

Dynamic, Coherent, and Lossless: The Iterative Structured Content Filling Framework for Infinite Document Parsing

Anonymous ACL submission

Abstract

Parsing long, unstructured documents into hierarchical machine-readable formats is critically challenged by the tension between document length and limited LLM context. Prevailing one-shot or stateless chunk-based methods suffer from the “lost-in-the-middle” phenomenon and fragmented structures. We argue for reframing the task as an iterative, stateful editing process. This paper introduces the **Iterative Structured Content Filling (ISCF)** framework, which incrementally builds a global tree through LLM-driven steps under partial observability. ISCF integrates three pillars: (1) **Dynamic Boundary Re-alignment (DBR)** for unsupervised chunk optimization via semantic coherence and uncertainty; (2) **Adaptive Structural Masking & Hierarchical-Aware Path Embedding (HAPE)** for efficient global awareness within finite context; (3) **Bijective Path-Content Mapping** ensuring lossless reconstruction. Extensive experiments on financial and legal documents demonstrate ISCF’s superior accuracy, robustness, and efficiency over strong baselines, offering a principled solution for unbounded document structuring.

1 Introduction

The digital transformation of industries such as finance, law, and academia has generated vast repositories of long-form unstructured documents. Automating the conversion of these documents into machine-readable hierarchical formats is crucial for downstream applications like intelligent review, knowledge graph construction, and semantic analysis. However, existing methods for document structure parsing face a fundamental tension between the unbounded length of real-world documents and the fixed context windows of Large Language Models (LLMs), as well as between localized text processing and the need for global structural consistency (Xu et al., 2020).

Prevailing approaches exhibit significant limitations. Full-text ingestion is immediately constrained by context ceilings and suffers from the well-documented “lost-in-the-middle” phenomenon, where model performance degrades sharply on information located in the middle of long inputs (Liu et al., 2023). Independent chunk parsing treats each segment in isolation, leading to catastrophic structural fractures at boundaries due to a lack of cross-chunk state. Overlapping sliding windows introduce redundant computation and often fail to resolve ambiguous content attribution across windows (Ding et al., 2024). Underlying these methods is a treatment of the task as a one-shot or stateless sequence-to-tree conversion (Sutskever et al., 2014; Dong and Lapata, 2016), which fundamentally ignores the inherent continuity and hierarchical dependencies within long documents.

We argue that parsing long documents should be modeled not as a single transduction, but as an iterative, stateful editing process. Analogous to how a human editor builds an outline, the system should maintain and refine a growing structural skeleton, integrating new content against it and making dynamic adjustments based on evolving context. This paradigm directly addresses the challenge of partial observability under limited context windows and ensures each local decision is grounded in a distilled, global structural view.

In this work, we formalize the task as **Structured Content Filling** and introduce the **Iterative Structured Content Filling (ISCF)** framework. ISCF iteratively refines an initially empty hierarchical tree into a complete structure through a series of LLM-driven editing steps. Each step is conditioned on a local text segment and a compressed, adaptive view of the current global structure. The framework is built upon three synergistic technical pillars:

1. **Dynamic Boundary Re-alignment (DBR):**

An online, unsupervised mechanism that optimizes chunk boundaries based on semantic coherence and local parsing uncertainty, preventing logical units from being severed by fixed segmentation (Duarte et al., 2024).

2. Adaptive Structural Masking & Hierarchical-Aware Path Embedding (HAPE):

A context compression and prompting strategy that retains the complete hierarchical skeleton and focuses attention on a relevant anchor point, effectively reconciling unbounded document structure with finite LLM context.

3. Bijective Path-Content Mapping:

A deterministic protocol that guarantees lossless, 100% accurate reconstruction of the original text from the final structured tree, eliminating content omission or duplication.

Our contributions are summarized as follows:

1. We propose ISCF, a novel, model-agnostic iterative paradigm for long-document structure parsing that ensures global coherence through stateful, incremental tree construction and editing.
2. We design DBR, an unsupervised chunking algorithm that dynamically realigns boundaries based on semantic cohesion and a prompting-based self-consistency measure, significantly improving local structural continuity.
3. We develop adaptive tree masking and HAPE prompting to efficiently maintain structural awareness within LLM context limits, mitigating information decay in long contexts.
4. Through comprehensive experiments on challenging financial and legal document datasets, we demonstrate ISCF’s superior accuracy, robustness, and efficiency over strong baselines, providing a principled solution for transforming unbounded unstructured text into precise hierarchical representations.

2 Method

We formalize the structural parsing of long, unstructured documents as a **Structured Content Filling** task. Given a long document $D = [s_0, s_1, \dots, s_{M-1}]$ consisting of a sequence of sentences where each s_i is a maximal textual unit ending with terminal punctuation and no explicit structural markers (e.g., headings or section delimiters)

are present our goal is to iteratively refine an initially empty hierarchical tree $T^{(0)}$, containing only a root node, into a complete, machine-readable structure tree $T^{(K)} = (\mathcal{N}, \mathcal{E})$, where \mathcal{N} denotes the set of nodes (including both structural and content nodes) and $\mathcal{E} \subseteq \mathcal{N} \times \mathcal{N}$ represents the set of directed parent-child edges that define the hierarchical structure (Wang et al., 2025), and K denotes the total number of LLM inference steps.

Unlike conventional one-shot approaches (e.g., sequence-to-tree or sequence-to-sequence models) (Sutskever et al., 2014; Dong and Lapata, 2016), our proposed **Iterative Structured Content Filling** (ISCF) framework formulates this as an iterative, stateful editing process. Constrained by the limited context window of large language models (LLMs), this process operates under partial observability: at each step, the system observes only a local text segment and a compressed view of the current structure, then executes discrete edit actions (e.g., creating nodes, assigning content) to update its internal state. This paradigm directly addresses the well-documented lost-in-the-middle phenomenon in long-context LLMs (Liu et al., 2023) where performance degrades on middle segments of long inputs by ensuring that each decision is grounded in both recent content and a distilled global outline.

The proposed solution is illustrated in Figure 1.

2.1 Problem Formulation

The input document D is a sequence of sentences from an unstructured source lacking predefined headings or sections. The target structure tree T comprises two types of nodes, both inferred purely from semantic content:

- **Structural nodes** $n_s \in \mathcal{N}_s$: represent inferred hierarchical sections (e.g., chapters, subsections). Each is formalized as a tuple (title, ℓ) , where the title is generated by the model based on semantic content, and $\ell \in \{1, 2, \dots, L\}$ denotes the inferred depth level. The span covered by n_s is implicitly defined as the union of spans of all its descendant content nodes: $\text{span}(n_s) := \bigcup_{n_c \in \text{Desc}(n_s) \cap \mathcal{N}_c} \text{span}(n_c)$.
- **Content nodes** $n_c \in \mathcal{N}_c$: represent contiguous sentence spans belonging to a specific structural node, formalized as (π, span) . Here, π is a unique path identifier from

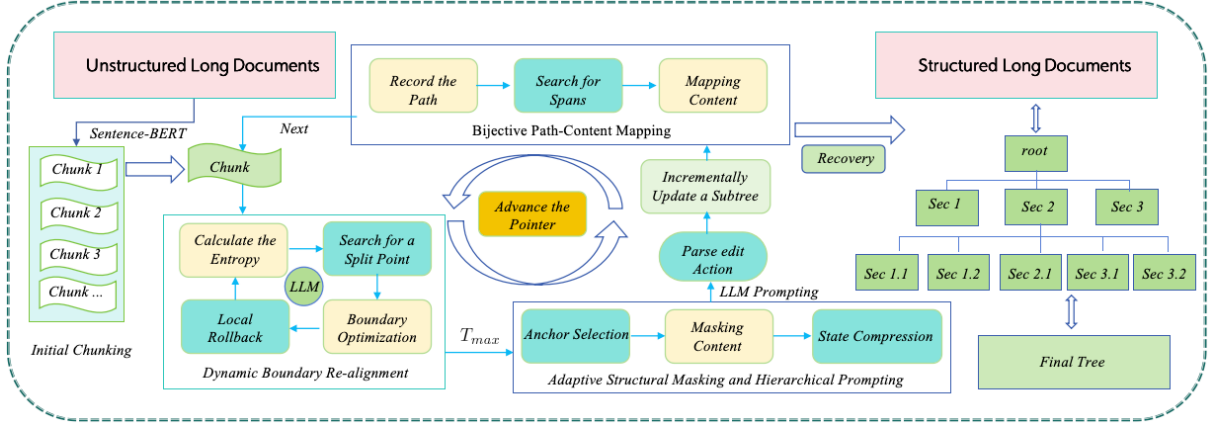


Figure 1: Architecture of the Iterative Structured Content Filling Framework.

the root to the parent structural node (e.g., root.chapter1.section2), and span = (i, j) denotes the inclusive sentence indices in D corresponding to the original text. This ensures deterministic recovery and avoids ambiguity from repeated phrases.

Constraints: The final tree $T^{(K)}$ must satisfy:

1. **Content Coverage:** $\bigcup_{n_c \in \mathcal{N}_c} D[\text{span}(n_c)] = D$;
2. **Content Disjointness:** for any $n_c \neq n'_c$, $\text{span}(n_c) \cap \text{span}(n'_c) = \emptyset$;
3. **Structural Containment:** if structural node n_i is an ancestor of n_j , then $\text{span}(n_i) \supseteq \text{span}(n_j)$.

At global step t , the system maintains a state $S^{(t)} = (T^{(t)}, \Phi^{(t)}, \text{ptr}^{(t)})$, where:

- $T^{(t)}$: the partially constructed structure tree with generated structural nodes and associated content assignments.
- $\Phi^{(t)}$: a bijective mapping from finalized path identifiers π to their exact original sentence spans in D , used for lossless reconstruction.
- $\text{ptr}^{(t)}$: the index of the next unprocessed sentence in D (an absolute position in $[0, M]$).

The system's observation $O^{(t)}$ the prompt fed to the LLM consists of:

1. The **current text chunk** $C^{(t)}$: a contiguous segment starting at $\text{ptr}^{(t)}$, obtained via dynamic chunking (Section 2.2).

2. A **compressed structural context**: a masked version $T_M^{(t)}$ of $T^{(t)}$ via adaptive masking (Section 2.3).

Based on $O^{(t)}$, the LLM generates a sequence of edit actions $A^{(t)} = \{a_1, a_2, \dots\}$, which are parsed into an edit script $\Delta T^{(t)}$ and applied to $T^{(t)}$ to yield $T^{(t+1)}$. The process terminates when $\text{ptr}^{(t)} = M$, producing a final tree $T^{(K)}$ that satisfies all constraints.

2.2 Dynamic Boundary Re-alignment (DBR)

Fixed-length chunking often severs semantically coherent units a critical failure mode for unstructured documents (Duarte et al., 2024). We propose **Dynamic Boundary Re-alignment (DBR)**, an online, unsupervised mechanism that adjusts chunk boundaries based on local parsing uncertainty.

2.2.1 Initial Chunking and Uncertainty Estimation

We first partition D into initial chunks represented as index intervals $\{(b_0, e_0), (b_1, e_1), \dots\}$ using a greedy merging algorithm based on Sentence-BERT embeddings (Reimers and Gurevych, 2019), respecting the token limit L_{\max} . Specifically, we iteratively merge consecutive sentences with the highest embedding similarity until adding the next sentence would exceed L_{\max} tokens.

Let c denote the current chunk index being processed. At retry attempt r for chunk c , let $\Delta \mathcal{N}_c^{(r)}$ denote the set of newly generated content nodes. To estimate structural uncertainty without requiring model fine-tuning, we use a prompting-based self-consistency strategy: we run the LLM R times with different random seeds and extract the empirical distribution over semantic roles (e.g.,

topic, evidence, conclusion) assigned to each content node. We define:

$$\mathcal{H}_c^{(r)} = \frac{1}{|\Delta\mathcal{N}_c^{(r)}|} \sum_{n_c \in \Delta\mathcal{N}_c^{(r)}} H(p_{\text{emp}}(\text{role} \mid n_c)), \quad (1)$$

where $H(\cdot)$ is entropy and p_{emp} is the empirical role distribution. High $\mathcal{H}_c^{(r)}$ suggests the current boundary cuts through a coherent semantic unit.

2.2.2 Boundary Optimization via Semantic Coherence

When $\mathcal{H}_c^{(r)} > \tau$ and $c \geq 1$, DBR triggers a boundary adjustment between chunks $c - 1$ and c . We search for a split point k (measured as the number of sentences from the start of chunk c) that maximizes intra-chunk coherence while encouraging smoothness across the boundary:

$$k^* = \arg \max_{1 \leq k \leq \min(w, |C_c| - 1)} \left[\underbrace{\text{Coh}(C_c^{<k}) + \text{Coh}(C_c^{\geq k})}_{\text{internal cohesion}} - \gamma \cdot \underbrace{\text{Dist}(C_{c-1}, C_c^{<k})}_{\text{boundary discontinuity}} \right], \quad (2)$$

where $w \in \mathbb{N}^+$ is the maximum number of initial sentences in chunk c eligible for reassignment, $\gamma > 0$ balances the trade-off, and $|C_c| = e_c - b_c$ denotes the number of sentences in chunk c . The coherence and distance functions are defined as:

$$\text{Coh}(X) = \frac{1}{|X|^2} \sum_{s, s' \in X} \cos(\mathbf{e}_s, \mathbf{e}_{s'}), \quad |X| \geq 1, \quad (3)$$

$$\text{Dist}(X, Y) = 1 - \cos(\bar{\mathbf{e}}_X, \bar{\mathbf{e}}_Y), \quad (4)$$

where $\bar{\mathbf{e}}_X$ is the mean embedding of all sentences in X . Here, $C_{c-1} = D[b_{c-1} : e_{c-1}]$ and $C_c = D[b_c : e_c]$ denote the current sentence sequences of chunks $c - 1$ and c , respectively.

We then update the chunk boundaries by moving the first k^* sentences of C_c to the end of C_{c-1} :

$$e_{c-1} \leftarrow b_c + k^*, \quad (5)$$

$$b_c \leftarrow b_c + k^*. \quad (6)$$

A **local rollback** is performed: T and Φ are reverted to the checkpoint saved after processing chunk $c - 2$, and chunks $c - 1$ and c are reparsed jointly. Checkpoints for the last two processed chunks are maintained, ensuring $O(1)$ rollback cost. The retry loop runs at most T_{max} times per chunk.

2.3 Adaptive Structural Masking and Hierarchical Prompting

To reconcile unbounded document structure with finite LLM context, we introduce adaptive compression and structured prompting.

2.3.1 Anchor Selection and Adaptive Masking

Given state $S^{(t)}$, we identify the **anchor structural node** $n_{\text{anchor}}^{(t)}$ as the deepest structural node whose span ends closest to but does not exceed $\text{ptr}^{(t)}$. If no such node exists (e.g., at initialization), we use the root node.

The masking function \mathcal{M} produces a compressed tree $T_M^{(t)}$ by:

- Retaining the full **skeleton**: all structural nodes (with titles and levels) and parent-child relations.
- For content nodes: retain full text only for those within the anchors subtree and overlapping with recent context; replace others with [CONTENT_MASK], preserving path π and role (if available).
- Collapsing distant subtrees (depth $> d_{\text{max}}$ and path distance $> \delta$ from anchor) into a single [SUBTREE_MASK] token.

This achieves $>60\%$ context reduction while preserving navigational cues.

2.3.2 Hierarchical-Aware Path Embedding (HAPE)

To ground the LLM in global structure, we design a hierarchical prompting strategy integrating:

1. **Path descriptors**: natural-language paths (e.g., under Section 3.1: Data Collection in Chapter 3: Methodology).
2. **Relative positioning**: explicit sequential context (e.g., this follows immediately after Experimental Setup).
3. **Structured delimiters**: XML-style tags for robust component separation.

The full prompt template is in Appendix A. This anchors working memory to recent content while preserving global awareness.

2.4 Bijective Path-Content Mapping and Deterministic Recovery

To guarantee 100% content recall:

- **Recording:** When a content node is finalized-defined as having low uncertainty or being outside the active editing window (i.e., not in the last two chunks) its (π, span) is stored in Φ .
- **Recovery:** At output time, replace each [CONTENT_MASK] by querying Φ using π and extracting $D[\text{span}]$.
- **Version control:** During DBR rollbacks, Φ is reverted to match the active structural state, ensuring bijectivity is preserved throughout execution.

2.5 Complete Algorithm

The dominant cost is LLM inference, invoked $O(N \cdot T_{\max})$ times with context length bounded by $L_{\max} + C_{\text{struct}}$ ($C_{\text{struct}} < 512$ tokens empirically). DBRs semantic computation adds negligible overhead. In practice, DBR activates sparsely ($< 15\%$ of chunks in our experiments), yielding near-linear scaling.

2.6 Future Work: Towards Enhanced Robustness

To further improve the robustness of ISCF against noisy inputs and implicit hierarchies, a natural extension is to integrate **structure-aware adversarial training** during the policy fine-tuning phase. This involves generating perturbed documents via transformations that either preserve (e.g., mild sentence reordering) or disrupt (e.g., injecting ambiguous topic shifts) the underlying structure. Training with a combined loss that enforces parsing consistency between original and perturbed documents would enhance the model’s stability against real-world variations, making the iterative process more reliable. This constitutes the immediate next step of our research. Furthermore, exploring how such robustness-enhancing techniques can enable remarkable performance even with smaller-parameter models presents a compelling long-term direction for efficient, deployable document structuring systems.

3 Experiments

We conduct a comprehensive evaluation of the proposed ISCF framework against state-of-the-

art baselines on challenging, real-world long-document datasets¹.

Our experiments are designed to answer the following research questions:

- **RQ1 (Overall Performance):** How does ISCF perform compared to existing methods in terms of holistic structural parsing accuracy and content fidelity on long, unstructured documents ?
- **RQ2 (Component Ablation):** What is the contribution of each core component (DBR, Adaptive Masking, Bijective Mapping) to the overall performance of ISCF ?
- **RQ3 (Robustness to Length):** How does the performance of ISCF scale with increasing document length compared to baseline methods ?
- **RQ4 (Efficiency):** What is the computational overhead of ISCFs iterative process and DBR mechanism, and how does it compare to other methods in terms of inference cost ?

3.1 Datasets

We evaluate on two long-document domains: finance and law, both characterized by implicit, semantically driven hierarchical structure.

- **FinDoc-Struct:** A collection of 500 annual reports from Chinese A-share listed companies (2011-2023), sourced from publicly disclosed scanned PDFs. After OCR processing, documents average $12,500 \pm 2,800$ sentences ($\sim 300\text{K}$ tokens) and contain no reliable structural markup (e.g., section headers or XML/HTML tags). We manually annotated 50 high-quality OCR documents with a 5-level hierarchy (Part \rightarrow Section \rightarrow Sub-section \rightarrow Sub-subsection \rightarrow Paragraph); the remaining 450 are used for qualitative and efficiency analysis.
- **LegalCase-Struct:** 200 Chinese court judgments (2000-2023) from public sources such as China Judgments Online², including decisions from the Supreme Peoples Court and higher peoples courts. Documents average $8,400 \pm 3,100$ sentences and loosely follow a standard legal structure (e.g., “Statement of Claims”, “Facts”, “Court’s Reasoning”,

¹We will publicly release our data and code.

²<https://wenshu.court.gov.cn>

419 “Judgment”), though with significant format- 469
420 ting variation. Expert annotators labeled dis- 470
421 course units (e.g., “Plaintiff Alleges”, “This 471
422 Court Holds”, “Judgment Order”) and their 472
423 hierarchical relationships in 30 representative 473
424 cases. 474

425 Both datasets lack consistent explicit structural 475
426 markers, requiring models to infer document hier- 476
427 archy solely from semantic and contextual cues. 477

428 3.2 Baselines and Implementation Details 478

429 We compare ISCF against three representative cat- 479
430 egories of baseline methods for long-document 480
431 processing. For all LLM-based methods, we use 481
432 DeepSeek-V3.2 as the backbone model due to its 482
433 strong performance and extensive context window 483
434 (128K tokens) (DeepSeek, 2024), ensuring fair 484
435 comparison of algorithmic paradigms rather than 485
436 raw model capabilities(Chalkidis et al., 2022). 486

- 437 • **Vanilla Full-Context (FC)**: The entire doc- 487
438 ument is truncated/padded to fit the model’s 488
439 maximum context window (we use 100K 489
440 tokens) and a single prompt instructs the 490
441 model to output the complete hierarchical 491
442 structure. This serves as an upper-bound for 492
443 what a “one-shot” long-context model can 493
444 achieve, and directly tests the “lost-in-the- 494
445 middle” phenomenon. 495
- 446 • **Fixed Chunking + Naive Merge (FCNM)**: 496
447 The document is split into non-overlapping 497
448 chunks of 800 tokens (approximately 40 498
449 sentences) using Sentence-BERT embed- 499
450 dings for semantic cohesion (Reimers and 500
451 Gurevych, 2019). Each chunk is parsed in- 501
452 dependently into a local tree. A simple, 502
453 rule-based post-processing step attempts to 503
454 merge consecutive chunks if their highest- 504
455 level nodes share high TF-IDF cosine similar- 505
456 ity (> 0.7). This mimics common production 506
457 pipelines. 507
- 458 • **Sliding Window with State Cache (SW- 508
459 SC)**: Inspired by recent streaming ap- 509
460 proaches (Zhang et al., 2024), we use a slid- 510
461 ing window of 800 tokens with a 200-token 511
462 stride. The model is prompted to continue ex- 512
463 panding the structure tree from the previous 513
464 state, which is summarized in a cache limited 514
465 to 300 tokens. This baseline tests incremen- 515
466 tal, stateful processing without our proposed 516
467 DBR and advanced masking. 517
- 468 • **ISCF**: We implement ISCF as described in

469 Section 2. Initial chunking uses Sentence- 470
471 BERT embeddings (Reimers and Gurevych, 472
473 2019) with a target of 40 sentences. DBR 474
475 parameters: uncertainty threshold $\tau = 1.2$ 476
477 (selected via grid search on a validation set), 478
479 trade-off $\gamma = 0.5$, search window $w = 15$, 480
481 max retries $T_{\max} = 2$. The entropy thresh- 482
483 old $\tau = 1.2$ corresponds to moderate uncer- 484
485 tainty in role assignment (approximately uni- 486
487 form distribution over 3-4 roles). For adap- 488
489 tive masking, we retain the anchor node’s 490
491 subtree and collapse subtrees beyond depth 492
493 $d_{\max} = 4$ or path distance $\delta = 5$. 494

495 3.3 Evaluation Metrics 495

496 We employ a multi-faceted evaluation protocol 497
498 covering both structural fidelity and content pres- 499
499 ervation. 500

- 501 • **Tree Edit Distance (TED)** (Zhang and 502
503 Shasha, 1989): Measures the minimum num- 504
505 ber of node insertions, deletions, and re- 506
507 labelings to transform the predicted tree 508
509 into the gold tree. We use normalized 509
510 TED: $N\text{-TED} = \text{TED}/(|\text{Nodes_Gold}| + |\text{Nodes_Pred}|)$. 511
- 512 • **Role Classification F1 (R-F1)**: For content 513
514 nodes, we classify their semantic role based 515
516 on the path context. We report the macro- 517
518 averaged F1 score against gold labels. 519
- 520 • **Content Recall (C-R) & Overlap (C-O)**: C- 521
522 R is the proportion of gold sentences cor- 523
524 rectly assigned to any node. C-O measures 525
526 precision by penalizing boundary duplication: 527
528 $C\text{-O} = 1 - \frac{\#\text{duplicated sentences}}{\#\text{sentences predicted}}$. ISCF guaran- 529
530 tees $C\text{-R}=1.0$ by design. 531
- 532 • **Boundary Coherence Score (BCS)**: For 533
534 each predicted boundary between content 535
536 nodes, we compute the cosine similarity be- 537
538 tween the embeddings of the last and first 539
540 sentences of adjacent nodes (Reimers and 541
542 Gurevych, 2019). Lower average similarity 543
544 indicates better boundary placement (coher- 545
546 ent units are not severed). 547

548 4 Results discussion 548

549 4.1 Overall Performance (RQ1) 549

550 We conducted the experiment multiple times and 551
552 took the average as the final result. Table 1 553
554 presents the main results. 555

Table 1: Overall performance comparison. Best results are in **bold**. ISCF significantly outperforms all baselines.

(a) FinDoc-Struct					(b) LegalCase-Struct		
Method	N-TED ↓	R-F1 ↑	C-R ↑	C-O ↑	Method	N-TED ↓	R-F1 ↑
Vanilla FC	0.412	0.538	0.923	0.881	Vanilla FC	0.385	0.621
FCNM	0.521	0.462	0.894	0.812	FCNM	0.503	0.498
SW-SC	0.368	0.602	0.981	0.942	SW-SC	0.341	0.685
ISCF (Ours)	0.287	0.724	1.000	0.988	ISCF (Ours)	0.258	0.763

Table 2: Ablation study on FinDoc-Struct.

Variant	N-TED ↓	R-F1 ↑	C-O ↑	BCS ↓
ISCF	0.287	0.724	0.988	0.15
– DBR	0.355	0.661	0.962	0.31
– AM [†]	0.332	0.643	0.990	0.18
– BM [†]	0.301	0.698	0.941	0.16
– HAPE	0.318	0.672	0.985	0.17

[†]AM = Adaptive Masking; BM = Bijective Mapping.

Table 3: Efficiency comparison.

Method	Calls ↓	Tok/Call ↓	Lat. (s) ↓
Vanilla FC	1	98,750	42.1
FCNM	15.6	1,200	38.5
SW-SC	48.3	1,050	112.7
ISCF (Ours)	22.4	950	65.8

Calls = Number of LLM calls; Tok/Call = Average number of tokens per LLM call; Lat. (s) = Average latency per document in seconds.

- ISCF achieves the best structural accuracy (lowest N-TED, highest R-F1), demonstrating the effectiveness of iterative, stateful editing with global awareness.
- ISCF achieves perfect Content Recall and near-perfect Content Overlap (0.988), validating the bijective mapping and DBR. In contrast, FCNM suffers from severe boundary duplication, and Vanilla FC loses content, evidencing the “lost-in-the-middle” effect (Liu et al., 2023).
- SW-SC performs second best but is outperformed by ISCF, highlighting the benefits of DBR and adaptive masking over fixed-stride sliding windows.

4.2 Ablation Study (RQ2)

Table 2 shows the contribution of each component on FinDoc-Struct.

- **DBR** is critical for local coherence, causing the largest N-TED degradation and BCS increase when removed.
- **Adaptive Masking** is essential for balancing global structure and local content, as its removal hurts R-F1.
- **Bijective Mapping** ensures precision, with heuristic merge leading to content overlap drop.

- **HAPE** improves structural awareness, as flattened context reduces performance.

4.3 Scalability with Document Length (RQ3)

We analyze performance by binning FinDoc-Struct documents by sentence count. ISCF’s N-TED remains stable across length bins, while SW-SC’s degrades noticeably beyond 10,000 sentences. This demonstrates ISCF’s robustness, stemming from adaptive masking (managing global skeleton size) and DBR (localizing boundary issues).

4.4 Efficiency Analysis (RQ4)

Table 3 reports computational cost on the full FinDoc-Struct dataset.

- ISCF requires fewer calls than SW-SC due to its non-overlapping, stateful iteration.
- ISCF uses the smallest context per call (**950 tokens**) thanks to adaptive masking, reducing per-step cost.
- Overall latency offers a favorable accuracy-efficiency trade-off. DBR retries added minimal overhead (< 10% of total time).

4.5 Conclusion of Experiments

The results comprehensively validate ISCF. It significantly outperforms strong baselines in structural accuracy and content fidelity on very long

569	documents, with each technical component contributing measurably. ISCF offers a practical and principled solution for transforming unbounded unstructured text into precise hierarchical representations.	617
570		618
571		
572		
573		
574	5 Related Work	
575	Our work on iterative structure parsing relates to several key areas in NLP.	
576		
577	5.1 Long-Context Modeling for LLMs.	
578	The predominant challenge of processing lengthy documents stems from the fixed context window of Transformer-based LLMs. To address this, early research focused on sparse attention mechanisms, such as the sliding window pattern in Longformer (Beltagy et al., 2020) or the global+local attention in BigBird (Zaheer et al., 2020), which reduce computational complexity. Recent advancements like StreamingLLM (Xiao et al., 2024) stabilize attention for infinite-length inputs. However, these methods primarily alleviate computational constraints and do not inherently solve the <i>semantic</i> challenge of information integration across distant text segments, a phenomenon characterized as the lost-in-the-middle problem (Liu et al., 2023). Our ISCF framework directly tackles this by adopting an iterative, stateful paradigm, ensuring each parsing decision is informed by a compressed global structural memory, thus bypassing the need for unreliable full-context ingestion.	619
579		620
580		621
581		622
582		623
583		624
584		625
585		626
586		627
587		628
588		629
589		630
590		631
591		632
592		633
593		
594		
595		
596		
597		
598	5.2 Document Structure and Hierarchy Parsing.	
599		
600	Parsing documents into a hierarchical format is a well-studied task. Traditional approaches relied on rule-based systems or sequence labeling models using layout features (Saijad et al., 2022). With the rise of LLMs, the task is often framed as a one-shot sequence-to-structure generation problem, using extensive prompting or fine-tuning (Dong and Lapata, 2016). A common production baseline involves processing fixed, independent chunks and merging results heuristically, which often leads to fractured structures at boundaries. Recent work like LumberChunker (Duarte et al., 2024) improves segmentation by considering semantic cohesion. ISCF fundamentally departs from these static or independent-chunk approaches by modeling parsing as a <i>stateful, iterative editing process</i> . This allows for dynamic correction and integration, ensuring global coherence akin to how an editor builds an outline incrementally.	617
601		618
602		
603		
604		
605		
606		
607		
608		
609		
610		
611		
612		
613		
614		
615		
616		
	5.3 Iterative and Stateful NLP Tasks.	
	The paradigm of iterative refinement is explored in related tasks. In long-form question answering, agents may iteratively gather evidence (Gao et al., 2023). In summarization, methods like those in (Zhang et al., 2024) process text in segments while maintaining a running summary. Our work adapts this iterative philosophy to the novel domain of hierarchical structure construction. Unlike simple sliding-window baselines that pass a limited text cache, ISCF actively maintains and refines an explicit, compressed structural skeleton, integrating novel techniques like uncertainty-driven boundary re-alignment (DBR) to manage state across chunk transitions effectively.	619
		620
		621
		622
		623
		624
		625
		626
		627
		628
		629
		630
		631
		632
		633
	6 Conclusion and Outlook	
		634
	This work addresses the critical challenge of parsing long, unstructured documents into precise hierarchical trees by proposing the Iterative Structured Content Filling (ISCF) framework. We reframe the task from a one-shot transduction to a stateful, iterative editing process, addressing fundamental tensions between document length and LLM context limits, and between local processing and global coherence. ISCF integrates three synergistic innovations: Dynamic Boundary Re-alignment optimizes chunk boundaries for semantic continuity; Adaptive Structural Masking with HAPE maintains global structural awareness within finite context; and Bijjective Path-Content Mapping guarantees lossless text reconstruction. Comprehensive experiments on financial and legal documents demonstrate ISCF’s superior structural accuracy and content fidelity over strong baselines, while maintaining favorable efficiency. This framework provides a principled, effective solution for transforming unbounded unstructured text into machine-readable hierarchical representations. Future work will focus on enhancing robustness to noisy inputs, extending to graph structures, optimizing efficiency, and testing broader domain generalization.	635
		636
		637
		638
		639
		640
		641
		642
		643
		644
		645
		646
		647
		648
		649
		650
		651
		652
		653
		654
		655
		656
		657
		658
		659
		660
	7 Limitations	
		661
	ISCF provides a robust framework through the co-design of segmentation, context management, and state tracking. Limitations include: 1) Depen-	662
		663
		664

665	dence on initial sentence segmentation quality. 2)	Nils Reimers and Iryna Gurevych. 2019. Sentence-	719
666	Assumption of tree structures; extending to graphs	bert: Sentence embeddings using siamese bert-	720
667	(DAGs) is future work. 3) Self-consistency check	networks. In <i>Proceedings of the Conference on</i>	721
668	for DBR adds cost; a lightweight validator could	<i>Empirical Methods in Natural Language Processing</i>	722
669	improve efficiency.	(<i>EMNLP</i>). Standard method for semantic sentence	723
		similarity, used in DBR.	724
670	References	Muhammad Saijad and 1 others. 2022. Parsing docu-	725
671	Iz Beltagy, Matthew E Peters, and Arman Cohan.	ments with guidance of layout and semantics. <i>Pro-</i>	726
672	2020. Longformer: The long-document transformer.	<i>ceedings of the Conference on Empirical Methods in</i>	727
673	<i>arXiv preprint arXiv:2004.05150</i> .	<i>Natural Language Processing</i> .	728
674	Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael	Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014.	729
675	Bommarito, Ion Androutsopoulos, Daniel Martin	Sequence to sequence learning with neural networks.	730
676	Katz, and Nikolaos Aletras. 2022. Lexglue: A	In <i>Advances in Neural Information Processing Sys-</i>	731
677	benchmark dataset for legal language understanding	<i>tems (NeurIPS)</i> . Introduced the encoder-decoder	732
678	in english. In <i>Proceedings of the 60th Annual Meet-</i>	RNN architecture for mapping sequences to se-	733
679	<i>ing of the Association for Computational Linguistics</i>	quences.	734
680	(<i>ACL</i>). Introduces a comprehensive benchmark for	Shu Wang, Yingli Zhou, and Yixiang Fang. 2025.	735
681	evaluating NLP systems on legal document under-	Bookrag: A hierarchical structure-aware index-	736
682	standing tasks, highlighting the need for hierarchical	-based approach for retrieval-augmented genera-	737
683	and long-context modeling in legal digitization.	tion on complex documents. <i>arXiv preprint</i>	738
684	DeepSeek. 2024. DeepSeek-V3 technical report .	<i>arXiv:2512.03413</i> .	739
685	<i>Preprint</i> , arXiv:2412.19437.	Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song	740
686	Yiran Ding, Li Lina Zhang, Chengruidong Zhang,	Han, and Mike Lewis. 2024. Efficient streaming lan-	741
687	Yuanyuan Xu, Ning Shang, Jiahang Xu, Fan Yang,	guage models with attention sinks. <i>arXiv preprint</i>	742
688	and Mao Yang. 2024. Longbench: A bilingual	<i>arXiv:2309.17453</i> .	743
689	benchmark for long context understanding. In <i>Pro-</i>	Yiheng Xu, Minghao Li, Lei Cui, Shaohan Huang,	744
690	<i>ceedings of the Annual Meeting of the Association</i>	Furu Wei, and Ming Zhou. 2020. Layoutlm: Pre-	745
691	<i>for Computational Linguistics (ACL)</i> . Provides em-	training of text and layout for document image	746
692	pirical evidence on the inconsistent performance of	understanding. In <i>Proceedings of the 26th ACM</i>	747
693	LLMs on long-context tasks.	<i>SIGKDD International Conference on Knowledge</i>	748
694	Li Dong and Mirella Lapata. 2016. Language to logi-	<i>Discovery & Data Mining (KDD)</i> . Presents a multi-	749
695	cal form with neural attention. In <i>Proceedings of the</i>	modal pre-trained model combining text and layout	750
696	<i>54th Annual Meeting of the Association for Compu-</i>	information, facilitating the conversion of scanned fi-	751
697	<i>tational Linguistics (ACL)</i> , pages 33–43. Pioneering	ancial or legal documents into structured, machine-	752
698	work on neural semantic parsing using sequence-to-	readable formats.	753
699	tree generation via stack-based decoding.	Manzil Zaheer, Guru Guruganesh, Kumar Avinava	754
700	André V. Duarte, João Marques, Miguel Graça,	Dubey, Joshua Ainslie, Chris Alberti, Santiago On-	755
701	Manuel F. Freire, Lei Li, and Arlindo L. Oliveira.	tanon, Philip Pham, Anirudh Ravula, Qifan Wang,	756
702	2024. Lumberchunker: Long-form narrative	Li Yang, and 1 others. 2020. Big bird: Transformers	757
703	document segmentation. In <i>arXiv preprint</i> .	for longer sequences. <i>Advances in neural informa-</i>	758
704	ArXiv:2406.17526; proposes a method for segment-	<i>tion processing systems</i> , 33:17283–17297.	759
705	ing long narrative documents.	Kaizhong Zhang and Dennis Shasha. 1989. Simple	760
706	Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia,	fast algorithms for the editing distance between trees	761
707	Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, Haofen	and related problems. <i>SIAM Journal on Computing</i> ,	762
708	Wang, and Meng Wang. 2023. Retrieval-augmented	18(6):1245–1262.	763
709	generation for large language models: A survey.	Yifan Zhang, Kento Sone, and Alexandre Alahi. 2024.	764
710	<i>arXiv preprint arXiv:2312.10997</i> .	Streaming llm: Efficient inference with fixed-size at-	765
711	Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paran-	tention windows. In <i>Proceedings of the 41st Inter-</i>	766
712	jape, Michele Bevilacqua, Fabio Petroni, and Percy	<i>national Conference on Machine Learning (ICML</i>	767
713	Liang. 2023. Lost in the middle: How language	2024), volume 235 of <i>Proceedings of Machine</i>	768
714	models use long contexts. In <i>Proceedings of the An-</i>	<i>Learning Research</i> . PMLR.	769
715	<i>annual Meeting of the Association for Computational</i>	Appendix A Prompt Template	770
716	<i>Linguistics (ACL)</i> . The Lost in the Middle phe-	The complete structured prompt used in our	771
717	nomenon is crucial for explaining the failure of Full-	Hierarchical-Aware Path Embedding (HAPE)	772
718	Text Ingestion.	framework is as follows:	773

<System Message>
 You are the parsing engine for the Iterative
 ↳ Structured Content Filling (ISCF) framework.
 You are at iteration $t=\{t\}$, having processed
 ↳ $\{ptr\}$ sentences so far (out of $\{total\}$
 ↳ $total$).
 Your task is to extend the current partial
 ↳ structure by processing the next chunk of
 ↳ text.
 </System Message>

<Global Structure Context>
 Below is the current document tree structure,
 ↳ compressed for context efficiency:
 - Nodes marked with [CONTENT_MASK] have their
 ↳ full text omitted but retain their
 ↳ hierarchical position.
 - Subtrees marked with [SUBTREE_MASK] are
 ↳ collapsed as they are distant from the
 ↳ current focus.

{Indented representation of $T_M^{\{t\}}$ }
 Example format:
 root (title: "Document Root")
 Chapter 1: Introduction (level: 1)
 Section 1.1: Background (level: 2)
 [CONTENT_MASK] (path:
 ↳ root.chapter1.section1.1, role:
 ↳ context-setting, span: 0-42)
 Section 1.2: Problem Statement (level: 2)
 ↳ [ANCHOR NODE]
 SentenceSpan_1 (role: definition, span:
 ↳ 43-58)
 [CONTENT_MASK] (path:
 ↳ root.chapter1.section1.2.s2, role:
 ↳ argument, span: 59-72)
 Chapter 2: Literature Review (level: 1)
 [SUBTREE_MASK] (depth: 3, distance from
 ↳ anchor: 3)
 </Global Structure Context>

<Sequential Context>
 The text you are about to process IMMEDIATELY
 ↳ FOLLOWS the content of the anchor node:
 Anchor Node: $\{n_anchor.title\}$ (level:
 ↳ $\{n_anchor.level\}$)
 Anchor Path: $\{n_anchor.path_descriptor\}$
 Last Processed Content:
 {If anchor has assigned content: [excerpt of last
 ↳ 2--3 sentences]; otherwise: (no content
 ↳ assigned yet)}
 </Sequential Context>

<Current Processing Chunk>
 You are now processing sentences $\{chunk_start\}$ to
 ↳ $\{chunk_end\}$ (absolute indices in document):
 $\{C_i\}$
 </Current Processing Chunk>

<Task Instructions>
 Based on the global structure and sequential
 ↳ context above:
 1. Identify if the new text continues existing
 ↳ structural nodes or begins new ones.
 2. Create new structural nodes with appropriate
 ↳ titles and levels when needed.
 3. Assign each contiguous sentence span in the
 ↳ new chunk to:
 - An existing structural node (specify by
 ↳ path), OR
 - A newly created structural node.

4. For each assigned content span, label its
 ↳ semantic role (topic, definition, assumption,
 ↳ method, evidence, result, conclusion,
 ↳ example).
 5. Output a JSON edit script with the following
 ↳ schema:

```
{
  "new_structural_nodes": [
    {"path": "root.chapter1.section1.3",
     ↳ "title": "Research Questions",
     ↳ "level": 2}
  ],
  "content_assignments": [
    {"span": [73, 85], "parent_path":
     ↳ "root.chapter1.section1.2", "role":
     ↳ "example"},
    {"span": [86, 94], "parent_path":
     ↳ "root.chapter1.section1.3", "role":
     ↳ "question"}
  ],
  "metadata": {
    "chunk_processed": {chunk_id},
    "estimated_coherence": null # Will be
     ↳ computed externally via DBR
  }
}
```

</Task Instructions>
 <Constraints>
 - Maintain bijective mapping: each sentence must
 ↳ be assigned to exactly one content node.
 - Ensure hierarchical consistency: child node
 ↳ spans must be within parent node spans.
 - Titles should be concise (38 words) and
 ↳ descriptive.
 - DO NOT generate new content or modify existing
 ↳ node titles/spans.
 - DO NOT output anything outside the specified
 ↳ JSON schema.
 </Constraints>