# CCF: A Context Compression Framework for Efficient Long-Sequence Language Modeling

**Wenhao Li**[*1, 2] **Bangcheng Sun**[*1] **Weihao Ye** [1] **Tianyi Zhang** [2] **Daohai Yu** [1] **Fei Chao** [1] **Rongrong Ji** [1]

[1]Key Laboratory of Multimedia Trusted Perception and Efficient Computing, Ministry of Education of China, Xiamen University, 361005, P.R. China
[2]Shanghai AI Laboratory

## Abstract

Scaling language models to longer contexts is essential for capturing rich dependencies across extended discourse. However, naïve context extension imposes significant computational and memory burdens, often resulting in inefficiencies during both training and inference. In this work, we propose CCF, a novel context compression framework designed to enable efficient long-context modeling by learning hierarchical latent representations that preserve global semantics while aggressively reducing input redundancy. CCF integrates segment-wise semantic aggregation with key-value memory encoding, forming compact representations that support accurate reconstruction and long-range understanding. To further enhance scalability, we introduce a training-efficient optimization strategy that couples incremental segment decoding with sparse reservoir sampling, substantially reducing memory overhead without degrading performance. Empirical results on multiple long-context language modeling benchmarks demonstrate that CCF achieves competitive perplexity under high compression ratios, and significantly improves throughput and memory efficiency compared to existing approaches. These findings highlight the potential of structured compression for scalable and effective long-context language modeling.

## Introduction

Large language models (LLMs) have achieved remarkable progress in recent years (Yang et al. 2025a; Grattafiori et al. 2024). As these models continue to grow in scale and capability, the ability to process extended contexts has become critical for supporting long-form reasoning (Ding et al. 2024; Wei et al. 2024), multi-document comprehension (Han et al. 2025), and structured knowledge integration (Wu et al. 2024b) in real-world scenarios. However, the deployment of LLMs under long-context conditions remains constrained by the quadratic complexity of the full attention mechanism, which leads to substantial computational and memory overhead as the input length increases (Liskavets et al. 2024). Consequently, a pressing need for a compression framework strikes a balance between preserving semantic fidelity, maintaining training scalability, and retaining architectural simplicity.

A growing line of research has explored context compression (Zhang et al. 2025a) to alleviate long-context language models' scalability bottlenecks. These methods typically employ lightweight encoders to summarize the input into condensed representations (Wang et al. 2024), which are subsequently integrated into the language model via semantic interdependencies (Piero et al. 2025), memory slots (Xie et al. 2025), or key-value caches (Wu et al. 2024a). By distilling salient information from long sequences, these approaches aim to reduce computational complexity while maintaining downstream task performance. Some methods leverage attention-based query mechanisms to extract context-aware features (Zhang et al. 2025c), while others apply cross-modal techniques to align compressed vectors with the language model's embedding space (Deng et al. 2025). Despite these advances, context compression remains fundamentally limited in long-context scenarios. A key challenge is the mismatch between training and inference conditions: many models are trained on short inputs but expected to generalize to significantly longer sequences during deployment, leading to instability and information loss. Furthermore, existing techniques often rely on fixed summarization heuristics or shallow representations that struggle to preserve long-range dependencies critical for reasoning.

To address these challenges, we introduce CCF, a context compression framework designed to enable efficient long-sequence language modeling with minimal compromise in performance. CCF is designed to compress extended input sequences while preserving essential semantics for high-quality language understanding and generation. The framework segments the input into fixed-length chunks, each augmented with a set of learnable special tokens. After being processed by a lightweight encoder, these tokens aggregate local semantic representations through attention, producing a compact summary of each segment. The condensed segment representations are then projected into a shared latent space and integrated into the language model's decoder as key-value caches, enabling effective long-context modeling without reintroducing the full sequence. CCF introduces targeted optimizations to both the encoder and decoder pathways to improve training efficiency and scalability. On the decoder side, it employs a segment-wise incremental training strategy that enables immediate gradient updates, re-

---

ducing memory consumption by avoiding global activation storage. On the encoder side, a sparse optimization mechanism based on reservoir sampling is adopted to retain only essential encoder activations and eliminate redundant computations across segments. This segment-level summarization and selective memory allocation combination enables CCF to support long-context tasks under realistic computational constraints while maintaining semantic fidelity and performance stability. The main contributions of this work are summarized as follows.

- We propose CCF, a context compression framework that enables efficient long-sequence language modeling by combining segment-wise semantic aggregation with key-value memory integration, reducing input redundancy while preserving contextual fidelity.

- We design a training-efficient optimization strategy that integrates segment-wise incremental decoding and reservoir sampling-based sparse encoding, substantially lowering memory and computation overhead without compromising model capacity.

- We conduct comprehensive empirical evaluations across a range of long-context benchmarks, demonstrating that CCF achieves competitive performance under high compression ratios while significantly improving training scalability and inference efficiency over existing approaches.

## Related Works

### Long-Context Large Language Models

Long-context modeling is increasingly critical for enabling extended reasoning (Ding et al. 2024; Wei et al. 2024), memory retention (Wu, Zhao, and Zheng 2024), and multi-document understanding (Han et al. 2025) in LLMs (Zhang et al. 2024b). Positional encoding (Ma et al. 2025) extensions improve input length tolerance without modifying model structure but retain the quadratic cost of full attention, limiting scalability. Architectural optimizations, including sparse attention (Yuan et al. 2025; Lu et al. 2025) and hierarchical designs (Zhu et al. 2024; Zhao, Wu, and Xu 2025), reduce computational overhead yet often face instability when modeling dense contextual interactions (Liu et al. 2025). Retrieval-augmented methods (Shi et al. 2023; Chen et al. 2023) alleviate memory demands by selecting relevant content through external mechanisms, although maintaining retrieval accuracy and global coherence remains challenging. Despite these advances, current solutions struggle to balance efficiency, stability, and semantic fidelity when handling extremely long sequences, highlighting the need for complementary compression-based strategies. Additionally, some parameter-efficient tuning methods are very suitable for optimizing long-context models (Zhang et al. 2025b; Hu et al. 2025; Zhong, Zhao, and Zhou 2025).

### Context Compression and Memory Optimization

Context compression techniques aim to alleviate the inefficiency of long-sequence modeling by summarizing input into compact representations (Dantas et al. 2024). Existing methods span several categories. Attention-based token selection (Fu et al. 2024; Hooper et al. 2024) and memory injection (Han et al. 2024; Kim et al. 2024) approaches distill salient information into a fixed number of vectors or memory slots, which are then integrated into the model through auxiliary tokens or key-value caches. Training-oriented strategies adopt segment-wise optimization (Yang et al. 2025b), or activation pruning (Jiang et al. 2024) to reduce memory footprint during gradient computation. Some methods employ learned compression modules to transform long inputs into latent representations aligned with the LLM embedding space (Leconte et al. 2024). Others leverage autoencoding techniques to pre-train bottleneck structures capable of reconstructing long-range dependencies from condensed signals (Verma 2024). More recent designs integrate retrieval-based priors or sampling-based sparsification to dynamically filter irrelevant context before encoding. While these approaches jointly address inference and training efficiency, they often struggle to maintain stability and generalization when applied to inputs substantially longer than during training. The effectiveness of compression also critically depends on how well the distilled representations preserve both local detail and global structure. These challenges highlight the need for a more robust framework.

## Methodology

### Overall Framework

The proposed architecture is designed to compress and integrate contextual information into a set of latent representations that guide downstream reasoning. The input context is first partitioned into segments, each of which is encoded into a group of task-specific latent tokens. These latent tokens are projected into the decoder's embedding space and serve as compact surrogates for the original input, enabling efficient information flow without overloading the sequence length capacity.

### Streamlined Encoder-Decoder Architecture

To encode context efficiently, we adopt a streamlined encoder-decoder architecture in which a decoder-only language model is repurposed as a segment-level encoder. Instead of channeling the encoder's outputs into the decoder's initial input layer, the architecture establishes a parallel connection between corresponding layers of the encoder and decoder. For each layer, the hidden states associated with special latent tokens are transformed into key-value pairs via a lightweight projection module, which directly populates the decoder's key-value cache. This allows the decoder to attend over compressed contextual memories with minimal latency.

The projection module is parameterized by a set of attention transformations derived from the encoder, augmented with lightweight adaptation modules that introduce task-specific flexibility while keeping parameter overhead minimal. To further enhance the encoder's ability to abstract and summarize long-range dependencies, additional adaptation modules are integrated into its attention flow, targeting the query and value pathways. These adaptation parameters are maintained separately from the base model weights, offering modular control over context summarization and memory abstraction. This dual adaptation mechanism—spanning

both the projection process and the encoder's summarization dynamics—enables precise, layer-wise control over the injection of contextual information into the decoder's reasoning process.

Formally, given an input sequence $X$ segmented into $k$ chunks $X = \{x_1, x_2, ..., x_k\}$, each chunk is processed independently to yield a set of latent representations $\{z_1, z_2, ..., z_k\}$. These representations are then layer-wise projected into key-value embeddings $\{KV_1, KV_2, ..., KV_L\}$, which populate the decoder's cache across $L$ layers. Any remaining fragment $x_{k+1}$, which may fall below the segment length, is processed analogously to maintain coherence across all context inputs.

To optimize the encoding of segment-level information, we design a two-stage process involving pre-filling and dynamic cache updating. In the pre-filling stage, each segment $x_i$ is appended with a sequence of special tokens $\pi$ of fixed length $c$, forming an extended sequence:

$$x_i \leftarrow x_i \oplus \pi, \quad \forall i \in [1, 2, ..., k], \tag{1}$$

where $\oplus$ denotes concatenation along the sequence dimension. These extended sequences are passed through the model's encoder layers to derive the hidden representations $h_i \in \mathbb{R}^{(l+c) \times d}$, where $l$ is the original segment length and $d$ is the embedding dimension. From $h_i$, we extract the final hidden states corresponding to the special tokens, denoted as $h_\pi$. These representations are projected into the key-value space via a learned projection function, producing segment-specific key-value pairs:

$$K_i, V_i \leftarrow \text{Projector}(h_\pi), \quad \forall i \in [1, 2, ..., k]. \tag{2}$$

The keys and values from all segments are concatenated to form the initial KV cache used in subsequent attention computation:

$$K_{\text{cache}} \leftarrow K_1 \oplus K_2 \oplus \cdots \oplus K_k, \quad V_{\text{cache}} \leftarrow V_1 \oplus V_2 \oplus \cdots \oplus V_k. \tag{3}$$

During decoding, new tokens are generated incrementally. Once the residual sequence $x_{k+1}$, formed by appending newly generated tokens to the current input, reaches segment length $l$, it is processed in the same manner as previous segments. The special tokens are appended, hidden states are computed, and new key-value pairs $K_{k+1}$, $V_{k+1}$ are extracted and integrated into the existing cache:

$$K_{\text{cache}} \leftarrow K_{\text{cache}} \oplus K_{k+1}, \quad V_{\text{cache}} \leftarrow V_{\text{cache}} \oplus V_{k+1}. \tag{4}$$

This dynamic update mechanism ensures that the cache evolves coherently with ongoing generation while maintaining efficient contextualization of earlier segments.

## Naive Optimization

**Segments Independence.** Though causal relationships exist among different segments, the compression process for each segment is independent and does not require the involvement of other segments. This allows segments to be processed separately, making a departure from previous approaches such as (Zhang et al. 2024a), where later segments could leverage the fused representations from earlier ones.
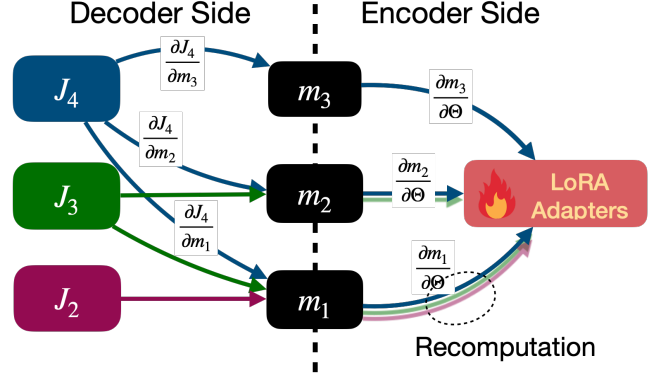


Figure 1: **Recomputation overhead in naive incremental computation.** The gradient of a preceding segment ($m_i$) is repeatedly calculated during the backpropagation, leading to computational redundancy.

**Gradient Expression.** The independent compression of each segment greatly simplifies the computational graph, enabling a concise expression for the parameter gradients. For the language modeling loss associated with the segment $x_j$, denoted as $J_j$, and considering the condensed representations $\{m_i\}_{i=1}^{i=j-1}$ of the $j-1$ preceding segments after passing through the encoder, the gradients can be efficiently calculated. These condensed representations are independent, acting as a relay during backpropagation to pass the gradient from the decoder to the trainable parameters in the encoder and the projectors.

Assuming the LoRA adapters' trainable parameters are $\Theta$, the gradient for the $j$-th segment is:

$$\nabla_\Theta J_j = \sum_{i=1}^{j-1} \frac{\partial J_j}{\partial m_i} \cdot \frac{\partial m_i}{\partial \Theta}. \tag{5}$$

The final gradient is the cumulative sum of the gradients from all segments:

$$\nabla_\Theta = \sum_{j=1}^{k} \nabla_\Theta J_j. \tag{6}$$

**Challenge in Training on Long Sequences.** The gradient calculation as shown in Eq. (5) requires storing $k$ independent forward pass caches due to the independent generation of different $m_i$. Given that our encoder is based on an LLM, these caches are significantly large. For example, even with the use of gradient checkpoint, our test reveals that the memory capacity of an $8\times$ RTX3090 machine is limited to $k \leq 8$, posing a challenge for optimizing long sequences.

## Sparse Optimization for Long Sequences

Running parallel forward and backward propagation for all segments from Eq. (6) can cause excessive GPU memory usage. To address this, we compute each segment one at a time and sum the results. This method, called incremental computation, balances the trade-off between time and memory, enabling training on longer sequences with limited resources.

**The Recomputation.** Upon reviewing Eq. (5) and Eq. (6), we find that naive incremental computation leads to considerable recomputation, as shown in Figure 1. The derivative $\frac{\partial m_i}{\partial \Theta}$ is repeatedly calculated across segments, causing overhead and reducing its benefits.

**Incremental Computation on Decoder Only.** To reduce the excessive recomputation from the naive incremental method on the encoder, we focus incremental computation on the decoder. For each segment, we complete forward passes for both the encoder and decoder, but only apply backpropagation to the decoder. After processing all segments sequentially, $m_i$ sums the gradients from each segment's loss, as shown in the equation below:

$$\nabla_{m_i} = \sum_{j=i+1}^{k} \frac{\partial J_j}{\partial m_i}, \quad i \in [1, 2, ..., k]. \tag{7}$$

After accumulating gradients, we backpropagate through the encoder to determine the final gradients for $\Theta$. This approach, by conducting a single backpropagation through the encoder at the end, eliminates all unnecessary computations. The resulting final gradient matches that shown in Eq. (6). We will now demonstrate their equivalence.

*Proof.* We start by defining an indicator function $I(i,j)$ as:

$$I(i,j) = \begin{cases} 1 & 1 \le j \le k, \ 1 \le i \le j-1, \\ 0 & \text{otherwise}, \end{cases} \tag{8}$$

where the condition can be inverted and explicitly written as:

$$I(i,j) = \begin{cases} 1 & 1 \le i \le k, \ i+1 \le j \le k, \\ 0 & \text{otherwise}. \end{cases} \tag{9}$$

After applying this indicator function, we can rewrite Eq. (6) as:

$$\begin{aligned} \nabla_\Theta &= \sum_{j=-\infty}^{+\infty} \left[ \sum_{i=-\infty}^{+\infty} I(i,j) \cdot \frac{\partial J_j}{\partial m_i} \cdot \frac{\partial m_i}{\partial \Theta} \right] \\ &= \sum_{i=-\infty}^{+\infty} \left[ \sum_{j=-\infty}^{+\infty} I(i,j) \cdot \frac{\partial J_j}{\partial m_i} \cdot \frac{\partial m_i}{\partial \Theta} \right] \\ &= \sum_{i=1}^{k} \left[ \sum_{j=i+1}^{k} \frac{\partial J_j}{\partial m_i} \cdot \frac{\partial m_i}{\partial \Theta} \right] \\ &= \sum_{i=1}^{k} \left[ \left( \sum_{j=i+1}^{k} \frac{\partial J_j}{\partial m_i} \right) \cdot \frac{\partial m_i}{\partial \Theta} \right] \\ &= \sum_{i=1}^{k} \left( \nabla_{m_i} \cdot \frac{\partial m_i}{\partial \Theta} \right). \end{aligned} \tag{10}$$

$\square$

**Sparse Optimization for Encoder Memory Reduction.** By applying incremental computation to the decoder, memory usage is stable since only the cache for one segment is needed, no matter the sequence length. This greatly reduces memory requirements, often leading to significant improvements. However, for the encoder, memory allocation remains high as all segment caches are stored, offering no reduction. Halting optimization here would only cut memory use in half without increasing FLOPs—a minor benefit, especially for very long sequences like 100K tokens. To cut encoder memory further, we suggest using a sparsity budget to limit the number of stored forward pass caches. if the cache limit is exceeded, we use an eviction strategy. For instance, if segment $i$ is removed, we first update $\nabla_\Theta$ by backpropagating $\nabla_{m_i}$, then apply a stop gradient to $m_i$ to halt further accumulation. We will now find the best eviction policy.

**Limitations of Local Window and Random Eviction Policies.** In Figure 2, the simplest eviction strategy is to keep only the most recent segment caches. This method ignores long-term dependencies and relies only on recent data for inference, reducing performance due to biased estimation. In contrast, random sparsity might ideally balance long-term and short-term dependencies, approaching the performance of dense optimization. However, since we cannot recover removed caches, we cannot randomly select segments from the complete set at each step.

**Reservoir Sampling for Eviction Policy.** To achieve efficient memory management, we need an eviction policy that integrates seamlessly with our incremental computation and provides unbiased gradient estimates. reservoir sampling meets both criteria. It offers a natural incremental processing mechanism that aligns with our sequential handling of segments. Additionally, its uniform sampling property ensures that the retained segments are a fair representation of the entire sequence, yielding unbiased estimates of the true gradients. A formal proof demonstrating that this method produces an unbiased gradient estimate is provided in the Appendix (which is in our supplementary materials).

Using reservoir sampling-based sparse optimization on the encoder side maintains constant memory allocation regardless of sequence length, significantly reducing memory requirements while preserving gradient fidelity. We provide the detailed process in Algorithm 1.

## Experimentation

### Datasets

To rigorously evaluate the proposed model's capacity for compressing, understanding, and generating from extended contexts, we employ a multifaceted benchmark suite. The lossless compression ability is quantified through an autoencoding task in which input sequences are passed through the context encoder to yield compact latent representations, subsequently decoded to reconstruct the original content. Higher reconstruction fidelity reflects more effective compression. For long-context language modeling, we measure perplexity across diverse corpora, including PG-19 (Rae
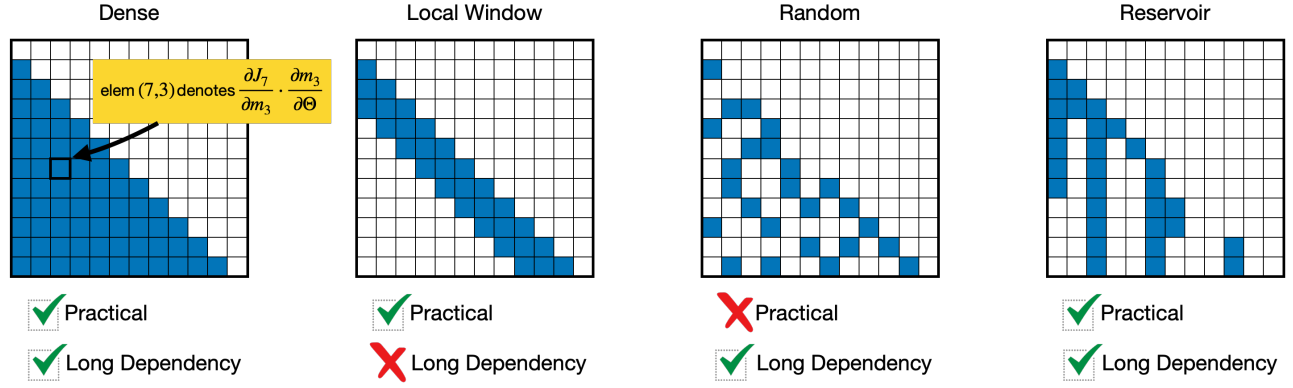
Figure 2: **A comparison of different gradient component eviction strategies.** *Local window* policies neglect long-range dependencies. *Random sampling* is theoretically sound but impractical for streaming data. *Reservoir sampling* offers a practical and unbiased strategy that effectively captures both short- and long-range dependencies.

---

**Algorithm 1: Sparse Optimization with budget size $S$**

1: $R \leftarrow \emptyset$      ▷ Initialize empty reservoir
2: **for** $i = 1$ to $T$ **do**
3:     $m_i \leftarrow \text{Encoder}(x_i)$     ▷ Perform forward pass through the encoder
4:     $J_i \leftarrow \text{Decoder}(x_i, \{m_j\}_{j=1}^{j=i-1})$   ▷ Perform forward pass through the decoder
5:     **backprop**$(J_i, \{m_j\}_{j \in R} \cup m_i)$     ▷ Backpropagate through the decoder only
6:     **if** $i < S$ **then**     ▷ If reservoir is not full, retain forward cache
7:        $R \leftarrow R \cup \{i\}$
8:     **else**     ▷ If reservoir is full, trigger eviction
9:        $j \leftarrow \text{randint}(1, i)$
10:      **if** $j < S$ **then**     ▷ Evict a previously saved segment $t$
11:          $t \leftarrow R[j]$
12:          $R[j] \leftarrow i$
13:          **stop_gradient**$(m_t)$
14:          **backprop**$(m_t, \Theta)$
15:      **else**     ▷ Or discard the incoming segment $i$
16:          **stop_gradient**$(m_i)$
17:          **backprop**$(m_i, \Theta)$
18:      **end if**
19:     **end if**
20: **end for**
21: **for** $i = 1$ to $S$ **do**     ▷ Backpropagate through the remaining segment
22:     $j \leftarrow R[i]$
23:     **backprop**$(m_j, \Theta)$
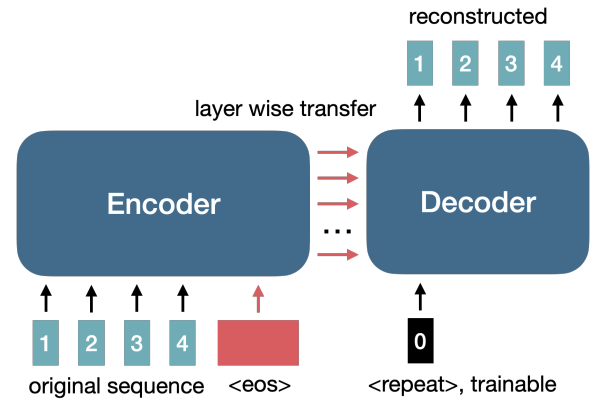24: **end for**

---



Figure 3: **The auto-encoding task.** The model is triggered by `<repeat>` to reconstruct the original input from its compressed representation to evaluate information loss.

formation embedded at arbitrary positions within lengthy sequences, testing both memory and precision. To further assess downstream task performance under long-context conditions, we adopt the LongBench benchmark (Bai et al. 2024) covering single-document and multi-document question answering, summarization, few-shot reasoning, and code generation.

**Implement Details**

The backbone model is built upon LLaMA2-7B (Touvron et al. 2023), which processes fixed-length input segments of 1,024 tokens. By applying compression ratios of 32 and 8, the model effectively extends its maximum context length to approximately 100,000 and 25,000 tokens, respectively. To adapt the pretrained weights efficiently without significantly increasing the number of parameters, LoRA (Hu et al. 2022) fine-tuning is uniformly applied to both the encoder and projection layers with consistent rank and scaling settings. The training procedure unfolds in two stages: first, the model is trained on a billion-token subset of SlimPajama (Weber et al.

et al. 2020), Proof-Pile, and four semantically distinct subsets from SlimPajama (Weber et al. 2024)—Arxiv, Books, Github, and StackExchange—capturing a wide spectrum of linguistic and structural characteristics. The retrieval capability is assessed using the Needle-in-a-Haystack task, which challenges the model to accurately extract target in-

| Model | Ratio | PG19 | Proof | Code |
|---|---|---|---|---|
| ICAE | 8 | 0.79 | 0.72 | 0.81 |
| | 32 | 0.56 | 0.54 | 0.60 |
| AutoCompressor | 8 | 0.67 | 0.64 | 0.71 |
| | 32 | 0.49 | 0.45 | 0.55 |
| CCF (Ours) | 8 | 0.99 | 0.97 | 1.00 |
| | 32 | 0.76 | 0.72 | 0.88 |

Table 1: **Comparison of reconstruction fidelity** (ROUGE-L) across models on 3 long-context corpora: PG-19, Proof-Pile (Proof), and CodeParrot (Code). Our model demonstrates superior performance at both compression ratios.

2024) to establish a strong language understanding foundation; subsequently, the decoder undergoes fine-tuning on a mixed dataset, formatted into conversational prompts following a standardized chat template to facilitate instruction learning. Optimization is performed using a sparse caching mechanism with a reservoir capable of holding up to three token segments, while parameter updates are managed by AdamW paired with a cosine decay learning rate schedule.

## Experimental Results and Analysis

**Auto-Encoding.** To quantify the fidelity of context compression, we evaluate our model on the auto-encoding task, where textual input sequences are compressed into latent representations and subsequently reconstructed by a decoder. Following standard protocol, a special "repeat" token is employed to condition decoding on reconstruction, as illustrated in Figure 3. We adopt ICAE (Ge et al. 2024) and AutoCompressor (Chevalier et al. 2023) as baselines, and evaluate on three domains—PG19 (Rae et al. 2019), Proof-Pile (Hoskinson-Center 2023), and CodeParrot (Unknown 2023)—at two compression levels: 8× and 32×. The results are summarized in Table 1, and a qualitative comparison of original and reconstructed samples is presented in Figure 4.

Our proposed CCF model consistently achieves the highest reconstruction fidelity across all domains and compression ratios. Notably, at 8× compression, it recovers nearly lossless representations, reaching Rouge-L scores of 0.99, 0.97, and 1.00 on PG19, ProofPile, and CodeParrot, respectively. This sharply outperforms ICAE and AutoCompressor, which suffer from significant degradation in semantic fidelity under the same setting.

At the extreme 32× compression, where informative tokens are heavily pruned, CCF still retains robust reconstruction capabilities, achieving a Rouge-L of 0.88 on CodeParrot and above 0.70 on the other two datasets. This performance gain stems from two pivotal innovations. First, the hierarchical latent representation, enhanced with cross-chunk dependency modeling, enables the model to retain semantically salient information across the entire input sequence, even under high compression ratios. Second, the introduction of a fine-tuned repeat-token decoding mechanism allows the decoder to explicitly leverage token recurrence patterns during reconstruction, thereby improving the fidelity of output sequences and ensuring better alignment with the original text.

The strong performance of CCF on both low and high compression settings validates its ability to minimize information loss, a prerequisite for reliable downstream reasoning over compressed context.

| Method | Ratio | 16K Ctx | | | 128K Ctx | | |
|---|---|---|---|---|---|---|---|
| | | PG19 | Proof | Code | PG19 | Proof | Code |
| Act. Beacon | 8 | 9.5 | 5.7 | 2.6 | 10.6 | 5.7 | 3.1 |
| | 32 | 9.9 | 4.8 | 2.7 | 13.2 | 7.2 | 3.5 |
| ICAE | 8 | 11.3 | 9.2 | 4.2 | 18.4 | 12.4 | 5.9 |
| | 32 | – | – | – | – | – | – |
| CCF (Ours) | 8 | 8.3 | 3.9 | 2.3 | 8.4 | 4.4 | 2.5 |
| | 32 | 8.6 | 4.3 | 2.5 | 8.7 | 4.7 | 2.6 |

Table 2: **Perplexity** (lower is better) across compression ratios and context lengths on PG19, ProofPile, and CodeParrot. Missing values are marked as –.

**Contiguous Pretraining.** To further assess the utility of our compressed representations beyond reconstruction, we evaluate their effectiveness in downstream language modeling tasks after the completion of Stage 1 contiguous pretraining. The experiment compares cross-entropy loss on three diverse benchmarks—PG19, ProofPile, and CodeParrot—under both 16K and 128K input contexts.

As shown in Table 2, CCF consistently achieves the lowest language modeling loss across all benchmarks and context lengths, particularly excelling at extreme compression. Under a 32× ratio with 128K context, our model achieves a loss of 8.74 on PG19, while ICAE (Ge et al. 2024) fails to converge, as indicated by the undefined (NaN) values. This robust convergence demonstrates the representational stability of our latent structure under long-range compression. In contrast, Activation Beacon (Zhang et al. 2024a) exhibits moderate performance but shows degradation at high compression levels and longer contexts, especially on PG19 and ProofPile.

The superior performance of CCF arises from its structurally aware compression pipeline. By integrating global semantic cues into the compression stage through cross-chunk interactions and by maintaining alignment during pretraining via an information-preserving latent space, the model ensures that downstream language modeling tasks benefit from reduced information bottlenecks. The stability of the loss across extended context lengths further suggests that the CCF architecture mitigates the compounding effect of information distortion often observed in extreme compression regimes.

**Needle-in-a-Haystack.** To further evaluate the model's capacity to retrieve specific semantic entities from compressed long-form contexts, we adopt the Needle-in-a-Haystack benchmark. This task requires models to locate and recall a uniquely embedded phrase within a long document, testing their ability to preserve token-level fidelity across extended spans. Figure 5 illustrates the pass rates across different compression ratios.
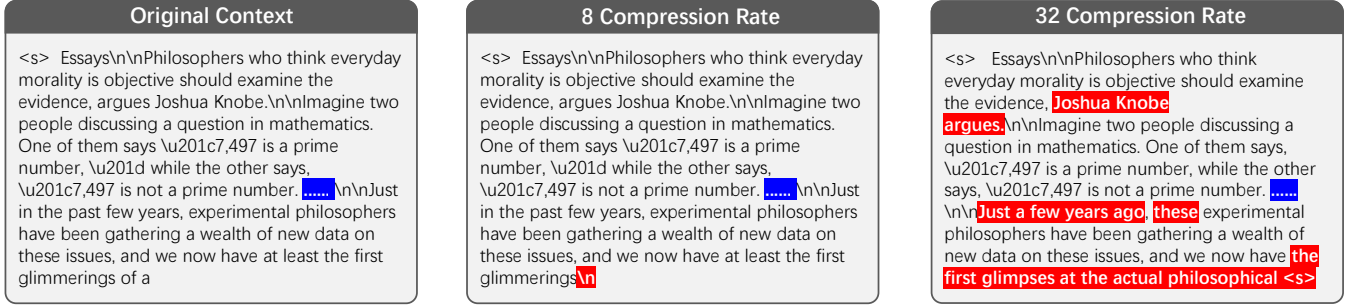
Figure 4: **Auto-encoding case study.** The 8x compression achieves near-perfect reconstruction, while the 32x compression preserves the core meaning with only minor wording variations.

The proposed CCF model achieves perfect retrieval accuracy at a compression ratio of 8×, demonstrating that its latent structure retains fine-grained positional and semantic information without degradation. Even under the more aggressive 32× compression setting, CCF maintains a substantially higher pass rate compared to Activation Beacon (Act. Beacon), revealing its robustness in preserving a retrievable signal across depth-reduced representations. This advantage is primarily driven by the model's compression-aware attention refinement and token recurrence conditioning, which enhance long-range retrieval precision despite drastic sequence shortening.

These results confirm that CCF not only reconstructs local structure accurately but also enables targeted decoding in sparse but semantically dense contexts, a key requirement for effective long-context reasoning.

**LongBench.** To assess the computational efficiency of the proposed compression scheme under realistic long-context conditions, we conduct evaluations on the Longbench. This setting emphasizes not only the quality of language modeling but also system-level metrics such as decoding throughput and memory footprint under various context lengths.

Table 3 compares our method with the standard LLaMA2-7B (Touvron et al. 2023) across 16K and 128K context lengths. At 128K, the baseline experiences a sharp drop in throughput and an exponential increase in KV cache size due to the quadratic growth of attention memory. In contrast, our model—at a 32× compression ratio—achieves over 3× faster decoding and reduces the KV cache memory consumption by more than 30×, while preserving end-task performance.

These results highlight that the model enables efficient inference in high-resolution contexts but also significantly alleviates memory bottlenecks, making it suitable for deployment in resource-constrained or latency-sensitive environments.

## Conclusion and Limitations

This paper introduced CCF, a Context Compression Framework designed to address the efficiency challenges of long-sequence language models. By integrating segment-wise semantic aggregation with efficient optimization strategies, including incremental decoding and sparse reservoir sam-
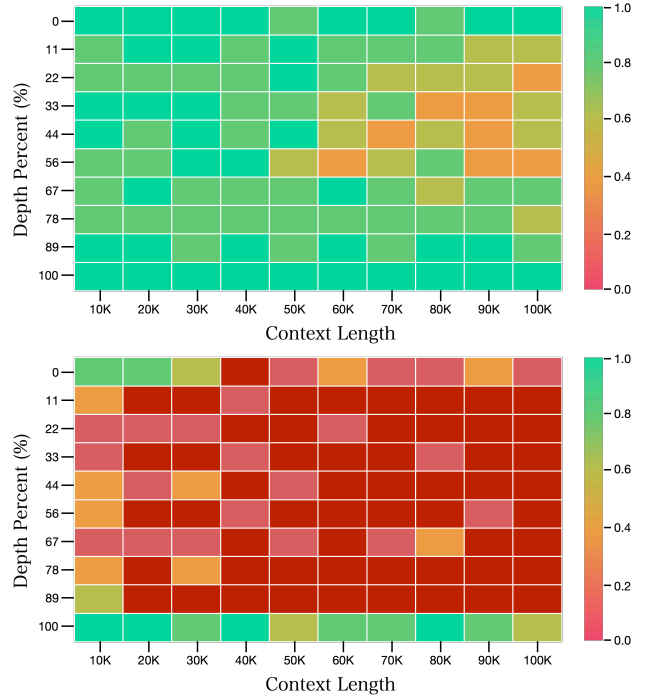


Figure 5: **Needle-in-a-Haystack retrieval accuracy at 32x compression.** CCF (top) demonstrates robust retrieval across all document depths, whereas Activation Beacon (bottom) shows significant performance degradation in the middle of the context.

| Model | Ratio | 16K Context | | 128K Context | |
|---|---|---|---|---|---|
| | | Thpt (tok/s) | KV (GB) | Thpt (tok/s) | KV (GB) |
| LLaMA2-7B | – | 45 | 8 | 18 | 64 |
| CCF (Ours) | 32 | 58 | 0.3 | 54 | 2 |

Table 3: **Inference efficiency on two context length.** CCF significantly improves throughput (Thpt) and reduces KV cache size (KV) compared to the baseline, especially at longer context lengths.

pling, CCF constructs a compact yet expressive representation space.

Despite the promising results, our approach has limitations. An inherent trade-off exists between high compression rates and precise token-level recovery, as aggressive compression can lead to information loss. Second, while effective on standard benchmarks, the model's generalization capacity on paraphrased or atypically structured inputs needs further verification.

# References

Bai, Y.; Lv, X.; Zhang, J.; Lyu, H.; Tang, J.; et al. 2024. LongBench: A Bilingual, Multitask Benchmark for Long Context Understanding. In *ACL*.

Chen, J.; Lin, H.; Han, X.; and Sun, L. 2023. Benchmarking Large Language Models in Retrieval-Augmented Generation. *arXiv preprint arXiv:2309.01431*.

Chevalier, A.; Wettig, A.; Ajith, A.; and Chen, D. 2023. Adapting Language Models to Compress Contexts. In *EMNLP*.

Dantas, P. V.; Sabino da Silva, W.; Cordeiro, L. C.; and Carvalho, C. B. 2024. A comprehensive review of model compression techniques in machine learning. *Applied Intelligence*.

Deng, J.; Jiang, Z.; Pang, L.; Chen, L.; Xu, K.; et al. 2025. Following the Autoregressive Nature of LLM Embeddings via Compression and Alignment. *arXiv preprint arXiv:2502.11401*.

Ding, Y.; Zhang, L. L.; Zhang, C.; Xu, Y.; Shang, N.; et al. 2024. LongRoPE: Extending LLM Context Window Beyond 2 Million Tokens. *arXiv preprint arXiv:2402.13753*.

Fu, T.; Huang, H.; Ning, X.; Zhang, G.; Chen, B.; et al. 2024. MoA: Mixture of Sparse Attention for Automatic Large Language Model Compression. *arXiv preprint arXiv:2406.14909*.

Ge, T.; Hu, J.; Wang, L.; Wang, X.; Chen, S.; et al. 2024. In-context Autoencoder for Context Compression in a Large Language Model. In *ICLR*.

Grattafiori, A.; Dubey, A.; Jauhri, A.; Pandey, A.; Kadian, A.; et al. 2024. The Llama 3 Herd of Models. *arXiv preprint arXiv:2407.21783*.

Han, S.; Xia, P.; Zhang, R.; Sun, T.; Li, Y.; et al. 2025. MDocAgent: A Multi-Modal Multi-Agent Framework for Document Understanding. *arXiv preprint arXiv:2503.13964*.

Han, W.; Zhou, P.; Poria, S.; and Yan, S. 2024. Two are better than one: Context window extension with multi-grained self-injection. *arXiv preprint arXiv:2410.19318*.

Hooper, C.; Kim, S.; Mohammadzadeh, H.; Maheswaran, M.; Paik, J.; et al. 2024. Squeezed Attention: Accelerating Long Context Length LLM Inference. *arXiv preprint arXiv:2411.09688*.

Hoskinson-Center. 2023. Github Repository: Proof-Pile. Github.

Hu, E. J.; Shen, Y.; Wallis, P.; Allen-Zhu, Z.; Li, Y.; et al. 2022. LoRA: Low-Rank Adaptation of Large Language Models. In *ICLR*.

Hu, Z.; Liu, Y.; Zhao, J.; Wang, S.; Wang, Y.; Shen, W.; et al. 2025. LongRecipe: Recipe for Efficient Long Context Generalization in Large Language Models. In *ACL*.

Jiang, B.; Lu, Y.; Lu, G.; and Zhang, B. 2024. QFormer: An Efficient Quaternion Transformer for Image Denoising. In *IJCAI*.

Kim, J.-H.; Yeom, J.; Yun, S.; and Song, H. O. 2024. Compressed Context Memory For Online Language Model Interaction. *arXiv preprint arXiv:2312.03414*.

Leconte, L.; Bedin, L.; Nguyen, V. M.; and Moulines, E. 2024. ReALLM: A general framework for LLM compression and fine-tuning. *arXiv preprint arXiv:2405.13155*.

Liskavets, B.; Ushakov, M.; Roy, S.; Klibanov, M.; Etemad, A.; et al. 2024. Prompt Compression with Context-Aware Sentence Encoding for Fast and Improved LLM Inference. *arXiv preprint arXiv:2409.01227*.

Liu, J.; Zhu, D.; Bai, Z.; He, Y.; Liao, H.; et al. 2025. A Comprehensive Survey on Long Context Language Modeling. *arXiv preprint arXiv:2503.17407*.

Lu, E.; Jiang, Z.; Liu, J.; Du, Y.; Jiang, T.; et al. 2025. MoBA: Mixture of Block Attention for Long-Context LLMs. *arXiv preprint arXiv:2502.13189*.

Ma, X.; Liu, W.; Zhang, P.; and Xu, N. 2025. 3D-RPE: Enhancing Long-Context Modeling Through 3D Rotary Position Encoding. In *AAAI*.

Piero, N.; Cromwell, Z.; Wainwright, N.; and Nethercott, M. 2025. Contextual Reinforcement in Multimodal Token Compression for Large Language Models. *arXiv preprint arXiv:2501.16658*.

Rae, J. W.; Potapenko, A.; Jayakumar, S. M.; Hillier, C.; and Lillicrap, T. P. 2019. Compressive Transformers for Long-Range Sequence Modelling. *arXiv preprint arXiv:1911.05507*.

Rae, J. W.; Potapenko, A.; Jayakumar, S. M.; Hillier, C.; and Lillicrap, T. P. 2020. Compressive Transformers for Long-Range Sequence Modelling. In *ICLR*.

Shi, W.; Min, S.; Yasunaga, M.; Seo, M.; James, R.; et al. 2023. REPLUG: Retrieval-Augmented Black-Box Language Models. *arXiv preprint arXiv:2301.12652*.

Touvron, H.; Martin, L.; Stone, K.; Albert, P.; Almahairi, A.; et al. 2023. Llama 2: Open Foundation and Fine-Tuned Chat Models. *arXiv preprint arXiv:2307.09288*.

Unknown. 2023. Huggingface Dataset: codeparrot/github-code. Huggingface.

Verma, S. 2024. Contextual compression in retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2409.13385*.

Wang, X.; Chen, Z.; Xu, T.; Xie, Z.; He, Y.; et al. 2024. In-Context Former: Lightning-fast Compressing Context for Large Language Model. In *EMNLP*.

Weber, M.; Fu, D. Y.; Anthony, Q.; Oren, Y.; Adams, S.; et al. 2024. RedPajama: an Open Dataset for Training Large Language Models. In *NeurIPS*.

Wei, J.; Yang, C.; Song, X.; Lu, Y.; Hu, N. Z.; et al. 2024. Long-form factuality in large language models. In *NeurIPS*.

Wu, J.; Wang, Z.; Zhang, L.; Lai, Y.; He, Y.; et al. 2024a. Scope: Optimizing key-value cache compression in long-context generation. *arXiv preprint arXiv:2412.13649*.

Wu, S.; Zhao, S.; Yasunaga, M.; Huang, K.; Cao, K.; et al. 2024b. STaRK: Benchmarking LLM Retrieval on Textual and Relational Knowledge Bases. In *NeurIPS*.

Wu, T.; Zhao, Y.; and Zheng, Z. 2024. An Efficient Recipe for Long Context Extension via Middle-Focused Positional Encoding. *arXiv preprint arXiv:2406.07138*.

Xie, R.; Haq, A. U.; Ma, L.; Fang, Y.; Engineer, Z. B.; et al. 2025. Reimagining Memory Access for LLM Inference: Compression-Aware Memory Controller Design. *arXiv preprint arXiv:2503.18869*.

Yang, A.; Li, A.; Yang, B.; Zhang, B.; Hui, B.; et al. 2025a. Qwen3 Technical Report. *arXiv preprint arXiv:2505.09388*.

Yang, Y.; Liu, S.; Rao, C.; An, B.; Shen, T.; et al. 2025b. Dynamic Context-oriented Decomposition for Task-aware Low-rank Adaptation with Less Forgetting and Faster Convergence. *arXiv preprint arXiv:2506.13187*.

Yuan, J.; Gao, H.; Dai, D.; Luo, J.; Zhao, L.; et al. 2025. Native Sparse Attention: Hardware-Aligned and Natively Trainable Sparse Attention. *arXiv preprint arXiv:2502.11089*.

Zhang, P.; Liu, Z.; Xiao, S.; Shao, N.; Ye, Q.; et al. 2024a. Long Context Compression with Activation Beacon. *arXiv preprint arXiv:2401.03462*.

Zhang, P.; Liu, Z.; Xiao, S.; Shao, N.; Ye, Q.; et al. 2025a. Long Context Compression with Activation Beacon. In *ICLR*.

Zhang, X.; Zhao, J.; Yang, Z.; Zhong, Y.; Guan, S.; Cao, L.; et al. 2025b. UORA: Uniform Orthogonal Reinitialization Adaptation in Parameter-Efficient Fine-Tuning of Large Models. In *ACL*.

Zhang, Y.; Huang, Y.; Cheng, N.; Guo, Y.; Zhu, Y.; et al. 2025c. Sentinel: Attention Probing of Proxy Models for LLM Context Compression with an Understanding Perspective. *arXiv preprint arXiv:2505.23277*.

Zhang, Z.; Chen, R.; Liu, S.; Yao, Z.; Ruwase, O.; et al. 2024b. Found in the Middle: How Language Models Use Long Contexts Better via Plug-and-Play Positional Encoding. In *NeurIPS*.

Zhao, Y.; Wu, H.; and Xu, B. 2025. Leveraging Attention to Effectively Compress Prompts for Long-Context LLMs. In *AAAI*.

Zhong, Y.; Zhao, J.; and Zhou, Y. 2025. Low-Rank Interconnected Adaptation across Layers. In *ACL*.

Zhu, Q.; Duan, J.; Chen, C.; Liu, S.; Li, X.; et al. 2024. SampleAttention: Near-Lossless Acceleration of Long Context LLM Inference with Adaptive Structured Sparse Attention. *arXiv preprint arXiv:2406.15486*.