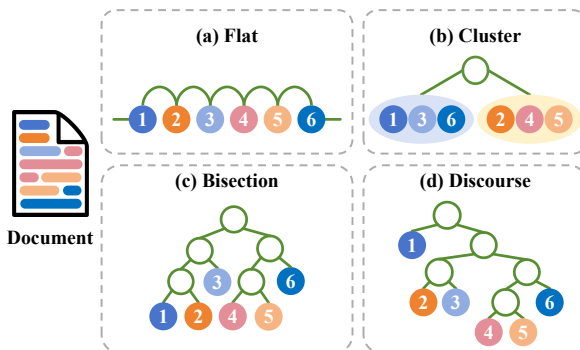


Figure 1: Comparison of document modeling approaches for long document QA. Numbers (1-6) show sentence order in original document, with similar colors indicating semantic relationships. Four approaches are compared: (a) Flat sequential modeling, (b) Bottom-up semantic clustering, (c) Bisection-based adjacent grouping, and (d) Our discourse-aware approach that preserves both semantic and discourse structures.

with less than 50% F1 scores (Shaham et al., 2022; Bai et al., 2024). This performance gap highlights long document question answering as a critical and underexplored research frontier.

The most straightforward approach to handling long documents involves flat sequential modeling, where texts are processed as linear sequences or divided into fixed-size chunks (Lewis et al., 2020; Guu et al., 2020a). This chunking-based paradigm has gained widespread adoption due to its simplicity and computational efficiency, enabling scalable processing of large document collections. Building upon this foundation, recent advances have introduced tree-based structures to better capture document organization and move beyond simple chunking, as illustrated in Figure 1. These innovations include semantic clustering methods like RAPTOR (Sarthi et al., 2024) that recursively group similar content, and bisection-based approaches that maintain local coherence through adjacent segment grouping (Liu and Lapata, 2019; Karpukhin et al., 2020). These hierarchical methods represent sig-

nificant improvements by capturing semantic similarity and textual coherence, offering more sophisticated document modeling capabilities.

Discourse structure provides a principled linguistic foundation for document organization that goes beyond surface-level similarities. Consider a document where the first sentence serves as a topic statement, followed by contrasting ideas (sentences 2-3), parallel supporting evidence (sentences 4-5), and a concluding statement (sentence 6). These discourse relationships naturally capture how humans organize and comprehend information in documents, providing guidance for more semantically appropriate chunking and hierarchical organization in retrieval. As a mature theoretical framework, discourse analysis has demonstrated effectiveness across various NLP tasks (Feng and Hirst, 2014; Cohan et al., 2018; Søgaard et al., 2021), presenting an opportunity to develop more principled retrieval methods that align with human document comprehension patterns.

Therefore, in this work, we present **DISRetrieval** (**DIS**course-aware hierarchical **Retrieval**), the first systematic approach that leverages rhetorical structure theory (RST) to enhance long document question answering through discourse-aware modeling with cross-lingual applicability. Our approach tackles three key technical challenges: First, we introduce RST adaptations along two dimensions: granularity adaptation that shifts processing to the sentence level for efficiency, and language adaptation that enables cross-lingual applicability through LLM-based data augmentation. Second, we introduce LLM-based node enhancement that enriches intermediate nodes with both discourse structure and semantic content. Third, we design structure-guided evidence retrieval mechanisms that leverage discourse organization for effective information extraction.

Comprehensive experiments across four challenging benchmarks: QASPER (Dasigi et al., 2021) for research paper understanding, QuALITY (Pang et al., 2022) for reading comprehension, NarrativeQA (Kočišký et al., 2018) for book-length narrative analysis and MultiFieldQA-zh (Bai et al., 2024) for Chinese documents question answering demonstrate substantial improvements over existing methods. Detailed ablation studies validate the effectiveness of our discourse-aware modeling and structure-guided processing, showing consistent gains across diverse scenarios. Our framework successfully captures both fine-grained semantic

details and document-level organizational patterns, providing a robust solution for discourse-informed question answering. Our main contributions can be summarized as follows:

- A novel discourse-aware hierarchical framework with granularity and language adaptations enabling efficient cross-lingual long document QA.
- An innovative LLM-enhanced hierarchical retrieval mechanism enabling multi-granularity evidence selection.
- Comprehensive empirical validation across diverse datasets, architectures, and languages.

Our code and datasets are publicly available at [github/DreamH1gh/DISRetrieval](https://github.com/DreamH1gh/DISRetrieval) to facilitate future research.

2 Related Work

Recent advances in LLMs have demonstrated impressive capabilities across diverse text comprehension tasks (Brown et al., 2020; Chowdhery et al., 2023; Zhang et al., 2025; Wu et al., 2025; Zhao et al., 2026; Wu et al., 2026). However, extended documents present challenges due to context length constraints and computational complexity (Tay et al., 2021; Zaheer et al., 2020; Beltagy et al., 2020; Ding et al., 2023; Ainslie et al., 2023). Traditional approaches operate on short segments without considering broader context (Liu and Lapata, 2019; Guo et al., 2022; Rae et al., 2020), while recent innovations explore chunking-free extraction (Zhao et al., 2024) and in-context retrieval (Qian et al., 2024; Wang et al., 2023; Izacard et al., 2023; Yu et al., 2023).

Retrieval-based methods segment documents before retrieving relevant evidence (Lewis et al., 2020; Karpukhin et al., 2020; Izacard and Grave, 2020). These have evolved from flat segmentation (Guo et al., 2020b) and early hierarchical approaches (Tang et al., 2017) to sophisticated methods including semantic clustering (Sarathi et al., 2024), bisection-based techniques (Zhang et al., 2020; Ivgi et al., 2023), hybrid systems (Liu et al., 2021; Arivazhagan et al., 2023), and structure conversion approaches (Jin et al., 2025). Discourse structure has been explored for QA through discourse-based systems (Santhosh and Ali, 2012), long-form answer analysis (Xu et al., 2022), and structure-discourse graphs (Nair et al., 2023; Du et al., 2023). However, existing approaches either rely on surface-level semantic similarity for re-

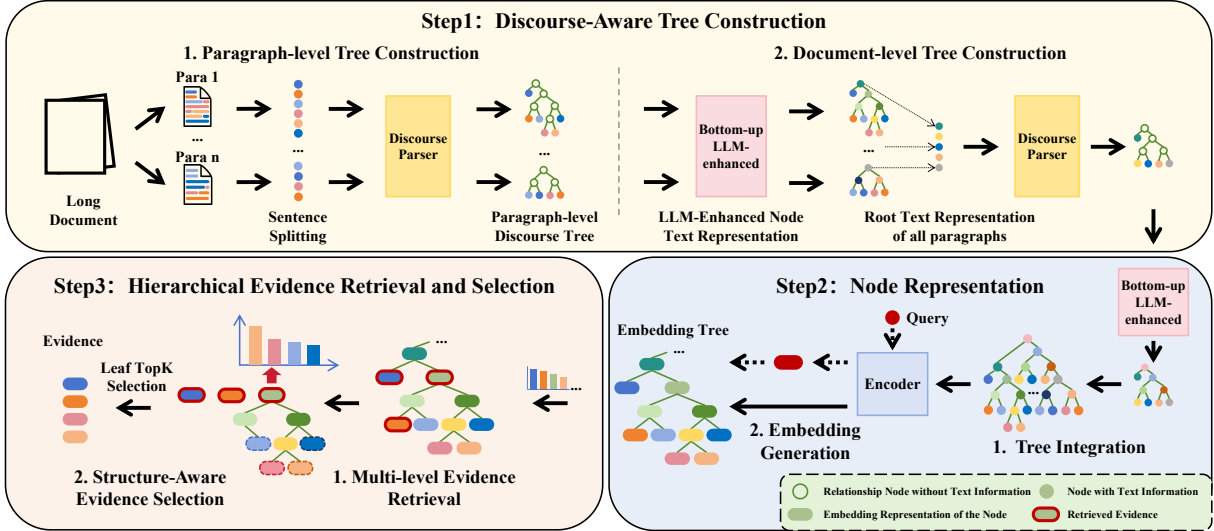


Figure 2: Overview of the DISRetrieval framework consisting of three main steps: (1) Discourse-Aware Tree Construction through paragraph-level parsing (left) and document-level integration (center), (2) Node Representation via tree integration and dense encoding, and (3) Hierarchical Evidence Retrieval with multi-level selection.

trieval or apply discourse primarily to answer generation tasks, overlooking linguistically-informed discourse structures for evidence retrieval (Liu et al., 2019; Wolf et al., 2020).

3 Method

Discourse analysis provides systematic frameworks for understanding how textual segments relate beyond surface-level semantics (Mann and Thompson, 1988; Carlson et al., 2001). RST formalizes this by representing documents as hierarchical trees where leaf nodes correspond to elementary discourse units (EDUs) and internal nodes encode rhetorical relations such as contrast, elaboration, and summary (Marcu, 2000). Recent advances in neural discourse parsing (Yu et al., 2018; Kobayashi et al., 2020; Yu et al., 2022a,b) have enabled automatic RST construction at scale (Maekawa et al., 2024; Yuan et al., 2025).

To address the challenge of bridging between discourse theory and document retrieval, we propose **DISRetrieval**, a discourse-aware framework that systematically incorporates RST structure into document retrieval as shown in Figure 2. Our framework operates in three stages: (1) constructing hierarchical discourse trees through sentence-level RST parsing, (2) enriching tree nodes with semantic representations via LLM enhancement and dense encoding, and (3) performing structure-guided evidence retrieval. Figure 5 provides a concrete example showing how a research paper paragraph is transformed into a hierarchical discourse structure.

3.1 Discourse-Aware Tree Construction

We construct a hierarchical tree structure that captures semantic content and organizational structure of documents through discourse analysis.

3.1.1 RST Adaptation

To effectively apply RST to long document question answering, we introduce two critical adaptations that address computational efficiency and cross-lingual applicability.

Granularity Adaptation. Traditional RST operates on fine-grained EDUs, creating computational overhead and semantic fragmentation for long documents. We address this challenge by shifting RST processing to the sentence level, achieving better efficiency while maintaining semantic coherence. Specifically, we train a sentence-level parser by converting existing EDU-based datasets through two operations: (1) merging intra-sentence EDUs into unified sentence units, and (2) determining inter-sentence relationships via lowest common ancestor analysis in the original EDU-level discourse trees. This enables our parser to capture meaningful discourse relations efficiently.

Language Adaptation. To enable cross-lingual applicability beyond English, we develop a language-universal discourse parser through LLM-based multilingual data augmentation. We employ GPT-4o to translate the RST-DT training corpus into target languages while preserving discourse structures at the sentence level. The translated

data is combined with the original corpus to train a unified parser $f_{\text{discourse}}$ that generalizes across languages without requiring language-specific annotations. Complete implementation details are provided in Appendix A.1 and B.2.

3.1.2 Complete Tree Construction Process

Building upon our adapted sentence-level discourse parser, we develop a hierarchical discourse modeling framework that constructs document-level discourse trees through a two-phase process.

Phase 1. Paragraph-Level Tree Construction.

We transform each paragraph’s sentence sequence $\mathbf{S}_i = \{s_{i,1}, s_{i,2}, \dots, s_{i,m}\}$ into a local discourse tree T_i using $f_{\text{discourse}}$. Each tree captures sentence connections within the paragraph, preserving fine-grained intra-paragraph discourse structure.

Phase 2. Document-Level Tree Construction.

We construct global document structure through bottom-up LLM enhancement. For each internal node v in paragraph-level trees, we apply adaptive processing:

$$v^* = \begin{cases} f_{\text{LLM}}(v_l, v_r), & \text{if } |v_l| + |v_r| \geq \tau \\ f_{\text{merge}}(v). & \text{otherwise} \end{cases} \quad (1)$$

When combined child length exceeds threshold τ , f_{LLM} generates concise summaries. Otherwise, f_{merge} directly concatenates content. Figure 5 illustrates how rhetorical relations transform into semantic representations.

The root representations from paragraph trees serve as input for document-level tree construction. We apply the same RST parser to these summaries, creating a hierarchical structure capturing document-wide discourse relationships.

3.2 Discourse-Aware Node Representation

A critical challenge in discourse-based retrieval lies in bridging the semantic gap between abstract rhetorical relations and concrete textual content required for neural retrieval. We develop a semantic enhancement framework through three key stages.

Bottom-up Semantic Enhancement. The document-level tree T_{doc} initially lacks concrete semantic content at internal nodes. We apply the same bottom-up LLM enhancement strategy (Equation 1) to generate meaningful textual representations for these structural nodes. Starting from leaf nodes and moving upward, this process systematically transforms each internal node from

abstract structural placeholders into semantically rich representations that capture the combined content of their subtrees. The result is a fully enhanced tree where all nodes possess concrete textual content suitable for semantic matching.

Multi-Level Tree Integration. We implement a substitution-based integration mechanism that combines the enhanced document-level tree T_{doc}^* with paragraph-level trees T_i into a unified discourse structure. Each leaf node in T_{doc}^* is systematically replaced with its corresponding paragraph-level discourse tree:

$$T_{\text{D}} = T_{\text{doc}}^*[l_i \leftarrow T_i, \forall i \in \{1, 2, \dots, n\}]. \quad (2)$$

This integration enables multi-granularity discourse modeling, where fine-grained intra-paragraph sentence relationships coexist with coarse-grained inter-paragraph discourse connections within a single hierarchical structure.

Node Encoding. To enable efficient retrieval, we transform the integrated discourse tree T_{D} into dense embeddings using a pre-trained encoder:

$$\mathbf{e}_v = f_{\text{enc}}(v), \quad \forall v \in \mathcal{N}(T_{\text{D}}). \quad (3)$$

The resulting embedding tree $\mathcal{T}_{\text{emb}} = (T_{\text{D}}, \{\mathbf{e}_v\})$ preserves both the hierarchical discourse structure and semantic information, enabling structure-aware retrieval through dense vector operations.

3.3 Structure-Guided Evidence Retrieval

We introduce a structure-aware retrieval mechanism that leverages the hierarchical nature of discourse trees to select evidence at multiple granularity levels. Unlike flat retrieval methods that operate on uniform text segments, our approach exploits discourse structure to implement a dual-selection strategy that balances local semantic relevance with global discourse coherence.

Retrieval Process. Given a query q , we first compute semantic similarities for all nodes in the embedding tree:

$$\text{score}(v) = \cos(f_{\text{enc}}(q), \mathbf{e}_v), \forall v \in \mathcal{N}(T_{\text{D}}). \quad (4)$$

Dual-Selection Strategy. Nodes are ranked by relevance and processed through two complementary selection mechanisms:

Algorithm 1: Structure-Guided Evidence Retrieval

Input: Query q , discourse tree T , evidence limit K , subtree selection size k
Output: Evidence set E
/* Compute node similarities and rank */
 $\mathbf{e}_q \leftarrow \text{Encoder}(q)$;
 $\text{scores}[v] \leftarrow \text{cosine}(\mathbf{e}_q, \mathbf{e}_v)$ for all $v \in T$;
 $V_{\text{ranked}} \leftarrow \text{sort}(\text{scores}, \text{descending})$;
/* Structure-aware selection */
 $E \leftarrow \{\}$, $\text{used} \leftarrow \{\}$;
for $v \in V_{\text{ranked}}$ **do**
 if v is leaf and $v \notin \text{used}$ **then**
 $E \leftarrow E \cup \{v\}$, $\text{used} \leftarrow \text{used} \cup \{v\}$;
 else
 $L \leftarrow \text{leaves}(v) \setminus \text{used}$; // Unused leaves in subtree
 if $L \neq \emptyset$ **then**
 $L_k \leftarrow \text{top}_k(L, \text{scores}, k)$;
 $E \leftarrow E \cup L_k$, $\text{used} \leftarrow \text{used} \cup L_k$;
 end
 end
 if $|E| \geq K$ **then**
 break
 end
end
return E

- **Direct leaf selection** for high-relevance sentences
- **Hierarchical expansion** with top- k subtree selection for internal nodes
- **Redundancy elimination** to prevent duplicates

For leaf nodes with high relevance, we select them directly. For internal nodes, we perform controlled subtree expansion by selecting the top- k most relevant unused leaves from their subtrees. This ensures both specific sentence-level evidence and coherent discourse segments are captured while eliminating redundancy.

Algorithm 1 details the complete process, balancing semantic relevance with discourse coherence. A concrete example illustrating how discourse structure guides evidence selection across granularity levels is provided in Appendix C.5.

4 Experiments

4.1 Experimental Setup

Datasets and Evaluation Metrics. We evaluate our approach on four challenging long document QA datasets: QASPER (Dasigi et al., 2021) for research papers (avg. 4170 words), QuALITY (Pang et al., 2022) for reading comprehension (avg. 5022 words), NarrativeQA (Kočíský et al., 2018) for narrative documents (avg. 51,372 words), and MultiFieldQA-zh (Bai et al., 2024) for Chinese documents (avg. 6701 words). We use F1-Match for QASPER and MultiFieldQA-zh, accuracy for QuALITY, and BLEU(B-1)/ROUGE/METEOR for NarrativeQA. Token-level F1/Recall evaluate retrieval quality on QASPER. We fix the retrieved

context length across methods for fair comparison.

Baselines. We evaluate our approach against several strong retrieval baselines:

- Flatten-chunk splits articles into chunks of maximum 100 words while preserving sentence boundaries for semantic coherence.
- Flatten-sentence adopts sentence-level splitting and direct retrieval, providing a direct comparison baseline for our hierarchical approach.
- RAPTOR constructs a semantic tree through recursive embedding, clustering, and summarization, with retrieval performed on a collapsed tree structure following Sarthi et al. (2024).
- Bisection shares our LLM-enhanced representations and hierarchical retrieval mechanism (Sections 3.2 and 3.3) but constructs trees by recursively dividing sentences into balanced binary subtrees. This isolates the specific contribution of discourse structure by keeping all other components identical to DISRetrieval.

Implementation details. We train our language-universal discourse parser on RST-DT (Carlson et al., 2001) following Yu et al. (2022b) with gte-multilingual-base¹ backbone (Zhang et al., 2024), augmented with GPT-4o-translated Chinese data. All summarization tasks, including both our approach (as defined in Equation 1) and RAPTOR baseline, utilize Llama3.1-8B-Instruct for consistency. We set threshold $\tau = 0$ for QASPER and $\tau = 50$ for QuALITY/NarrativeQA (Appendix C.4). For sentence encoder (Section §3.2), we adopt a completely training-free design, directly applying pre-trained embedding models without any fine-tuning to ensure zero training cost and plug-and-play flexibility across different deployment settings. Specifically, we test SentenceBERT² and OpenAI text-embedding-3-large³. Top-K is set to 5 (Appendix C.3). Generation models include UnifiedQA-3B, GPT-4.1-mini, and Deepseek-v3. Details in Appendix B.

4.2 Main Results

Generation Performance. We evaluate our approach across different context lengths (200-400 words) and embedding models (SBERT and OpenAI) on QASPER and QuALITY datasets. Table 1

¹huggingface.co/Alibaba-NLP/gte-multilingual-base

²multi-qa-mpnet-base-cos-v1

³https://platform.openai.com/docs/guides/embeddings

Method	F1-Match (QASPER) / %						Accuracy (QuALITY) / %					
	200 (s)	200 (o)	300 (s)	300 (o)	400 (s)	400 (o)	200 (s)	200 (o)	300 (s)	300 (o)	400 (s)	400 (o)
UnifiedQA-3B												
flatten-chunk	33.97	37.50	35.41	38.46	36.13	39.03	52.97	56.28	54.46	57.53	55.23	57.62
flatten-sentence	35.16	<u>38.43</u>	37.24	39.31	37.99	<u>40.06</u>	53.16	56.04	54.55	56.71	55.27	57.38
RAPTOR	33.57	37.46	34.95	<u>39.96</u>	37.00	39.53	53.60	55.99	54.75	57.00	55.51	58.53
<i>Ours</i>												
Bisection	<u>36.17</u>	37.41	<u>37.66</u>	39.49	<u>38.84</u>	39.70	<u>55.13</u>	<u>57.00</u>	<u>56.04</u>	<u>58.53</u>	<u>57.24</u>	<u>60.21</u>
DISRetrieval	37.36	38.49	39.56	40.03	40.65	40.74	55.56	57.67	57.62	59.64	58.87	60.64
GPT-4.1-mini												
flatten-chunk	37.37	41.13	41.38	43.34	42.72	44.78	60.16	65.77	63.66	69.27	67.69	71.05
flatten-sentence	39.82	43.08	41.84	<u>45.28</u>	42.44	<u>45.78</u>	60.12	64.43	63.09	66.68	65.24	69.94
RAPTOR	37.55	40.88	39.95	43.26	42.50	43.85	61.31	64.77	63.81	67.79	67.26	70.71
<i>Ours</i>												
Bisection	<u>40.98</u>	<u>43.12</u>	<u>42.74</u>	44.54	<u>43.49</u>	45.69	<u>63.71</u>	<u>67.88</u>	<u>65.68</u>	<u>70.71</u>	<u>67.93</u>	<u>72.00</u>
DISRetrieval	40.99	43.45	43.01	45.19	44.95	46.31	64.57	69.27	66.63	72.44	69.37	73.54
Deepseek-v3												
flatten-chunk	35.29	38.14	37.99	40.78	40.57	42.26	65.68	<u>71.52</u>	69.65	<u>75.02</u>	73.25	76.56
flatten-sentence	36.51	<u>39.27</u>	39.41	41.57	40.82	43.04	64.14	68.74	68.41	72.24	69.65	73.63
RAPTOR	34.66	37.84	37.30	40.50	39.77	42.17	65.53	69.13	68.65	72.24	71.19	75.22
<i>Ours</i>												
Bisection	<u>37.28</u>	<u>39.27</u>	<u>39.80</u>	<u>41.81</u>	<u>40.96</u>	<u>43.57</u>	<u>67.83</u>	71.09	<u>70.81</u>	74.40	<u>73.39</u>	<u>76.94</u>
DISRetrieval	37.79	39.77	40.39	42.19	41.51	43.65	68.89	72.28	72.15	76.46	73.68	77.71

Table 1: Generation performance comparison of different methods under varying retrieved context lengths (200-400 words) and different embedding models (SBERT: s, OpenAI text-embedding-3-large: o) across two datasets. Bisection shares our LLM-enhanced node representations and hierarchical retrieval mechanism but uses binary tree construction instead of discourse structure. Best results are **bolded**, runners-up are underlined.

presents the comparative results. The results reveal three key insights:

Consistent superiority across settings. DISRetrieval outperforms all baselines regardless of context length, embedding model, or generation architecture (UnifiedQA-3B, GPT-4.1-mini, and Deepseek-v3). For instance, with 400-word contexts and UnifiedQA-3B, we achieve +2.66% F1-Match on QASPER and +3.60% accuracy on QuALITY over the flatten-sentence baseline.

Discourse structure surpasses semantic clustering. Compared to RAPTOR’s semantic clustering approach, our discourse-aware method demonstrates clear advantages (40.03% vs. 39.96% F1-Match on QASPER with OpenAI embeddings), confirming that linguistic discourse structure provides more principled document organization than purely semantic-based methods.

Linguistic discourse essential for hierarchical modeling. The Bisection ablation validates our core hypothesis: while hierarchical organization helps (Bisection > flatten baselines), incorporating discourse structure is crucial (DISRetrieval > Bisection consistently), providing substantial benefits beyond simple tree-based document modeling.

Retrieval Performance. We evaluate DISRetrieval’s retrieval effectiveness on QASPER using token-level F1 and Recall metrics. The Retrieval Result in Table 2 reveals three key findings:

First, DISRetrieval consistently outperforms baselines across all settings, achieving the highest F1 and Recall scores with both SBERT and OpenAI embeddings. This demonstrates that our discourse-aware context modeling effectively captures the semantic relationships within documents.

Second, for longer contexts (300-400 words), while all methods show some F1 score degradation, DISRetrieval maintains superior performance, particularly in Recall metrics. This robust performance on longer contexts validates the capability of our method in handling complex document structures through discourse-guided retrieval.

The ablation with Bisection shows that while hierarchical organization helps, the full discourse-aware approach provides additional benefits. Additionally, OpenAI embeddings consistently outperform SBERT across all settings, suggesting that strong semantic representations are fundamental to the effectiveness of discourse-aware retrieval.

Cross-Lingual Effectiveness. To verify the language-universal capability of our approach, we

Retrieval Result	200 (SBERT)		200 (OpenAI)		300 (SBERT)		300 (OpenAI)		400 (SBERT)		400 (OpenAI)	
	F1 / %	Recall / %	F1 / %	Recall / %	F1 / %	Recall / %	F1 / %	Recall / %	F1 / %	Recall / %	F1 / %	Recall / %
flatten-chunk	26.13	56.05	29.17	62.16	23.10	63.92	25.12	69.04	20.38	68.75	21.91	73.38
flatten-sentence	26.15	57.80	28.25	62.78	22.68	64.63	24.04	68.43	19.98	69.02	21.06	72.42
RAPTOR	24.57	52.42	27.18	58.49	21.71	60.00	23.57	65.63	19.27	65.11	20.64	70.04
<i>Ours</i>												
Bisection	<u>27.63</u>	<u>59.82</u>	<u>29.29</u>	<u>63.11</u>	<u>23.83</u>	<u>66.69</u>	<u>25.16</u>	<u>69.94</u>	<u>21.10</u>	<u>71.52</u>	<u>21.98</u>	<u>74.18</u>
DISRetrieval	28.13	60.75	30.27	65.33	24.62	67.98	26.00	71.71	21.58	72.61	22.79	75.95

Generation Result	Full Document + Retrieved Evidence						Full Document	Gold Evidence
	200 (SBERT)	200 (OpenAI)	300 (SBERT)	300 (OpenAI)	400 (SBERT)	400 (OpenAI)	Avg. 4170 words	Avg. 129 words
F1 / %	49.65 (+0.84)	50.05 (+1.24)	49.77 (+0.96)	50.11 (+1.30)	49.54 (+0.73)	50.20 (+1.39)	48.81	50.71 (+1.90)

Table 2: Retrieval and generation results (token-level F1 and Recall) on the QASPER dataset. Embedding models: SBERT and OpenAI text-embedding-3-large. Generation model: Llama3.1-8B-Instruct.

evaluate on MultiFieldQA-zh, a Chinese long-document QA benchmark. As shown in Table 3, DISRetrieval consistently outperforms all baselines across different context lengths with both GPT-4.1-mini (35.25% F1) and Deepseek-v3 (29.54% F1). The improvements are particularly notable at longer contexts, where our method gains +1.30% to +1.56% over Bisection with Deepseek-v3. These results validate that linguistically-grounded discourse structures transcend language boundaries, with our language-universal parser effectively capturing discourse relationships in Chinese documents. This cross-lingual effectiveness highlights the language-agnostic nature of rhetorical structure theory and demonstrates the extensibility to other languages through similar data augmentation strategies.

Method	GPT-4.1-mini (F1 / %)			Deepseek-v3 (F1 / %)		
	200	300	400	200	300	400
flatten-chunk	<u>28.77</u>	31.25	33.46	22.69	26.91	26.70
flatten-sentence	28.59	31.95	34.41	26.15	27.04	27.06
RAPTOR	26.61	31.70	34.32	23.11	25.96	27.01
<i>Ours</i>						
Bisection	28.23	<u>32.31</u>	<u>34.80</u>	<u>26.71</u>	<u>27.05</u>	<u>28.24</u>
DISRetrieval	29.60	32.97	35.25	26.76	28.61	29.54

Table 3: Performance on MultiFieldQA-zh (Chinese).

	BLEU / %	ROUGE / %	METEOR / %
flatten-chunk	24.24	28.42	19.70
flatten-sentence	21.53	28.22	18.33
RAPTOR	<u>25.05</u>	<u>30.24</u>	<u>20.92</u>
<i>Ours</i>			
Bisection	24.71	28.85	20.41
DISRetrieval	25.39	30.31	21.16

Table 4: Performance on the NarrativeQA dataset.

4.3 Discussions

RQ1: Is DISRetrieval effective for handling extremely long documents? **Tab 4** We evaluate our method on NarrativeQA, which contains exceptionally long documents (average: 51,372 words, maximum: 346,902 words) that exceed most generative models’ context limits. DISRetrieval achieves superior performance across all metrics, outperforming the widely-used flatten-chunk baseline by +1.15% BLEU, +1.89% ROUGE, and +1.46% METEOR. Notably, both DISRetrieval and RAPTOR substantially outperform flatten-based methods, confirming that structured approaches are more effective than flat for extremely long document retrieval. However, DISRetrieval’s discourse-aware structure provides consistent advantages over RAPTOR’s semantic clustering approach.

RQ2: Is precise evidence retrieval essential for effective question answering? **Tab 2** To evaluate the importance of precise evidence retrieval, we compare three context settings using Llama3.1-8B-Instruct: (i) golden evidence, (ii) full document, and (iii) full document augmented with retrieved evidence (highlighted with [EVIDENCE] ... [/EVIDENCE] markers). Golden evidence achieves the highest F1 score of 50.71% with only 129 words on average, significantly outperforming the full document approach (48.81% with 4,170 words). This demonstrates that concise yet accurate retrieval is more effective than using substantially larger context. Furthermore, augmenting the full document with our retrieved evidence yields consistent improvements of 0.73-1.39% across all settings, confirming that precise retrieval enhances performance even when combined with complete document access. These results underscore that effective long-document QA systems should prioritize retrieval

quality over context quantity.

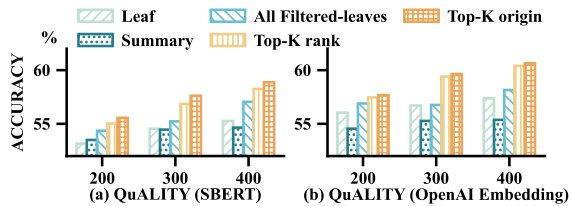


Figure 3: Ablation results of different variants.

RQ3: Is the hierarchical retrieval strategy effective? **Fig 3** We conduct ablation experiments on the QuALITY dataset to evaluate our hierarchical retrieval strategy. We compare five variants: leaf-only baseline, summary-based retrieval, all filtered-leaves, Top-K with ranking order, and our Top-K with original order. Three key findings emerge: (1) Summary-based retrieval underperforms the leaf baseline, confirming that preserving original text details is crucial for effective retrieval. (2) While all filtered-leaves shows marginal improvements, selective Top-K methods are superior due to reduced noise from irrelevant content. (3) Our Top-K origin method consistently achieves the best performance by preserving natural document flow. This advantage intensifies with longer contexts and stronger embeddings, validating that maintaining document structure is essential for hierarchical retrieval. Notably, the proportion of internal nodes versus leaf nodes in retrieved evidence varies substantially across document types: for narrative documents with short sentences such as QuALITY, the system predominantly selects internal nodes at deeper levels, whereas for academic papers with longer self-contained sentences such as QASPER, it relies more on leaf-level nodes. This adaptive behavior across granularity levels is the key mechanism driving performance gains, and we provide a detailed quantitative analysis in Appendix C.2.

RQ4: Does the scale of LLMs affect the quality of the discourse-aware tree structures in node text enhancement? **Tab 5** We investigate whether model scale affects discourse-aware tree construction by comparing different LLMs. As shown in Table 5, smaller models like Llama-3.1-8B, Qwen2.5-7B, and Mistral-7B achieve comparable performance to larger models, with differences under 0.5%. Notably, Llama-3.1-8B achieves the best recall at 400-word context length. These results indicate that our approach does not heavily

depend on LLM scale, enabling low-cost deployment with 7B models. Moreover, our method is $3\times$ faster than RAPTOR (e.g., 50K words: 103s vs. 338s), with preprocessing costs amortized across multiple queries (detailed in Appendix C.1).

	200		300		400	
	F1 / %	Recall / %	F1 / %	Recall / %	F1 / %	Recall / %
Llama-3.1-8B	28.13	60.75	24.62	67.98	21.58	72.61
Qwen2.5-7B	28.54	61.08	24.73	68.13	21.76	72.71
Mistral-7B	28.51	61.57	24.68	68.53	21.54	72.50
GPT-4o-mini	28.29	61.63	24.65	68.42	21.49	72.65
Deepseek-v3	28.66	61.80	24.77	68.66	21.72	72.69

Table 5: Effect of different LLMs for node text enhancement on QASPER dataset retrieval performance (token-level F1 and Recall) across varying context lengths.

RQ5: Does the capability of the discourse parser have a significant impact? **Fig 4** We examine how discourse parser capabilities affect downstream performance by training parsers on varying amounts of data (0-100%). As shown in Figure 4, both retrieval recall and answer F1 scores improve consistently as parser training data increases across all context lengths. Notably, despite our parser being trained only on RST-DT (news-text-dominant), DISRetrieval consistently outperforms all baselines across diverse genres (research papers, fiction, movie scripts), demonstrating practical robustness and extensible potential as parsing technology advances.

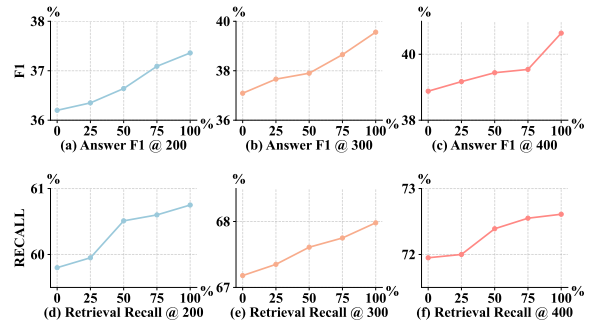


Figure 4: Impact of discourse parser capability on subsequent retrieval and question answering performance. Parsers used for comparative evaluation are trained on varying data scales, ranging from 0% to 100%.

RQ6: What drives the performance improvement, relation type labels or tree structure? **Tab 6** To clarify the mechanism underlying DISRetrieval’s improvements, we analyze RST relation type distributions between full documents and retrieved Top-20 nodes on QASPER, as shown in

Table 6. The near-identical distributions across all categories confirm that relation type labels carry no discriminative signal for retrieval correctness, as labels such as *Elaboration* and *Contrast* encode structural connectivity rather than semantic relevance to a given query. The true driver is the hierarchical grouping structure induced by RST parsing, which identifies sentences that naturally form coherent semantic units across granularity levels. This is validated by two observations: DISRetrieval consistently outperforms the Bisection baseline that replaces discourse boundaries with arbitrary binary splits while keeping all other components identical, and parser quality directly correlates with downstream QA performance as shown in Figure 4. For retrieval, knowing that a set of sentences forms a meaningful discourse unit is fundamentally more informative than knowing the rhetorical label of their connection.

	Elab.	List	Attr.	Evid.	Contr.	Others
Full Documents	40.95%	10.14%	9.18%	5.95%	5.54%	28.24%
Top-20 Nodes	40.70%	10.28%	8.79%	5.91%	5.77%	28.55%

Table 6: Distribution of RST relation types in full documents versus retrieved Top-20 nodes on QASPER. The near-identical distributions indicate that relation type labels do not correlate with retrieval correctness.

5 Conclusion

In this paper, we presented DISRetrieval, a discourse-aware hierarchical retrieval framework that systematically incorporates rhetorical structure theory into long document question answering. Our approach features three key innovations: language-universal discourse parsing, LLM-enhanced node representations, and structure-aware evidence selection. Comprehensive experiments on QASPER, QuALITY, NarrativeQA, and MultiFieldQA-zh demonstrate substantial improvements over existing methods across varying context lengths, embedding models, generation architectures, and languages. Ablation studies confirm that discourse structure provides more principled document organization than semantic clustering approaches, with the framework naturally adapting to different document types and information needs through multi-granularity retrieval. Our cross-lingual experiments validate that linguistically-grounded discourse structures effectively transcend language boundaries, demonstrating robust generalization from English to Chinese documents. Beyond per-

formance improvements, this work demonstrates that linguistically-grounded approaches remain valuable in the era of LLMs, offering principled alternatives to purely data-driven methods. Our framework opens new directions for incorporating structured linguistic knowledge into neural retrieval systems, with potential applications in legal document analysis, scientific literature review, and educational content processing.

Limitations

While DISRetrieval demonstrates substantial improvements across multiple datasets and architectures, we acknowledge several limitations that point to promising future directions. First, our discourse parser’s performance inherently bounds the overall system effectiveness, as shown in Figure 4, where parser quality directly impacts downstream performance. However, our language-universal parser, trained with LLM-based multilingual data augmentation, demonstrates effective cross-lingual generalization from English to Chinese and consistent cross-genre improvements (research papers, fiction, movie scripts), showing practical robustness and extensible potential as parsing technology advances. While we currently demonstrate effectiveness on English and Chinese, extending to more languages would require additional translated training data through similar augmentation strategies. We note that the scarcity of suitable long-document QA datasets in other languages currently limits broader multilingual evaluation. Future work will explore expanding language coverage through collaboration with multilingual NLP communities, developing domain-adaptive parsers, and investigating multi-granularity discourse analysis. Second, our adaptive summarization strategy using threshold τ is relatively simple but proves effective across diverse document types ($\tau=0$ for academic papers, $\tau=50$ for narratives). While this approach successfully balances computational efficiency with representation quality, more sophisticated dynamic thresholding based on content complexity and hierarchical position could further optimize the trade-off. Third, current evaluation metrics, though comprehensive across four challenging datasets spanning multiple languages, may not fully capture the nuanced benefits of discourse-aware retrieval, such as structural coherence preservation and hierarchical information flow. We plan to develop evaluation frameworks that better assess

discourse-aware retrieval quality across multiple dimensions and languages. Fourth, from a task applicability perspective, DISRetrieval is most effective for long-document evidence retrieval and question answering tasks that require assembling multiple evidence pieces, as demonstrated by consistent gains across structured documents such as scientific papers, unstructured narratives such as stories and movie scripts, extremely long documents, and cross-lingual scenarios in our experiments. The method is less suited to short-document question answering or multi-document retrieval settings, where the structural advantages of discourse hierarchy are less pronounced and the overhead of discourse tree construction may outweigh the retrieval benefits. We regard extending the framework to these settings as a direction for future work.

Ethical Considerations

This research adheres to ethical principles in natural language processing research and does not raise significant ethical concerns. Our work focuses on improving document retrieval and question answering systems through discourse-aware hierarchical modeling, which has potential positive societal impacts by enhancing information access and comprehension across multiple languages. The datasets used in our experiments (QASPER, QuALITY, NarrativeQA, and MultiFieldQA-zh) are publicly available benchmark datasets that have been previously vetted by the research community and do not contain sensitive personal information or harmful content. Our discourse parsing approach operates on linguistic structures rather than content semantics, minimizing risks of bias amplification or misuse. The language-universal parser is trained using LLM-based translation for data augmentation, which may inherit potential biases from the translation model, though our focus on structural discourse relations rather than semantic content helps mitigate such risks. The proposed DISRetrieval framework is designed as a general-purpose retrieval enhancement technique that can benefit various applications requiring long document understanding across languages, such as academic research assistance, educational content processing, and legal document analysis. We acknowledge that like any information retrieval system, our method could potentially be misused if applied to spread misinformation or manipulate access to information, but such concerns are not specific to our

approach and apply broadly to information retrieval technologies. We encourage responsible deployment of our system with appropriate safeguards and human oversight, particularly in high-stakes applications and cross-lingual scenarios. All experiments were conducted using publicly available computational resources and open-source tools, ensuring reproducibility and transparency in our research process.

Acknowledgements

We thank the anonymous reviewers and chairs for their helpful comments. This work is supported by National Natural Science Foundation of China (Grant No. 62336008), Shenzhen Basic Research Program (Grant No. JCYJ20241202123503005), and Shenzhen Basic Research Program (Grant No. JCYJ20240813105111016).

References

- Joshua Ainslie, James Lee-Thorp, Michiel de Jong, Yury Zemlyanskiy, Federico Lebron, and Sumit Sanghai. 2023. GQA: Training generalized multi-query transformer models from multi-head checkpoints. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4895–4901.
- Manoj Ghuhan Arivazhagan, Lan Liu, Peng Qi, Xinchu Chen, William Yang Wang, and Zhiheng Huang. 2023. Hybrid hierarchical retrieval for open-domain question answering. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 10680–10689.
- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024. LongBench: A bilingual, multi-task benchmark for long context understanding. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3119–3137.
- Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- Lynn Carlson, Daniel Marcu, and Mary Ellen Okurovsky. 2001. Building a discourse-tagged corpus in the framework of Rhetorical Structure Theory.

- In *Proceedings of the Second SIGdial Workshop on Discourse and Dialogue*.
- Danqi Chen, Adam Fisch, Jason Weston, and Antoine Bordes. 2017. Reading Wikipedia to answer open-domain questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1870–1879.
- Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2023. Palm: Scaling language modeling with pathways. *Journal of Machine Learning Research*, 24(240):1–113.
- Arman Cohan, Franck Dernoncourt, Doo Soon Kim, Trung Bui, Seokhwan Kim, Walter Chang, and Nazli Goharian. 2018. A discourse-aware attention model for abstractive summarization of long documents. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 2 (Short Papers)*, pages 615–621.
- Pradeep Dasigi, Kyle Lo, Iz Beltagy, Arman Cohan, Noah A. Smith, and Matt Gardner. 2021. A dataset of information-seeking questions and answers anchored in research papers. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4599–4610.
- Jiayu Ding, Shuming Ma, Li Dong, Xingxing Zhang, Shaohan Huang, Wenhui Wang, Nanning Zheng, and Furu Wei. 2023. Longnet: Scaling transformers to 1,000,000,000 tokens. *arXiv preprint arXiv:2307.02486*.
- Haowei Du, Yansong Feng, Chen Li, Yang Li, Yunshi Lan, and Dongyan Zhao. 2023. Structure-discourse hierarchical graph for conditional question answering on long documents. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 6282–6293.
- Vanessa Wei Feng and Graeme Hirst. 2014. A linear-time bottom-up discourse parser with constraints and post-editing. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 511–521.
- Mandy Guo, Joshua Ainslie, David Uthus, Santiago Ontanon, Jianmo Ni, Yun-Hsuan Sung, and Yinfei Yang. 2022. LongT5: Efficient text-to-text transformer for long sequences. In *Findings of the Association for Computational Linguistics: NAACL 2022*, pages 724–736.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Ming-Wei Chang. 2020a. REALM: retrieval-augmented language model pre-training. *CoRR*.
- Kelvin Guu, Kenton Lee, Zora Tung, Panupong Pasupat, and Mingwei Chang. 2020b. Retrieval augmented language model pre-training. In *International conference on machine learning*, pages 3929–3938.
- Maor Ivgi, Uri Shaham, and Jonathan Berant. 2023. Efficient long-text understanding with short-text models. *Transactions of the Association for Computational Linguistics*, 11:284–299.
- Gautier Izacard and Edouard Grave. 2020. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282*.
- Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-Yu, Armand Joulin, Sebastian Riedel, and Edouard Grave. 2023. Atlas: Few-shot learning with retrieval augmented language models. *J. Mach. Learn. Res.*, 24:251:1–251:43.
- Jiajie Jin, Xiaoxi Li, Guanting Dong, Yuyao Zhang, Yutao Zhu, Yongkang Wu, Zhonghua Li, Qi Ye, and Zhicheng Dou. 2025. Hierarchical document refinement for long-context retrieval-augmented generation. *arXiv preprint arXiv:2505.10413*.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781.
- Naoki Kobayashi, Tsutomu Hiraio, Hidetaka Kamigaito, Manabu Okumura, and Masaaki Nagata. 2020. Top-down rst parsing utilizing granularity levels in documents. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 8099–8106.
- Tomáš Kočiský, Jonathan Schwarz, Phil Blunsom, Chris Dyer, Karl Moritz Hermann, Gábor Melis, and Edward Grefenstette. 2018. The NarrativeQA reading comprehension challenge. *Transactions of the Association for Computational Linguistics*, 6:317–328.
- 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.
- Haijun Li, Tianqi Shi, Zifu Shang, Yuxuan Han, Xueyu Zhao, Hao Wang, Yu Qian, Zhiqiang Qian, Linlong Xu, Minghao Wu, et al. 2025. Transbench: Benchmarking machine translation for industrial-scale applications. *arXiv preprint arXiv:2505.14244*.
- Yang Liu and Mirella Lapata. 2019. Hierarchical transformers for multi-document summarization. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5070–5081.
- Ye Liu, Kazuma Hashimoto, Yingbo Zhou, Semih Yavuz, Caiming Xiong, and Philip Yu. 2021. Dense hierarchical retrieval for open-domain question answering. In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 188–200.

- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized BERT pretraining approach. *CoRR*, abs/1907.11692.
- Aru Maekawa, Tsutomu Hirao, Hidetaka Kamigaito, and Manabu Okumura. 2024. Can we obtain significant success in RST discourse parsing by using large language models? In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2803–2815.
- Sadhika Malladi, Tianyu Gao, Eshaan Nichani, Alex Damian, Jason D Lee, Danqi Chen, and Sanjeev Arora. 2023. Fine-tuning language models with just forward passes. *Advances in Neural Information Processing Systems*, 36:53038–53075.
- William C Mann and Sandra A Thompson. 1988. Rhetorical structure theory: Toward a functional theory of text organization. *Text-interdisciplinary Journal for the Study of Discourse*, 8(3):243–281.
- Daniel Marcu. 2000. *The theory and practice of discourse parsing and summarization*. MIT press.
- Inderjeet Nair, Shwetha Somasundaram, Apoorv Saxena, and Koustava Goswami. 2023. Drilling down into the discourse structure with llms for long document question answering. *arXiv preprint arXiv:2311.13565*.
- Richard Yuanzhe Pang, Alicia Parrish, Nitish Joshi, Nikita Nangia, Jason Phang, Angelica Chen, Vishakh Padmakumar, Johnny Ma, Jana Thompson, He He, and Samuel Bowman. 2022. QuALITY: Question answering with long input texts, yes! In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 5336–5358.
- Hongjin Qian, Zheng Liu, Kelong Mao, Yujia Zhou, and Zhicheng Dou. 2024. Grounding language model with chunking-free in-context retrieval. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1298–1311.
- Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, Chloe Hillier, and Timothy P. Lillicrap. 2020. Compressive transformers for long-range sequence modelling. In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392.
- Siva Reddy, Danqi Chen, and Christopher D. Manning. 2019. CoQA: A conversational question answering challenge. *Transactions of the Association for Computational Linguistics*, 7:249–266.
- Soumya Santhosh and Jahfar Ali. 2012. Discourse based advancement on question answering system. *Journal on Soft Computing*.
- Parth Sarthi, Salman Abdullah, Aditi Tuli, Shubh Khanna, Anna Goldie, and Christopher D. Manning. 2024. RAPTOR: recursive abstractive processing for tree-organized retrieval. In *The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024*.
- Uri Shaham, Elad Segal, Maor Ivgi, Avia Efrat, Ori Yoran, Adi Haviv, Ankit Gupta, Wenhan Xiong, Mor Geva, Jonathan Berant, and Omer Levy. 2022. SCROLLS: Standardized CompaRison over long language sequences. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 12007–12021.
- Anders Søgaard, Sebastian Ebert, Jasmijn Bastings, and Katja Filippova. 2021. We need to talk about random splits. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1823–1832.
- Yixuan Tang, Weilong Huang, Qi Liu, Anthony KH Tung, Xiaoli Wang, Jisong Yang, and Beibei Zhang. 2017. Qalink: enriching text documents with relevant q&a site contents. In *Proceedings of the 2017 ACM on Conference on Information and Knowledge Management*, pages 1359–1368.
- Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. 2021. Long range arena : A benchmark for efficient transformers. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*.
- Liang Wang, Nan Yang, Xiaolong Huang, Binxing Jiao, Linjun Yang, Daxin Jiang, Rangan Majumder, and Furu Wei. 2023. SimLM: Pre-training with representation bottleneck for dense passage retrieval. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2244–2258.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45.
- Jinyang Wu, Guocheng Zhai, Ruihan Jin, Jiahao Yuan, Yuhao Shen, Shuai Zhang, Zhengqi Wen, and Jianhua Tao. 2026. Atlas: Orchestrating heterogeneous models and tools for multi-domain complex reasoning. *arXiv preprint arXiv:2601.03872*.

- Jinyang Wu, Shuai Zhang, Feihu Che, Mingkuan Feng, Pengpeng Shao, and Jianhua Tao. 2025. Pandora’s box or aladdin’s lamp: A comprehensive analysis revealing the role of RAG noise in large language models. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5019–5039, Vienna, Austria. Association for Computational Linguistics.
- Fangyuan Xu, Junyi Jessy Li, and Eunsol Choi. 2022. How do we answer complex questions: Discourse structure of long-form answers. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 3556–3572.
- Nan Yu, Meishan Zhang, and Guohong Fu. 2018. Transition-based neural rst parsing with implicit syntax features. In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 559–570.
- Nan Yu, Meishan Zhang, Guohong Fu, and Min Zhang. 2022a. Rst discourse parsing with second-stage edu-level pre-training. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4269–4280.
- Nan Yu, Meishan Zhang, Guohong Fu, and Min Zhang. 2022b. RST discourse parsing with second-stage EDU-level pre-training. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4269–4280.
- Wenhao Yu, Dan Iter, Shuohang Wang, Yichong Xu, Mingxuan Ju, Soumya Sanyal, Chenguang Zhu, Michael Zeng, and Meng Jiang. 2023. Generate rather than retrieve: Large language models are strong context generators. In *The Eleventh International Conference on Learning Representations, ICLR 2023, Kigali, Rwanda, May 1-5, 2023*. OpenReview.net.
- Mengying Yuan, WenHao Wang, Zixuan Wang, Yujie Huang, Kangli Wei, Fei Li, Chong Teng, and Donghong Ji. 2025. Cross-document cross-lingual NLI via RST-enhanced graph fusion and interpretability prediction. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 31607–31629.
- Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. 2020. Big bird: Transformers for longer sequences. *Advances in neural information processing systems*, 33:17283–17297.
- Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter J. Liu. 2020. PEGASUS: pre-training with extracted gap-sentences for abstractive summarization. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, volume 119 of *Proceedings of Machine Learning Research*, pages 11328–11339. PMLR.
- Meishan Zhang, Xin Zhang, Xinping Zhao, Shouzheng Huang, Baotian Hu, and Min Zhang. 2025. On the role of pretrained language models in general-purpose text embeddings: A survey. *arXiv preprint arXiv:2507.20783*.
- Xin Zhang, Yanzhao Zhang, Dingkun Long, Wen Xie, Ziqi Dai, Jialong Tang, Huan Lin, Baosong Yang, Pengjun Xie, Fei Huang, et al. 2024. mgte: Generalized long-context text representation and reranking models for multilingual text retrieval. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing: Industry Track*, pages 1393–1412.
- Xinping Zhao, Xinshuo Hu, Zifei Shan, Shouzheng Huang, Yao Zhou, Xin Zhang, Zetian Sun, zhenyu liu, Dongfang Li, Xinyuan Wei, Youcheng Pan, Yang Xiang, Meishan Zhang, Haofen Wang, Jun Yu, Baotian Hu, and Min Zhang. 2026. KaLM-embedding-v2: Superior training techniques and data inspire a versatile embedding model. In *The Fourteenth International Conference on Learning Representations*.
- Xinping Zhao, Dongfang Li, Yan Zhong, Boren Hu, Yibin Chen, Baotian Hu, and Min Zhang. 2024. SEER: Self-aligned evidence extraction for retrieval-augmented generation. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 3027–3041.

A Extension of Technical Details

In this section, we introduce the specific details of our method which we cannot present in the main article due to the space limit.

A.1 Discourse Parser Details

The discourse parser architecture builds upon a transition-based system that constructs discourse trees through a series of well-defined actions. This system is particularly effective for capturing both local and global discourse relationships while maintaining computational efficiency.

The transition system maintains two primary data structures: a stack σ that holds partially built trees, and a queue β that contains unprocessed sentences. This design follows the intuition that discourse relationships often exist between adjacent text spans and can be built incrementally from bottom to top.

The system builds the tree step by step using three basic operations:

1. A “shift” action moves a sentence from the queue to the stack when we need new content to process.

Dataset	Avg. Sentence Length	Avg. Mid Node Depth		Avg. Leaf Num		Avg. Mid Node Percentage / %	
		SBERT	OpenAI	SBERT	OpenAI	SBERT	OpenAI
QASPER	22.77	5.37	5.44	12.89	13.61	54.95	66.02
QuALITY	14.64	13.98	17.34	68.84	87.92	74.86	89.95

Table 7: Statistical analysis of retrieved intermediate node characteristics across QuALITY and QASPER datasets. The table compares various key metrics: *Avg. Sentence Length* means the average sentence length of all documents; *Avg. Mid Node Depth* is the average depth of retrieved intermediate nodes; *Avg. Leaf Num* represents the average number of leaf nodes that each retrieved intermediate node maps to; *Avg. Mid Node Percentage* is the average percentage of intermediate nodes among the Top-20 retrieved nodes.

2. A “reduce” action combines two adjacent subtrees on top of the stack into a new subtree by identifying their discourse relationship.
3. A “pop root” action concludes the process when we have successfully built a complete tree.

Each state of the system is represented as $c = (\sigma, \beta)$, starting from $c_0 = ([], S_i)$ with all sentences in the queue, and ending at $c_f = ([T_i], [])$ with a complete discourse tree T_i .

The transition system follows a deterministic process guided by the neural scoring model:

1. Initialize $\sigma = []$ and $\beta = S_i$.
2. While β is not empty or $|\sigma| > 1$: (a) If $|\sigma| < 2$ and β is not empty, perform a “shift” action to move the next sentence from β to σ ; (b) Else if β is empty, perform a “reduce” action to combine the top two subtrees in σ ; (c) Else, use the neural scoring model to decide between a “shift” or “reduce” action based on the current state of σ and β .
3. Return the single tree T_i remaining on the stack σ .

The scoring model considers the three topmost subtrees on the stack (s_1, s_2, s_3) and the next sentence in the queue q_1 . This design is motivated by several factors:

1. s_1 and s_2 are the immediate candidates for the next potential "reduce" action.
2. s_3 provides crucial context about the recently built structure.
3. q_1 helps determine if we should introduce new content via a "shift" action.

For each tree node v , we compute its representation h_v recursively:

$$\mathbf{h}_v = \begin{cases} \text{PLM}(s_i), & \text{if } v \text{ is sentence} \\ \frac{1}{|C(v)|} \sum_{u \in C(v)} \mathbf{h}_u, & \text{if } v \text{ is relationship} \end{cases} \quad (5)$$

where $C(v)$ denotes the set of child nodes of v , and $\text{PLM}(\cdot)$ is a pre-trained language model that encodes the semantic meaning of individual sentences. The action scores are then computed as:

$$\mathbf{y}(a) = \mathbf{W}(\mathbf{h}_{s_1} \oplus \mathbf{h}_{s_2} \oplus \mathbf{h}_{s_3} \oplus \mathbf{h}_{q_1}) + \mathbf{b}, \quad (6)$$

where \oplus concatenates the representations to capture their interactions, and \mathbf{W} and \mathbf{b} are learnable parameters. The probability of taking action a is computed using a softmax function over the action scores:

$$p(a|c) = \frac{\exp(\mathbf{y}(a))}{\sum_{a' \in A} \exp(\mathbf{y}(a'))}, \quad (7)$$

where A is the set of valid actions at state c . We train the model using supervised learning with gold-standard discourse trees. The objective function combines cross-entropy loss for action prediction with L2 regularization:

$$\mathcal{L}(\theta) = -\log p(a^*|c) + \frac{\lambda \|\theta\|_2}{2}, \quad (8)$$

where a^* is the correct action derived from gold-standard trees. During inference, we greedily select the highest-scoring action at each step, effectively building the tree in a bottom-up manner while maintaining the discourse relationships between text spans.

A.2 Iterative Tree Construction Process

The complete algorithm of our DISRetrieval is presented in Algorithm 2. Below we detail the iterative tree construction process, which corresponds to Stage 1 of the algorithm. The construction of discourse trees for long documents presents

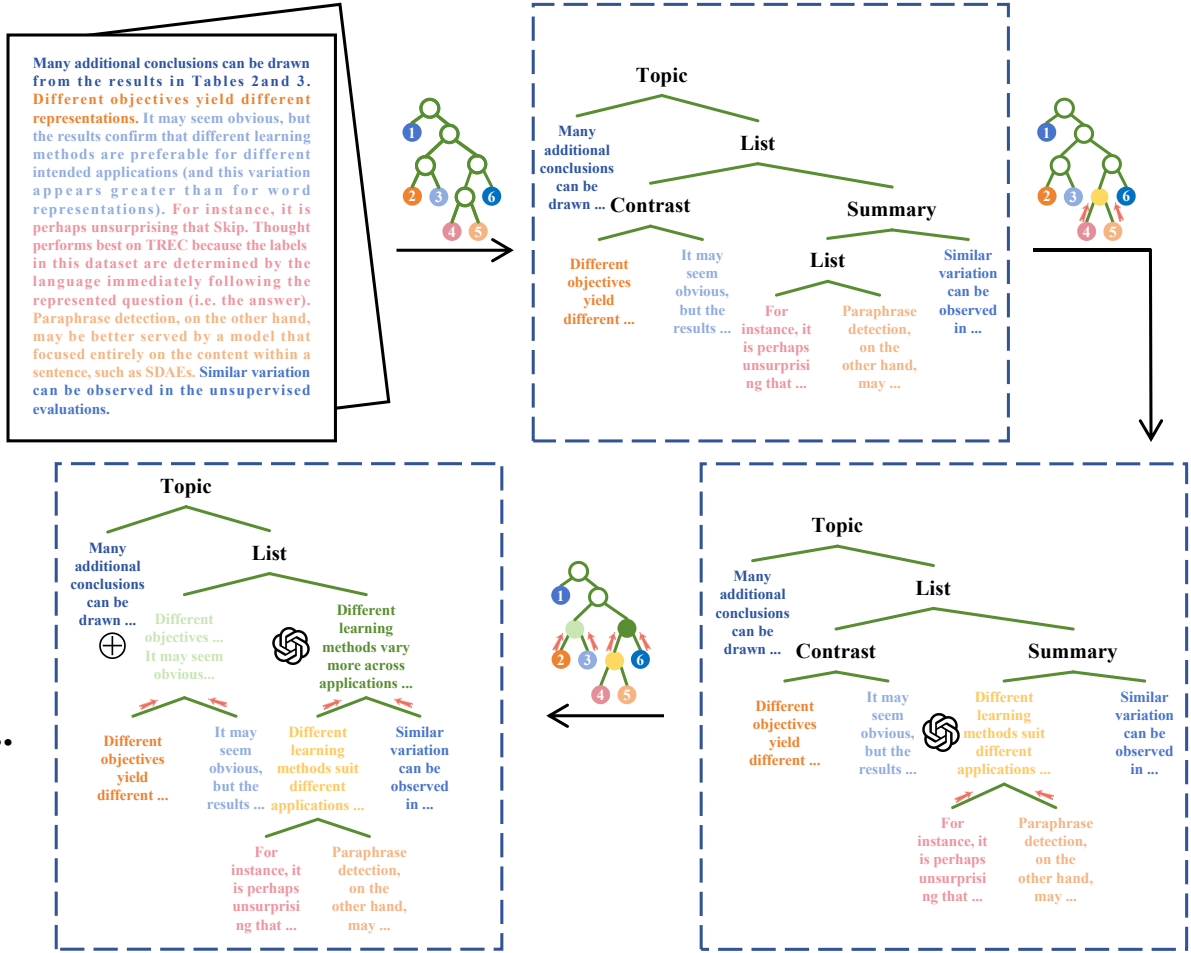


Figure 5: Illustration of bottom-up LLM enhancement in Phase 2 of discourse tree construction. Top: Input paragraph with its initial discourse tree structure. Center: Two-way processing strategy based on text length - using direct concatenation (\oplus) when combined length is below threshold τ , and LLM-based summarization when above threshold. Bottom: Enhanced discourse tree with progressively generated semantic summaries following Equation 1, demonstrating the transformation from rhetorical relations to concrete semantic representations.

unique challenges in balancing computational efficiency with structural integrity. We propose an iterative construction strategy that addresses these challenges through a hierarchical, phase-wise approach. This method effectively manages computational resources while preserving discourse relationships at multiple granularity levels. Our iterative process consists of three distinct phases, each designed to handle specific aspects of the tree construction:

Phase 1: Paragraph-level Tree Construction.

The first phase focuses on building local discourse trees for individual paragraphs. For each paragraph p_i containing sentences $S_i = \{s_{i,1}, s_{i,2}, \dots, s_{i,m}\}$, we construct a local discourse tree T_i using our transition-based parsing system. This phase is particularly efficient as it:

1. Initializes each paragraph’s parsing state with an empty stack and sentence queue: $c_0 =$

$([], S_i)$.

2. Processes paragraphs independently, enabling parallel computation.
3. Applies transition actions iteratively until a complete tree is formed.
4. Stores both the resulting paragraph-level tree T_i and its root representation h_{T_i} .

Phase 2: Document-level Tree Construction.

The second phase focuses on capturing document-level discourse structure. After obtaining all paragraph-level trees T_1, T_2, \dots, T_n , we:

1. For each paragraph-level tree T_i , apply bottom-up LLM-enhanced summarization:
 - For each non-leaf node v with children

v_l and v_r :

$$v^* = \begin{cases} f_{\text{LLM}}(v_l, v_r), & \text{if } |v_l| + |v_r| \geq \tau \\ f_{\text{merge}}(v), & \text{otherwise} \end{cases} \quad (9)$$

where $|v_l|$ and $|v_r|$ are the textual content length of child nodes

- Continue until reaching root node to obtain semantic unit u_i

2. Form the semantic units set $U = \{u_1, u_2, \dots, u_n\}$ from root representations
3. Apply the discourse parser to these units to construct a document-level tree T_{doc} using the same transition-based parsing system:

$$T_{\text{doc}} = f_{\text{discourse}}(U) \quad (10)$$

4. Apply bottom-up LLM-enhanced summarization to T_{doc} :

- For each non-leaf node $v \in T_{\text{doc}}$ with children v_l and v_r :

$$t_v = \begin{cases} f_{\text{LLM}}(t_l \oplus t_r), & \text{if } |t_l \oplus t_r| \geq \tau \\ t_l \oplus t_r, & \text{otherwise} \end{cases} \quad (11)$$

- Process nodes level by level from bottom to top until reaching the root of T_{doc}

This step effectively captures the high-level discourse relationships between paragraphs while maintaining computational efficiency by working with LLM-enhanced condensed representations at both paragraph and document levels.

Tree Integration. The final phase integrates the local and global structures into a unified discourse tree. This integration is accomplished through a careful replacement process:

$$T_{\text{D}} = \text{Replace}(T_{\text{doc}}, \{(u_i, T_i) | i \in [1, n]\}) \quad (12)$$

where $\text{Replace}(\cdot)$ replaces each unit u_i in the document-level tree with its corresponding paragraph-level tree T_i . This integration phase carefully preserves both local and global discourse relationships by maintaining the internal structure of paragraph-level trees while retaining the document-level relationships established in Phase 2, ultimately creating a seamless hierarchical structure that spans the entire document and effectively captures discourse relationships at all levels of granularity.

The resulting tree structure effectively captures discourse relationships at multiple levels of granularity, from sentence-level connections within paragraphs to broader document-level organizational patterns. This approach not only ensures computational efficiency through its phase-wise processing but also maintains the integrity of discourse relationships throughout the document hierarchy.

A.3 Detailed LLM Enhancement Process

This section provides a detailed illustration of the bottom-up LLM enhancement process described in Phase 2 of our discourse tree construction methodology (Section 3.1.2).

A.3.1 LLM Enhancement Workflow

Figure 5 demonstrates the complete workflow of our LLM enhancement process using a concrete example from a research paper paragraph. The process consists of three main components:

Input Structure. The top portion shows the input paragraph with its initial sentence-level discourse tree structure. Each sentence is represented as a leaf node, connected through discourse relations such as *Contrast*, *List*, *Summary*, and *Topic*. This hierarchical structure captures the rhetorical organization of the paragraph content.

Adaptive Processing Strategy. The center portion illustrates our two-way processing strategy based on text length thresholds. When the combined length of child nodes ($|t_l| + |t_r|$) falls below the threshold τ , we apply direct concatenation (\oplus) to preserve the original textual details. Conversely, when the combined length exceeds τ , we employ LLM-based summarization to generate more concise representations while retaining essential semantic information.

Enhanced Tree Structure. The bottom portion shows the resulting enhanced discourse tree with progressively generated semantic summaries. Each internal node now contains meaningful textual representations that capture the essence of its subtree content. For example, the *Contrast* relation between different learning methods is transformed into the concrete summary "Different learning methods suit different applications", while maintaining the hierarchical discourse structure.

A.3.2 Key Benefits

This LLM enhancement process provides several advantages:

- **Semantic Preservation:** Maintains essential meaning while reducing textual complexity
- **Computational Efficiency:** Adaptive thresholding prevents unnecessary LLM calls for short text segments
- **Hierarchical Coherence:** Preserves discourse relationships at multiple granularity levels
- **Downstream Compatibility:** Generates representations suitable for sentence-level discourse parsing

The transformation from abstract rhetorical relations to concrete semantic representations enables our framework to effectively bridge the gap between discourse structure and semantic understanding, facilitating improved retrieval performance in long document question answering tasks.

B Detailed Experiment Settings

B.1 Dataset Specifications

The QASPER dataset is constructed from NLP research papers, containing questions that require deep understanding of technical content. All answerable questions are annotated with multiple reference answers to account for different valid expressions of the same information. The ground-truth evidence spans are carefully annotated by domain experts to ensure the reliability of retrieval evaluation.

Dataset	Used Set	Question Num	Doc. Num	Avg. words	Max. words
QASPER	test	1456	416	4170	21,165
QuALITY	dev	2086	115	5022	5967
NarrativeQA	test	10,577	355	51,372	346,902
MultiFieldQA-zh	test	200	200	6701	14918

Table 8: Detailed information of the datasets used in our experiments.

The QuALITY dataset consists of long passages primarily drawn from fiction stories and magazine articles, with an average length longer than typical QA datasets. Each question is accompanied by multiple choice options that test comprehensive understanding of the passage. We utilize the validation set with publicly available labels for our experiments to enable thorough analysis and comparison. Notably, to enable more extensive experimentation, we evaluated on the validation set with publicly available labels rather than the submission-required test set.

The NarrativeQA dataset is a significant advancement in evaluating machine reading comprehension. Unlike traditional datasets that focus on short

passages, NarrativeQA emphasizes understanding long-form narratives, such as entire books or full-length movie scripts. One of the defining features of NarrativeQA is the length of its documents. On average, each document contains 51,372 words. Notably, some individual documents extend up to 346,902 words, making it one of the longest reading comprehension datasets available.

The MultiFieldQA-zh Chinese dataset comprises long-form documents from various sources, including legal texts, government reports, encyclopedias, and academic papers, each accompanied by relevant questions and answers. Compared to the English documents, the Chinese data exhibits unique linguistic structures, expressions, and domain-specific terminology, posing higher demands on models’ capabilities for long-text comprehension and cross-domain information integration. In our experiments, we use the test set for evaluation and comparison. More detailed information is presented in Table 8.

B.2 RST-DT Data Translation

To train a cross-lingual discourse parser, we translate the RST-DT training data into Chinese to augment the available Chinese data. Given the syntactic and word-order differences across languages, we perform translation at the sentence level, which is consistent with our training objective of sentence-level RST parsing. Sentence-level translation preserves inter-sentence discourse relations while deliberately ignoring finer-grained intra-sentence structures, ensuring that the translated corpus aligns with the granularity at which our parser operates. We then combine the original RST-DT data, the sentence-level RST data, and the Chinese-translated data to form the final training corpus.

Training Data	Test Accuracy / %
English only	50.63
English & Chinese	50.51

Table 9: Parser performance on the English RST-DT test set under two training configurations. Adding Chinese translated data does not degrade English parsing accuracy.

For the translation process, we employ the GPT-4o model. GPT-4o ranked first among all translation systems in the TransBench benchmark (Li et al., 2025), surpassing both traditional machine translation systems and all large language models

released before 2025, providing strong evidence for its translation capability. To further validate translation reliability, we randomly sampled 300 translated sentences for human evaluation, finding that over 95% of translations were correct. The rare deviations observed were limited to very short sentences or proper nouns lacking sufficient context, and such cases have negligible impact on parser training. To confirm that the augmented Chinese data does not degrade English parsing ability, we compare parser performance on the original English RST-DT test set under two training configurations, as shown in Table 9. The results show that adding Chinese translated data does not degrade parser performance, further validating the quality of the bilingual training corpus.

The prompt is provided below:

Prompt for Translation

System: You are a helpful assistant.

User: Please translate the following sentences into Chinese. Translate each sentence individually while preserving the original order, and return the results in JSON format.
For example:
Original sentences:
1. Sentence 1
2. Sentence 2
3. Sentence 3
Translation:
{ 1: translated sentence, 2: translated sentence, ... }

Strictly return the output in valid JSON format. Do not include any additional text, and ensure that the output can be directly parsed as JSON.

Original sentences: {sentences}

LLM: JSON-formatted Chinese translations.

B.3 Model Specifications

For semantic embeddings, the Sentence-BERT model (multi-qa-mpnet-base-cos-v1) is based on the MPNet architecture with 768-dimensional representations, specifically optimized for question-answering tasks. The OpenAI embedding model (text-embedding-3-large) represents their latest advancement in semantic representation capabilities.

B.4 Implementation Details

In the discourse analysis process, we preserve sentence boundaries throughout all splitting operations to maintain semantic coherence. For the RAPTOR baseline implementation, we follow the original paper’s collapsed tree approach where all nodes are considered simultaneously during retrieval. The tree construction process in our Bisection baseline ensures a nearly balanced binary structure through recursive division of the sentence set.

B.5 Experiments Compute Resources

For the training of the sentence-level Discourse Parser, we used 1 NVIDIA A100-40G GPU. All other experiments including document discourse parsing, LLM summarization and node embedding, as well as the retrieval and generation processes, are conducted using 4 NVIDIA A800-80G GPUs.

B.6 Prompts

We provides the specific prompts used in our approach to ensure reproducibility.

In the prompt templates: (1) fixed prompts are displayed in black. (2) Input text is highlighted in deep red. (3) The retrieved context is colored in purple. (4) The output generated by the LLM is presented in green.

Prompt for Intermediate Node Text Summarization

System: You are a helpful assistant.

User: Write a summary of the given sentences, keeps as more key information as possible. Only give the summary without other text. Make sure that the summary no more than 200 words.
Given text: {left child node text} {input child node text}

LLM: Summerization result.

Prompt for Question Answering on QASPER Dataset

System: You are a helpful assistant.

User: Using the following information: {context}. Answer the following question in less than 5-7 words, if possible. For yes or no question, only return 'yes' or 'no'.
question: {question}

LLM: Question Answering result.

Prompt for Question Answering on QuALITY Dataset

System: You are a helpful assistant.

User: Given context: {context}.
Answer the following multiple-choice question: {question}

LLM: The correct answer is (A). The context provided...

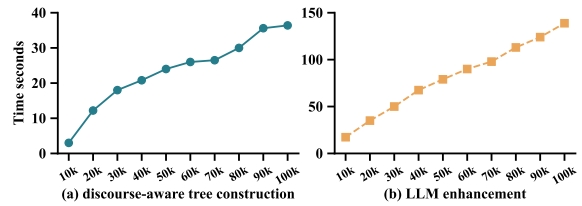


Figure 6: Average processing time of discourse-aware tree construction and LLM enhancement across varying document lengths.

behavior indicates that our method maintains consistent per-word processing efficiency regardless of document length. Moreover, as discussed in Section 4.3, the size and performance of the summarization model have no significant impact on the results, which further suggests that our method can be deployed with a modest computational cost.

Crucially, both tree construction and LLM enhancement are query-independent preprocessing operations performed once per document, allowing costs to be amortized across multiple queries on the same document. For instance, processing a 50,000-word document requires approximately 90 seconds of preprocessing but can subsequently serve unlimited queries with minimal additional overhead. In practical deployment scenarios where documents receive multiple queries over time, this one-time preprocessing cost becomes negligible compared to the cumulative benefits of improved retrieval accuracy.

Compared to traditional chunking approaches that require no preprocessing, our method introduces modest upfront costs but delivers significantly better retrieval quality. The preprocessing overhead is comparable to other hierarchical methods like RAPTOR, while providing superior performance through linguistically-grounded discourse modeling.

Method	10K words	30K words	50K words	70K words	90K words
RAPTOR	63.4s	194.3s	338.5s	488.0s	611.5s
DISRetrieval	20.3s	68.1s	103.0s	124.4s	159.9s
Speedup	3.1×	2.9×	3.3×	3.9×	3.8×

Table 10: Computational efficiency comparison between DISRetrieval and RAPTOR using NVIDIA A800 GPU and LLaMA3.1-8B-Instruct.

To provide a comprehensive evaluation of computational efficiency, we conducted a direct comparison between our method and RAPTOR, another preprocessing-based hierarchical retrieval approach. Both methods were evaluated under identi-

C Extended Experiment Analysis

C.1 Computational Efficiency and Scalability

To evaluate the computational efficiency and scalability of our approach, we conducted comprehensive timing analysis across varying document lengths. We systematically selected documents ranging from 10,000 to 100,000 words, performed discourse-aware tree construction and LLM enhancement, and measured the computation time for each process. All experiments were conducted on a single A800 80G GPU, with LLM enhancement performed using the LLaMA3.1-8B-Instruct model accelerated by the vLLM framework. To ensure statistical reliability, we sampled 10 documents for each length range and report the average processing time with standard deviation. The detailed results are presented in Figure 6.

The processing time for both discourse-aware tree construction and LLM enhancement exhibits near-linear scaling with document length, demonstrating favorable computational complexity for long documents. Specifically, tree construction for documents up to 100,000 words requires less than 40 seconds, while LLM enhancement completes within 140 seconds, confirming practical feasibility for real-world applications. The linear scaling

cal hardware and software configurations to ensure fair comparison. Specifically, we used a single NVIDIA A800 GPU with LLaMA3.1-8B-Instruct as the language model for both discourse tree construction in our method and recursive summarization in RAPTOR.

We tested both methods on documents of varying lengths ranging from 10,000 to 90,000 words to assess scalability. For our method, the total time includes both discourse parsing and LLM-based node enhancement. For RAPTOR, the time includes the complete recursive summarization process that builds the hierarchical tree structure. As shown in Table 10, our method demonstrates consistent computational advantages across all document lengths, achieving approximately $3\times$ speedup compared to RAPTOR. For instance, processing a 50,000-word document requires 103.0 seconds for our method versus 338.5 seconds for RAPTOR, representing a $3.3\times$ speedup. This efficiency advantage becomes even more pronounced for longer documents, with speedup factors reaching $3.9\times$ for 70,000-word documents.

The computational efficiency of our method stems from two key factors. First, discourse parsing operates at the sentence level rather than requiring recursive processing of all text segments, which significantly reduces the number of LLM calls needed. Second, our LLM enhancement process is selective and targeted, focusing only on internal nodes that require summarization, whereas RAPTOR performs recursive summarization across all levels of the hierarchy. Despite this substantial speedup, our method achieves superior retrieval performance as demonstrated in Tables 1 2 4, indicating that computational efficiency does not come at the cost of effectiveness.

It is important to note that both discourse parsing and LLM enhancement in our framework are query-independent preprocessing steps performed once per document. In real-world deployment scenarios where multiple queries are issued against the same document collection, these preprocessing costs are amortized across all queries, making the per-query computational overhead negligible. This characteristic makes our approach particularly suitable for applications such as scientific literature review, legal document analysis, and enterprise knowledge management, where document collections are relatively stable but query patterns are diverse and frequent.

C.2 Analysis of Multi-granularity Adaptive Retrieval

We conducted a thorough statistical analysis to evaluate the multi-granularity adaptive retrieval capability of our method. For each query, we retrieved the Top-20 nodes from the full discourse-aware tree, and analyzed the characteristics of retrieved intermediate nodes across datasets.

Retrieved Node Characteristics Across Datasets.

Table 7 reveals significant differences in the retrieved intermediate nodes between datasets. QuALITY exhibits substantially higher values for intermediate node metrics compared to QASPER: the average depth of retrieved intermediate nodes is 2.6 times greater with SBERT embeddings (13.78 vs. 5.29) and 3.1 times greater with OpenAI embeddings (17.09 vs. 5.53). Similarly, the average number of leaf nodes that each retrieved intermediate node maps to is 5.4 times higher with SBERT (67.67 vs. 12.64) and 6.6 times higher with OpenAI (88.67 vs. 13.37). The percentage of intermediate nodes among Top-20 retrieved results is also higher for QuALITY across both embedding methods (73.93% vs. 54.15% with SBERT; 88.53% vs. 65.13% with OpenAI). These differences exist despite QASPER having a 50% longer average sentence length (22.32 vs. 14.86 words).

Distributional Analysis. Figure 7 further illustrates these distinctions through distributional analysis. The sentence length distribution (Fig. 7a) shows that QuALITY is heavily skewed toward shorter sentences (5-15 words), while QASPER has a broader, more even distribution extending to 40+ words. The node depth distributions (Fig. 7b,c) demonstrate that retrieved intermediate nodes from QuALITY frequently reach depths of 10-15, whereas those from QASPER are predominantly shallower (depths below 10), consistent across both embedding methods.

Retrieval Strategy Adaptation. These patterns reveal how our retrieval approach adapts to different document structures. For narrative fiction in QuALITY with abundant dialogue and shorter sentences, the retrieval system favors higher-level intermediate nodes that aggregate related content, as evidenced by the higher percentage of intermediate nodes and greater node depths. This suggests that hierarchical composition is particularly beneficial for documents where meaning is distributed across multiple short sentences. In contrast, for

research papers in QASPER with longer, more self-contained sentences, the retrieval system relies less on deep hierarchical structures, as shown by the shallower node depths and lower percentage of intermediate nodes. This indicates that relevant information in scientific documents can often be retrieved effectively with less compositional processing.

These findings empirically validate that our approach effectively adapts to diverse document structures and information needs. Our discourse-aware tree construction naturally induces appropriate segmentations based on document characteristics rather than arbitrary chunking, enabling multi-level text composition that overcomes the limitations of fixed-granularity methods. Through these mechanisms, DISRetrieval provides a multi-granularity retrieval framework that adaptively selects appropriate granularity levels across different document types, accommodating their inherent structural variations.

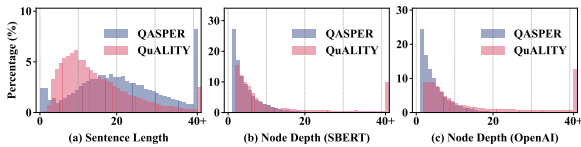


Figure 7: Comparative analysis of distribution difference of two datasets. Figure (a) shows the difference on Sentence Length, while (b) and (c) demonstrate the distribution of retrieved intermediate node depth across different embedding models.

C.3 Ablation study on different Top-K settings

We conduct comprehensive experiments to investigate the impact of varying the value of K from 1 to 20, with results presented in Figure 8.

Our analysis reveals several noteworthy patterns:

Performance Trends. Across all configurations, performance initially improves as K increases, typically peaking around $K = 5$ (marked by the vertical dashed line), after which it either plateaus, gradually declines, or exhibits minor fluctuations. This consistent pattern suggests an optimal balance point where sufficient context is provided without introducing excessive noise.

Dataset-Specific Behaviors. QuALITY demonstrates more pronounced performance variations with changing K values compared to QASPER, with performance differences of up to 2 percentage

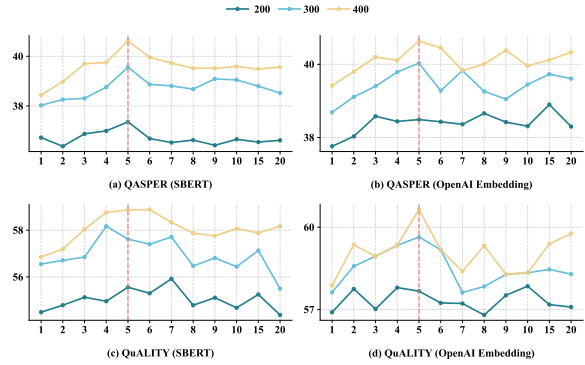


Figure 8: Ablation results with different values of K . The horizontal axis represents different choices of K , and the vertical axis indicates generation performance (F1-match for QASPER and accuracy for QuALITY). All question answering tasks are conducted on the UnifiedQA-3B model.

points between optimal and suboptimal settings. This higher sensitivity likely reflects QuALITY’s more complex narrative structure, where precise evidence selection is particularly crucial.

Context Length Impact. Longer context windows (300 and 400 words) consistently outperform shorter ones (200 words) across all K values. Notably, the performance advantage of longer contexts is most significant when K is small. This suggests that when the system retrieves fewer evidence segments, the comprehensiveness of each individual segment becomes critical, as the model must extract all necessary information from a limited number of passages.

Based on these observations, we adopt $K = 5$ as our optimal setting for all main experiments, as it consistently delivers strong performance across datasets and embedding methods while maintaining computational efficiency.

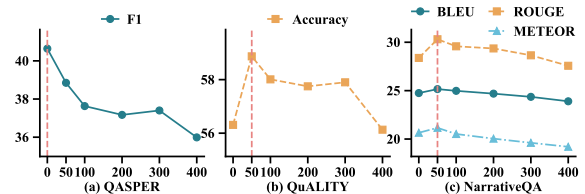


Figure 9: Ablation results with different values of τ across three datasets. Sentence-BERT is used as the embedding model and UnifiedQA-3B as the question-answering model, with a context length of 400 words.

C.4 Ablation study on parameter τ across different datasets

Equation 1 employs a threshold parameter τ to decide whether use an LLM for summarization or directly merge the subtree. To assess the impact of this parameter and identify the optimal settings, we conduct experiments on three datasets. The results are shown in figure 9.

Performance Trends. For QASPER, the F1 score is highest when $\tau = 0$ and gradually declines as τ increases, indicating that summarization is more beneficial than direct merging sentences of subtrees. In contrast, QuALITY and NarrativeQA reach their peak performance at $\tau = 50$ and gradually decline afterward, only approaching the level observed at $\tau = 0$ when τ reaches 300 and 200 respectively, suggesting that a moderate merging strategy is essential for these datasets.

Dataset Analysis These results are consistent with the characteristics of the datasets. QASPER consists of academic papers, where sentences are relatively independent and logically complete. So keeping sentences as the minimal unit allows for more precise information extraction. In contrast, QuALITY and NarrativeQA mainly contain novels and movie scripts, which include numerous dialogues and short sentences that require a relatively complete context to convey the correct semantics. For these two datasets, merging subtrees based on rhetorical structure parsing results enables logically related sentences to form more coherent text segments, preventing semantic fragmentation. At the same time, moderate merging is crucial, as excessively long text segments may introduce redundancy. Therefore, we set $\tau = 0$ for QASPER, $\tau = 50$ for QuALITY and NarrativeQA respectively. This can also provide guidance for other datasets: for documents with more formal and complete sentences, smaller τ values are preferable, whereas for documents with shorter and more fragmented sentences, relatively larger τ values should be used.

C.5 Illustrative Example of Discourse-Guided Evidence Retrieval

To illustrate how discourse structure guides evidence selection in an intuitive manner, we present a concrete example in Figure 10. Sentences 4 and 5 jointly answer the query and are organized under the same internal node *List 2* in the discourse tree.

Query.

What NLP technologies does the system use?

Document Discourse Tree.

```
ROOT
+-- (Topic)
+-- (List 1)
|   +-- 1. In this report we present a system that can
|       |   generate political speeches for a desired
|       |   political party.
|   +-- 2. Furthermore, the system allows to specify
|       |   whether a speech should hold a supportive
|       |   or opposing opinion.
+-- (Elaboration)
+-- 3. The system relies on a combination of several
|   NLP methods which are discussed in this report.
+-- (List 2)
+-- 4. These include n-grams, Justeson&Katz POS
|   tag filter, recurrent neural networks, and
|   latent Dirichlet allocation.
+-- 5. Sequences of words are generated based on
|   probabilities obtained from two underlying
|   models: a language model takes care of
|   grammatical correctness while a topic model
|   aims for textual consistency.
```

Figure 10: An illustrative example of discourse-guided evidence retrieval. Given the query *What NLP technologies does the system use?*, sentences 4 and 5 are grouped under the same discourse node *List 2* and retrieved together as complementary evidence, whereas flat chunking and semantic clustering methods may fail to preserve this grouping.

The hierarchical retrieval mechanism identifies this intermediate node and retrieves both sentences together as complementary evidence. Flat chunking methods may split sentences 4 and 5 into separate chunks depending on the chunk boundary, while semantic clustering methods may separate them due to surface-level lexical differences despite their shared discourse function. This example demonstrates that discourse structure captures the grouping of semantically related sentences in a linguistically principled manner, enabling more complete and coherent evidence retrieval than methods based solely on fixed boundaries or surface similarity.

Algorithm 2: DISRetrieval: Discourse Structure-based Long Document Retrieval

```
Input: Document  $D$  containing sentences  $S_i$ , paragraphs  $P_i$ ; Query  $q$ 
Output: Retrieved evidence segments  $E$ 
/* Stage 1: Discourse-Aware Tree Construction */
for each paragraph  $P_i \in D$  do
  Initialize stack  $\sigma$  and sentence queue  $\beta \leftarrow S_i$ ;
  while  $\beta$  not empty OR  $|\sigma| > 1$  do
    /* RST parsing operations */
     $a_t \leftarrow$  Select action from {shift, reduce, pop_root}; // Determine next parsing action
    switch  $a_t$  do
      case shift do
        Move next sentence from  $\beta$  to  $\sigma$ ; // Add new sentence to stack
      end
      case reduce do
         $s_2, s_1 \leftarrow$  Pop top two elements from  $\sigma$ ;
        Determine discourse relation  $r$  between  $s_1$  and  $s_2$ ; // Using RST relations
        Create new node with relation  $r$ ;
        Push new node to  $\sigma$ ;
      end
      case pop_root do
        Final tree node  $\leftarrow$  Pop from  $\sigma$ ; // Complete the tree
      end
    end
  end
   $T_i \leftarrow$  Resulting discourse tree;
  /* LLM Enhancement for Paragraph Tree */
  for each non-leaf node  $v \in T_i$  (bottom-up order) do
     $v_l, v_r \leftarrow$  Get left and right children;
    if  $|v_l| + |v_r| \geq \tau$  then
       $v^* \leftarrow f_{LLM}(v_l, v_r)$ ; // Generate concise summary using LLM
    else
       $v^* \leftarrow f_{merge}(v)$ ; // Merge the subtree rooted at v
    end
  end
  Store root text representation  $t_{root}^i$  for  $T_i$ ;
end
Initialize stack  $\sigma$  and queue  $\beta \leftarrow \{t_{root}^1, \dots, t_{root}^n\}$ ;
while  $\beta$  not empty OR  $|\sigma| > 1$  do
  /* Similar RST parsing at document level */
  Apply RST parsing operations as in Phase 1;
  Build document tree  $T_{doc}$ ;
end
/* Stage 2: Node Representation */
for each non-leaf node  $v \in T_{doc}$  (bottom-up order) do
   $t_l, t_r \leftarrow$  Get content from left and right children;
  if  $|t_l| + |t_r| \geq \tau$  then
     $t_v \leftarrow f_{LLM}(t_l, t_r)$ ; // LLM-based enhancement for document-level nodes
  else
     $t_v \leftarrow t_l \oplus t_r$ ; // Direct concatenation for short text, don't need to merge subtree
  end
end
 $T_D \leftarrow$  Replace leaf nodes in  $T_{doc}$  with corresponding  $T_i$ ; // Integrate trees into unified structure
for each node  $v \in T_D$  do
   $e_v \leftarrow$  Encoder( $t_v$ ); // Generate dense vector representations
end
/* Stage 3: Hierarchical Evidence Retrieval and Selection */
 $e_q \leftarrow$  Encoder( $q$ ); // Transform query to embedding space
scores  $\leftarrow \{\}$ ,  $E \leftarrow \{\}$ ;
for node  $v \in T_D$  do
  scores[ $v$ ]  $\leftarrow$  cosine( $e_q, e_v$ ); // Compute relevance scores
end
Apply Algorithm 1 for evidence selection;
return  $E$ 
```
