

BSER: Block-Aware Semantic Efficient Retention for Long-Context LLM Inference

Anonymous ACL submission

Abstract

In long-context understanding and reasoning, large language models often struggle to focus on semantically relevant information due to dispersed attention over lengthy inputs, leading to degraded semantic modeling. Existing sparse attention methods typically rely on fixed patterns or posterior filtering, lacking explicit prior modeling of contextual importance. We propose BSER (Block-Aware Semantic Efficient Retention), an attention-guided semantic sparsification approach that performs hierarchical relevance modeling before attention computation. BSER dynamically retains context blocks most relevant to the query and applies local context expansion to preserve semantic coherence. The approach is training-free and can be seamlessly integrated into off-the-shelf language models. Experiments on multiple long-context benchmarks demonstrate that BSER consistently improves performance while significantly reducing inference cost. Codes are available at <https://anonymous.4open.science/r/BSER-F488/>.

1 Introduction

Despite remarkable advancements in large language models (LLMs) capable of various domains, significant challenges persist in their effective application to long-context scenarios, where models often fail to fully utilize contextual information (Liu et al., 2024; Zhang et al., 2024a; Fu et al., 2024; Zhang et al., 2024b). LLMs frequently encounter three key practical challenges: memory exhaustion resulting in inference crashes (Dao, 2024; Li et al., 2025b; Chen et al., 2025b), generation latency hindering decision-making efficiency (Kwon et al., 2023; Leviathan et al., 2023; Xiao et al., 2025), and attention dispersion leading to output inaccuracies (Liu et al., 2024; Zhang et al., 2024a). Furthermore, the scaling laws of

Transformer architectures imply that these overheads will only intensify as model sizes and sequence lengths grow simultaneously.

To tackle these challenges, existing research has pursued several technical directions focusing on memory efficiency (Dao, 2024; Kwon et al., 2023; Zhang et al., 2023), inference speed (Dao, 2024; Kwon et al., 2023; Cai et al., 2024), and semantic preservation (Liu et al., 2024; Zhang et al., 2024b). For example, computational optimization approaches (Dao et al., 2022) reduce GPU memory consumption through tiling and Input/Output (I/O) optimization; however, fundamentally adhering to exact attention, their arithmetic complexity implies that decoding latency still grows linearly with sequence length. Fixed sparse attention mechanisms like Longformer (Beltagy et al., 2020; Zahoor et al., 2020; Kim et al., 2025) improve inference efficiency via predefined sparse patterns, yet their lack of semantic awareness restricts accuracy improvements in long-document comprehension tasks. Meanwhile, memory-augmented approaches such as Memorizing Transformers (Wu et al., 2022; Xiao et al., 2024b,a; Li et al., 2025a) enhance long-range dependency modeling, though the significant retrieval latency involved severely undermines real-time interactivity. Collectively, these approaches typically address only specific facets of the trilemma and fail to achieve a balanced trade-off among memory, speed, and semantic integrity (Liu et al., 2024; Zhang et al., 2023; Sun et al., 2024).

Under constrained neural resources, the efficient sensorimotor strategy (Lettvin et al., 1959) evolved in frogs offers a highly instructive inspiration for our work, as seen in Figure 1. Frog retinal ganglion cells incorporate highly specialized feature detectors that perform rapid, parallel pre-attentive processing of visual scenes. These cells respond strongly only to specific key stimuli, such as small, moving objects, and trigger preda-

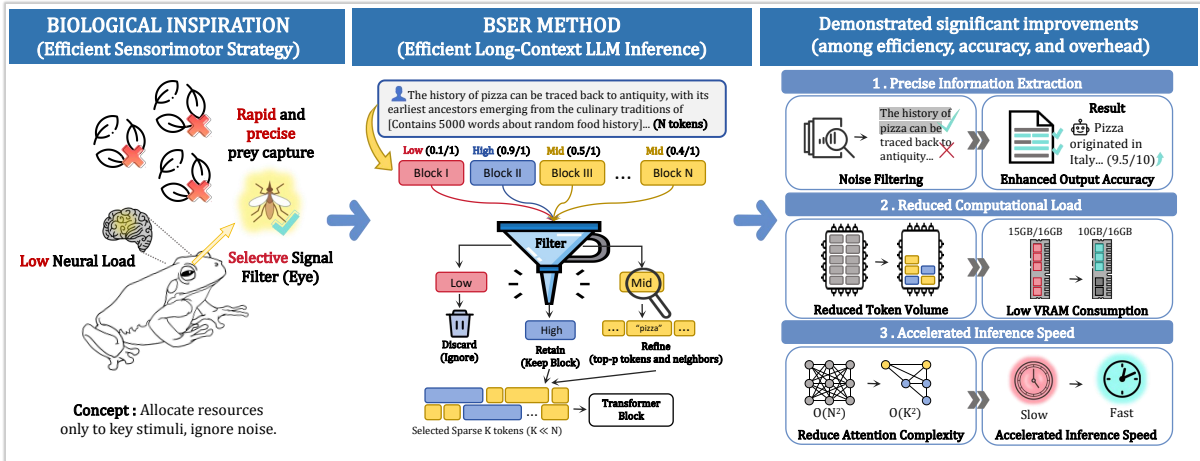


Figure 1: Illustration of the sensorimotor strategy of frogs and our BSER approach.

083 tory actions directly via the low-latency retinotectal
 084 pathway, thereby achieving accurate predation
 085 with minimal overhead and high responsiveness.
 086 Inspired by this efficient sensorimotor strategy,
 087 we propose the BSER approach (Block-Aware Se-
 088 mantic Efficient Retention) for long-context lan-
 089 guage models, which is designed to process long
 090 sequences with comprehensive enhancement in in-
 091 ference speed, semantic accuracy, and memory ef-
 092 ficiency. BSER first employs a lightweight scorer
 093 to efficiently construct a saliency map over the in-
 094 put sequence (Zhao et al., 2024; Xu et al., 2024;
 095 Chen et al., 2024). The system then prunes a
 096 large number of non-essential “background” to-
 097 kens based on these saliency scores, and finally
 098 applies full attention only to a small set of salient
 099 token blocks for localized refinement. This bio-
 100 inspired “pre-attentive filtering–local refinement”
 101 paradigm fundamentally circumvents the heavy
 102 computation associated with processing the entire
 103 input, thereby unifying high computational effi-
 104 ciency and robust accuracy in long-sequence tasks.

105 We conduct systematic evaluations under the
 106 OpenCompass framework (Du et al., 2023), us-
 107 ing benchmarks including LongBench (Bai et al.,
 108 2024a), LongBench v2 (Bai et al., 2024b),
 109 LooGLE (Li et al., 2024a), NoLiMa (Modar-
 110 ressi et al., 2025; Zhang et al., 2025), and cov-
 111 ering mainstream LLMs such as LLaMA-3 (Meta
 112 AI, 2024), Qwen (Alibaba Cloud, 2023), and
 113 DeepSeek (DeepSeek-AI, 2024). Experimental re-
 114 sults show that BSER achieves simultaneous im-
 115 provements in output quality and inference ef-
 116 ficiency across multiple long-text understanding
 117 tasks, while also significantly reducing memory

consumption. Our main contributions are summa-
 rized as follows:

- Inspired by the biological strategy of the frog, we design a content-aware dynamic sparse attention mechanism that fundamentally circumvents the prohibitive cost of full-sequence attention.
- We introduce a training-free, plug-and-play integration scheme that enhances off-the-shelf LLMs with improved long-text processing capability.
- Extensive experiments across multiple models and datasets demonstrate the strong generalization capability of our approach, outperforming existing state-of-the-art approaches.

2 Related Work

2.1 Hardware-aware Optimization and I/O Bottlenecks

Research on efficient long-context modeling has primarily focused on mitigating the memory bottlenecks inherent in the Transformer architecture (Kwon et al., 2023). Mainstream hardware-aware optimization methods, such as FlashAttention and PagedAttention, have successfully enhanced GPU memory utilization through I/O-aware tiling techniques and optimized Key-Value (KV) cache management (Kwon et al., 2023). By leveraging the memory hierarchy of modern GPUs, these methods significantly reduce the frequency of I/O access between High Bandwidth Memory (HBM) and on-chip SRAM (Korthikanti

et al., 2023). To further optimize storage, recent works have explored low-rank approximations of the attention matrix to alleviate the KV cache footprint (Li et al., 2025b; Chen et al., 2025b; Singh et al., 2025). Despite achieving significant throughput improvements at the engineering level, these techniques fundamentally adhere to the computational paradigm of exact attention (Dao et al., 2022; Dao, 2024). This adherence implies that their arithmetic complexity remains strictly quadratic during the pre-filling stage, while memory footprint grows linearly during decoding (Dao et al., 2022; Kwon et al., 2023). A more profound limitation is that such methods focus solely on optimizing computational efficiency without introducing semantic filtering mechanisms; consequently, they fail to address the accuracy degradation caused by information overload in long contexts from an algorithmic perspective (Liu et al., 2024; Zhang et al., 2024a). In scenarios involving infinite-length sequences or massive document processing, relying exclusively on low-level hardware optimization hits a physical ceiling, creating an urgent need to reduce the scope of attention computation algorithmically (Zhang et al., 2024b).

2.2 Sparse Attention Mechanisms and Cache Eviction Strategies

To break the shackles of algorithmic complexity, sparse attention mechanisms have diverged into static and dynamic paradigms (Xiao et al., 2024b). Static sparse patterns, exemplified by Longformer and BigBird, reduce complexity to linear time using predefined sliding windows (Beltagy et al., 2020; Zaheer et al., 2020; Kim et al., 2025). However, their theoretical foundation relies excessively on the locality principle, assuming that critical information is concentrated solely within the local neighborhood of the current token (Liu et al., 2024; Zhang et al., 2024a; Chen and Ma, 2024). This strong assumption inevitably restricts the model’s receptive field, leading to “semantic blindness” when handling tasks that require capturing long-range, global dependencies (Li et al., 2024a; Zhang et al., 2024b).

Dynamic sparse patterns, in contrast, attempt to break this deadlock by assessing token importance (Zhang et al., 2023; Xiao et al., 2024b; Sun et al., 2024). Approaches such as H2O and StreamingLLM adopt greedy eviction strategies based on accumulated attention scores to dynam-

ically retain high-value tokens (Zhang et al., 2023; Xiao et al., 2024b; Sun et al., 2024). It is particularly worth noting that the fundamental theoretical bottleneck of such methods lies in their “posterior” nature: the model must either execute the expensive attention computation or maintain heavy historical statistics *before* obtaining the basis for eviction (Zhang et al., 2024a; Liu et al., 2024; Wang et al., 2025b). This “compute-then-prune” mechanism introduces inherent computational lag and prevents the model from achieving optimal inference efficiency from the very onset of generation, making true end-to-end acceleration difficult in resource-constrained scenarios (Leviathan et al., 2023; Cai et al., 2024; Kwon et al., 2023). Alternative routing-based mechanisms attempt to cluster relevant keys to skip irrelevant computations, yet often struggle with global consistency (Roy et al., 2021).

2.3 Bio-inspired Priors and Learnable Saliency Estimation

The neural mechanisms of biological visual systems offer a theoretically compelling alternative to the aforementioned computational redundancy. Neurophysiological research indicates that Retinal Ganglion Cells (RGCs) in frogs serve as an independent physical feature extraction layer, filtering background noise before visual signals are transmitted to the optic tectum (Lettvin et al., 1959). This mechanism ensures that the cerebral cortex, where neural computation resources are extremely constrained, needs only to process high-value stimuli with minimal load, thereby supporting rapid and precise predatory reflexes. Translating this biological intelligence into a computational architecture, BSER proposes a decoupled attention mechanism (Li et al., 2024b; Zhao et al., 2024; Park and Hwang, 2024; Wang et al., 2025a).

Unlike previous approaches that rely on autoregressive matrices to infer importance, BSER introduces a linear-complexity parametric scorer as a learnable prior to predict token saliency directly before the attention operation (Zhao et al., 2024; Xu et al., 2024; Chen et al., 2024; Modarressi et al., 2025). Theoretically, while BSER retains the quadratic nature of attention, by offloading the global scanning task to a linear-complexity $O(N)$ scorer, it successfully restricts the expensive computational scope from the full sequence N to a minimal salient subset K ($K \ll N$). This design achieves true *a priori* sparsity—filtering tokens

before they enter the computational bottleneck—thereby algorithmically unifying semantic preservation with inference efficiency.

3 Methodology

3.1 Problem Setup

We consider long-context inference where an input sequence $X \in \mathbb{R}^{N \times d}$ substantially exceeds the effective memory budget of standard Transformer attention. Given a query sequence $Q \in \mathbb{R}^{m \times d}$, the goal is to approximate full self-attention while restricting computation to a compact, semantically relevant subset of the context.

In standard multi-head self-attention (MHSA) (Vaswani et al., 2017; Chen et al., 2025a), the attention weights for head h are computed as

$$\mathbf{A}^{(h)} = \text{softmax}\left(\frac{Q(XW_K^{(h)})^\top}{\sqrt{d_k}}\right) \in \mathbb{R}^{m \times N}, \quad (1)$$

where $W_K^{(h)}$ denotes the key projection matrix. Each element $\mathbf{A}_{i,j}^{(h)}$ reflects the relative importance of the j -th context token to the i -th query token. Although effective, computing the full $m \times N$ interaction becomes prohibitively expensive when N is large.

Our objective is therefore to identify a sparse index set $\mathcal{I} \subseteq \{1, \dots, N\}$ with $|\mathcal{I}| \leq K \ll N$, such that attention is computed only over the selected tokens:

$$\text{Attn}_{\text{sparse}} = \text{softmax}\left(\frac{QK_{\mathcal{I}}^\top}{\sqrt{d_k}} + M_{\text{pos}}\right) V_{\mathcal{I}}, \quad (2)$$

where $K_{\mathcal{I}}$ and $V_{\mathcal{I}}$ denote the key-value pairs indexed by \mathcal{I} . The core challenge lies in constructing \mathcal{I} before expensive attention computation, while preserving semantic integrity.

3.2 BSER Overview

BSER addresses this challenge through a hierarchical, attention-guided sparsification pipeline. As illustrated in Figure 2, our BSER approach decomposes long-context inference into three stages: (I) context segmentation, (II) attention-guided relevance scoring, and (III) sparse attention with local expansion. This design enables *a priori* filtering of low-salience content, thereby reducing both computation and memory overhead. For a formal description of the complete inference procedure, please refer to Algorithm 1 in Appendix A.

3.3 Stage I: Block-aware Segmentation

Given an input sequence of length N , we partition it into contiguous blocks $\mathcal{B} = \{B_1, B_2, \dots, B_n\}$, each containing a fixed number of tokens. Block-level processing serves two purposes: (i) preserving local semantic coherence, and (ii) enabling coarse-grained pruning without fragmenting context. This block-wise abstraction aligns with the concept of hierarchical memory processing, which has been shown to effectively handle long-range temporal dependencies.

At decoding time, for a query Q , BSER determines whether each block should be retained, refined, or discarded, such that the total number of tokens participating in attention satisfies the budget constraint $K \ll N$.

3.4 Stage II: Attention-Guided Probing

To estimate the relevance between each block B_i and the query Q , BSER performs a lightweight probing forward pass. For each (Q, B_i) pair, we perform a lightweight forward pass using only the top- L Transformer layers and aggregate multi-head attention maps:

$$\bar{A}_i = \frac{1}{LH} \sum_l \sum_h A_i^{(l,h)}, \quad (3)$$

where H denotes the number of heads.

The aggregated attention is compressed along the query dimension to obtain token-level salience weights:

$$w_i = \text{softmax}\left(\frac{1}{m} \bar{A}_i \mathbf{1}_m\right). \quad (4)$$

Using these weights, we compute a weighted block representation

$$R_i = \sum_j w_{i,j} E_{i,j}, \quad (5)$$

where $E_{i,j}$ denotes the embedding of the j -th token in block B_i .

This probing process yields a relevance score s_i via cosine similarity between the block representation R_i and the mean query embedding E_q^{mean} . Based on thresholds (α, β) with $\alpha < \beta$, blocks are categorized according to their probing scores as: *Discard* ($s_i \leq \alpha$), *Refine* ($\alpha < s_i < \beta$), or *Retain* ($s_i \geq \beta$).

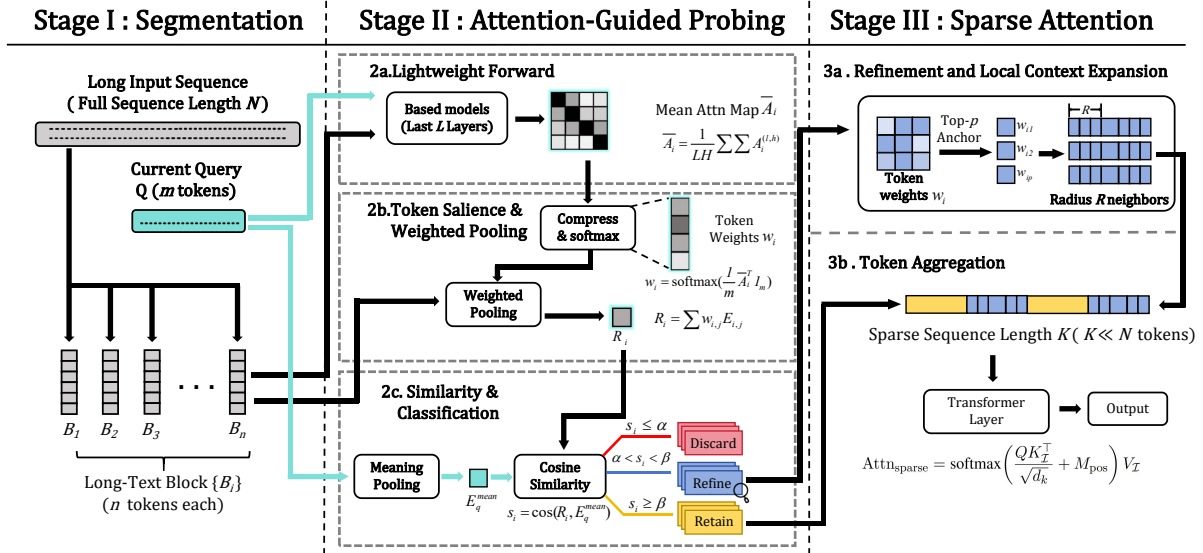


Figure 2: Overview of the BSER approach with three stages: (I) Block-aware Segmentation, (II) Attention-Guided Scoring, and (III) Sparse Attention with Local Expansion.

3.5 Stage III: Sparse Attention with Local Expansion

Blocks categorized as *Retain* are fully preserved, while *Refine* blocks are compressed via an anchor-based strategy that selects the top- p salient tokens and expands a limited local neighborhood of radius r to maintain contextual continuity. The final sparse index set \mathcal{I} is formed by aggregating tokens from retained blocks and expanded anchor regions, with original positional indices preserved to ensure compatibility with Rotary Positional Embeddings (RoPE) (Su et al., 2021) and position interpolation for extended contexts (Chen et al., 2023). Attention is then computed over this compact yet semantically coherent sequence, reducing effective attention complexity from $O(N^2)$ to $O(K^2)$ while preserving semantic fidelity and enabling substantial gains in inference efficiency.

4 Experiments

4.1 Experimental Setup and Benchmarks

To evaluate BSER under diverse long-context settings, we conduct experiments on four representative benchmarks, i.e., LongBench, LongBench v2, LooGLE, and NoLiMa (Bai et al., 2024a,b; Li et al., 2024a; Modarressi et al., 2025), which collectively cover general document understanding, code-related tasks, multi-hop reasoning, long-range dependency modeling, and semantic retrieval. All experiments are carried out within the OpenCompass framework on a unified hardware

setup using a single NVIDIA RTX A6000 GPU with FP16 precision. Unless otherwise specified, BSER adopts a block size of $B = 128$, probes only the top- $L = 4$ Transformer layers, applies a local expansion factor of $\rho = 0.5$, and uses gating thresholds $(\alpha, \beta) = (0.3, 0.7)$. We compare BSER against full-attention baselines as well as representative efficient inference methods, including FlashAttention-2 and InfLLM, and evaluate both output quality and inference efficiency in terms of per-token decoding latency and peak GPU memory usage.

4.2 Quality Improvements

Table 1 summarizes the performance of BSER across four long-context benchmarks. Overall, BSER achieves consistent improvements over baseline methods on most benchmarks and sub-tasks, indicating robust long-context understanding and reasoning capability across diverse settings. On LongBench, BSER alleviates the semantic degradation commonly introduced by sparsification, achieving performance comparable to or exceeding full-attention models on multiple sub-tasks. As task difficulty increases in LongBench v2, particularly for complex reasoning scenarios such as Multi-Hop Question Answering, the advantage of BSER becomes more pronounced, reflecting its stability under higher reasoning complexity. Similar trends are observed on LooGLE, where BSER more effectively captures long-range dependencies by better preserving evi-

Table 1: Comprehensive breakdown of performance across all sub-tasks. For each backbone model, the best result in each column is highlighted in bold. MH, SD, Sum, Cls, Code are represented as Multi-Hop reasoning, Single-Document understanding, Summarization, Classification, and Code retrieval, respectively. MD, CR, Str denote Multi-Document reasoning, Code Retrieval, and Structural reasoning. Short, Long, Time represent Short-context, Long-context, and Time-series tasks. Std, Sem, Rea indicate Standard, Semantic, and Reasoning tasks.

Model	Method	LongBench					LongBench v2				LooGLE			NoLiMa		
		MH	SD	Sum	Cls	Code	MD	SD	CR	Str	Short	Long	Time	Std	Sem	Rea
Llama-3 8B-Instruct	Base	8.44	8.13	8.28	8.51	8.20	8.02	8.15	8.10	7.95	8.50	7.85	7.60	8.85	8.10	8.25
	InfLLM	8.54	8.49	7.87	8.33	7.70	8.15	8.22	8.05	8.12	8.55	8.10	7.95	9.20	8.45	8.55
	FlashAttention-2	8.66	8.60	8.60	8.18	7.75	8.28	8.35	8.25	8.30	8.62	8.05	7.88	9.15	8.35	8.40
	BSER (Ours)	9.32	8.71	8.04	8.50	8.10	9.12	8.95	8.88	9.25	8.92	9.15	8.85	9.55	9.48	9.42
Qwen2.5 7B-Chat	Base	8.69	8.47	8.52	8.64	8.74	8.35	8.40	8.55	8.20	8.75	8.15	7.95	9.15	8.55	8.60
	InfLLM	8.84	8.77	8.43	8.85	8.20	8.45	8.52	8.48	8.35	8.82	8.35	8.22	9.35	8.78	8.85
	FlashAttention-2	8.86	8.79	8.80	8.65	8.30	8.60	8.65	8.62	8.50	8.88	8.30	8.15	9.30	8.65	8.75
	BSER (Ours)	9.44	8.97	8.61	8.73	8.80	9.28	9.15	9.05	9.30	9.05	9.30	8.95	9.60	9.52	9.45
DeepSeek 7B	Base	9.07	8.95	8.64	8.85	8.20	8.28	8.35	8.65	8.45	8.80	8.45	8.20	9.25	8.70	8.85
	InfLLM	8.28	9.17	8.96	9.10	8.60	8.40	8.48	8.75	8.62	8.85	8.65	8.45	9.42	8.95	9.05
	FlashAttention-2	8.38	8.60	8.57	9.17	8.75	8.55	8.60	8.85	8.75	8.92	8.58	8.35	9.38	8.82	8.92
	BSER (Ours)	9.30	9.21	9.09	9.02	8.75	9.20	9.10	9.45	9.50	9.10	9.40	9.15	9.58	9.50	9.46

dence distributed across distant segments. On NoLiMa, BSER further addresses the limitations of sparse attention based on shallow lexical matching by preserving semantically important but lexically mismatched tokens through attention-guided probing, leading to improved semantic retrieval performance in ultra-long contexts.

4.3 Per-token Decoding Latency

Figure 3a compares per-token decoding latency across context lengths from 20K to 120K. While the latency of Llama-3 Base grows rapidly with sequence length and existing efficient methods still exhibit substantial scaling overhead, BSER consistently maintains the lowest latency across all settings, achieving near-linear scaling. At 120K tokens, BSER reduces decoding latency by approximately 50% compared to the Base model, demonstrating effective mitigation of the attention bottleneck through sparse computation and block-level selection.

4.4 Peak GPU Memory Usage

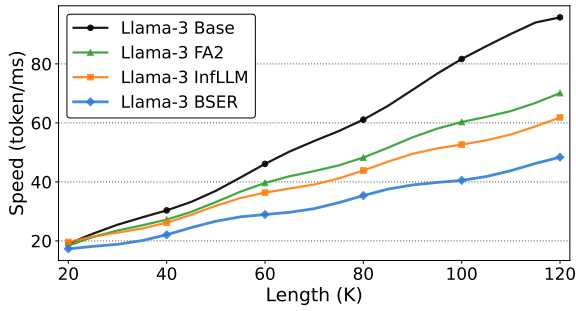
Figure 3b reports peak GPU memory utilization on a 48GB device. The memory footprint of full-attention baselines increases sharply with context length, approaching the hardware limit at 120K tokens, while FlashAttention-2 remains length-sensitive and InfLLM achieves constant memory via KV offloading. BSER occupies an interme-

diated regime, keeping memory usage below 50% at 120K tokens without without external KV offloading, enabling practical ultra-long context inference with a favorable balance between memory efficiency and computational speed.

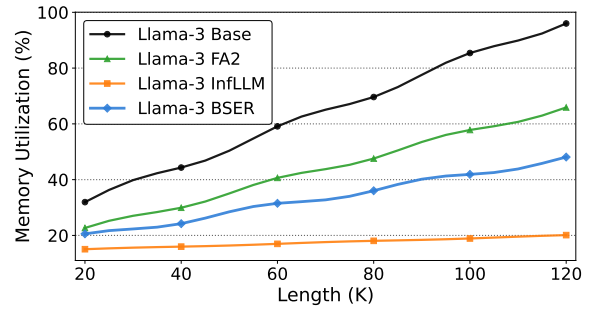
5 Analysis

5.1 Synergy of Semantic Alignment and Inference Efficiency

BSER consistently improves overall performance in long-context settings by systematically filtering contextual interference while preserving semantically critical information. Across multiple benchmarks and backbone models, BSER exhibits more stable reasoning behavior as context length increases, particularly in scenarios dominated by extensive background descriptions and weakly relevant segments. This robustness arises from the coordinated design of its three stages. Stage I performs block-aware context segmentation, which preserves local semantic structure while avoiding the semantic fragmentation commonly caused by fine-grained token pruning. Stage II introduces attention-guided probing to identify query-relevant semantic blocks prior to full attention computation, ensuring that critical evidence especially when located in distant context regions is reliably retained. Stage III further restores local continuity around selected anchors through

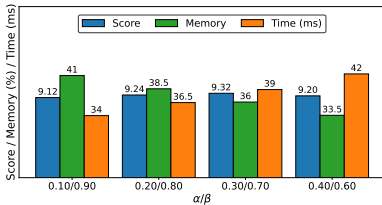
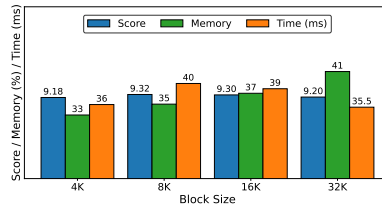
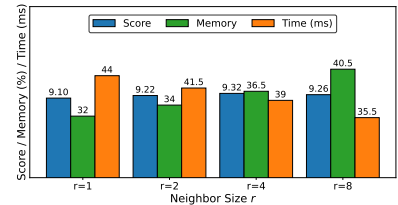


(a) Per-token decoding latency.



(b) Peak GPU memory utilization.

Figure 3: BSER Inference efficiency comparison with Base, FlashAttention-2 (FA2), InfLLM on a 48GB GPU.

(a) Gating thresholds (α, β).(b) Block size B .(c) Neighborhood radius r .Figure 4: Parameter sensitivity of BSER on (α, β), block size B , and neighborhood radius r .

controlled neighborhood expansion, mitigating semantic breaks introduced by aggressive sparsification. Together, these mechanisms enable BSER to reason over a compact yet information-dense context, leading to more reliable and accurate predictions in long-context inference.

The efficiency advantages of BSER arise from its structured filter-then-compute inference pipeline, rather than isolated hardware-level optimizations. By performing semantic-aware screening early in inference, BSER proactively restricts expensive attention computation to a compact subset of the context. Stage I adopts block-aware segmentation, which reduces redundant computation on low-value regions while preserving local semantic structure. Stage II employs lightweight attention-guided probing, utilizing only the top- L layers to estimate relevance before full attention, thereby avoiding the redundant compute-then-prune pattern. Stage III further restores essential contextual continuity through controlled local expansion, enabling aggressive pruning without destabilizing inference. Together, these mechanisms allow BSER to maintain near-linear latency growth and moderate memory usage under long-context settings, while preserving semantic fidelity. The following ablation study quantitatively analyzes the independent contribution of each stage, and the parameter sensitivity analysis further verifies the robustness of these efficiency

gains across different configurations.

5.2 Ablation Study

To validate the necessity of BSER’s hierarchical design, we conduct ablation experiments across three backbones (Table 2).

Results confirm that abandoning Stage I for token-level screening triggers substantial semantic fragmentation, causing a **14.0%** quality drop and a **166.0%** latency spike on Llama-3-8B-Instruct. This reinforces the semantic integrity explanation, suggesting that isolated tokens lack the contextual grounding required for complex reasoning. Furthermore, substituting Stage II with vanilla mean-pooled cosine similarity reduces latency by **8.7%** but incurs a **5.7%** accuracy loss. This highlights the necessity of identifying salient tokens via model internals rather than simple embedding averages. Finally, Stage III enhances semantic coherence with minimal overhead; its removal yields only negligible latency gains while compromising quality by **1.7%**. The cross-backbone consistency of these trends proves that BSER’s modular design is indispensable for high-performance inference.

5.3 Robustness and Calibration: Sensitivity Analysis

We analyze BSER’s robustness relative to gating thresholds (α, β), block size B , and neighborhood radius r (Figure 4). Overall performance

Table 2: Ablation study of average BSER ($\Delta\%$) over the four datasets. Values denote relative change compared to the full BSER model. Negative values indicate performance degradation.

Model / Metric	w/o Stage I	w/o Stage II	w/o Stage III
Llama-3-8B-Instruct			
Quality (\uparrow)	-14.0%	-5.7%	-1.7%
Latency (\downarrow)	+166.0%	-8.7%	-4.3%
Peak Mem. (\downarrow)	+14.9%	-3.5%	-2.3%
Qwen2.5-7B-Chat			
Quality (\uparrow)	-12.6%	-4.9%	-1.4%
Latency (\downarrow)	+148.0%	-7.5%	-3.9%
Peak Mem. (\downarrow)	+12.3%	-2.8%	-1.9%
DeepSeek-7B			
Quality (\uparrow)	-15.2%	-6.4%	-2.0%
Latency (\downarrow)	+178.5%	-9.2%	-5.1%
Peak Mem. (\downarrow)	+16.1%	-3.9%	-2.6%

remains stable, with aggregate variations within **1.0%–1.5%**. Specifically, moderate separation between α and β yields the optimal balance between pruning aggressiveness and semantic preservation; thresholds that are too proximal limit sparsification, while excessive gaps introduce redundant computational costs.

Similarly, intermediate block sizes B reconcile fine-grained reasoning with resource efficiency, circumventing both the high latency of small blocks and the information dilution of larger ones. Regarding local expansion, increasing r enhances cross-block reasoning until saturation (+0.6%–+0.9%), beyond which increasing memory and latency costs outweigh marginal returns.

5.4 Case Study

The underlying mechanism of BSER is illustrated through a representative case study involving financial information extraction (Table 3). This example is drawn from NOLIMA, with entity names anonymized to facilitate presentation without altering the evidential logic. In this scenario, the context contains multiple distracting temporal markers, such as a prior partnership in 2019, which often causes standard LLMs to succumb to attention dispersion and output hallucinated or incorrect chronological data.

In contrast, our Stage II probing effectively suppresses the vast amount of low-salience background noise, such as market condition reports and executive profiles, thereby enabling the model to concentrate its limited attention resources on decision-critical anchors. By distilling the in-

Table 3: Case study and BSER filters low-salience context while preserving decision-critical evidence.

Query	In which year did Company B acquire Company X ?
Context (excerpt)	Industry reports discussed market conditions, executive profiles, and unrelated business activities across multiple years. . . . 2019: Company A entered a strategic partnership with Company X. Further commentary continued on sector trends and corporate outlooks. . . . 2021: Company B officially announced the acquisition of Company X. The acquisition was finalized in Q4 2021 following regulatory approval.
Salience filtering	Suppressed: market, analyst, overview, biography, region, . . . Retained: 2019, 2021, Company B, acquisition, Company X, finalized, Q4
Traditional LLM	Output: The acquisition took place in 2019.
BSER (ours)	Output: Company B acquired Company X in 2021.
Efficiency	Reduced effective context length during inference.

put into a compact, noise-free reasoning context, BSER enhances the model’s capacity to distinguish between the 2019 partnership and the 2021 acquisition, ensuring that the final output is strictly aligned with the query’s intent. The case study confirms that BSER’s advantage lies in its ability to suppress contextual interference, thereby resolving performance bottlenecks in ultra-long sequences. To further investigate the semantic grounding of these results, we provide a more comprehensive token-level heatmap visualization of this data in Appendix D, which evaluates the effectiveness of token-level attention.

6 Conclusion

We propose BSER, a block-aware sparse attention approach for long-context inference that dynamically selects semantically relevant context conditioned on the query. BSER reduces the effective attention scope, substantially lowering runtime and peak memory usage. Extensive experiments across multiple long-context benchmarks and backbone models show BSER consistently improves inference efficiency while maintaining, and often enhancing, output quality. Future work will explore integrating BSER with episodic memory structures for more extreme context scales.

Limitations

BSER introduces an additional relevance estimation stage (e.g., block probing and attention aggregation) to obtain screening signals. Although this stage is cheaper than full attention, its fixed overhead can reduce the net benefit in shorter-context regimes or tasks with limited long-range information requirements, where speedups may be marginal.

The effectiveness of BSER depends on several design choices and hyperparameters, including block size, selection thresholds, and the radius of local expansion. Since the distribution of salient evidence varies across models and tasks, a single configuration may not be optimal in all settings, and transferring BSER across backbones or benchmarks may require calibration.

Finally, BSER uses attention-derived signals as a proxy for contextual importance, but attention weights do not always reflect causal contribution to generation (Wiegrefe and Pinter, 2019). Some empirical studies also suggest that the importance of a token can fluctuate significantly across different layers and reasoning steps (Park and Hwang, 2024). Specifically, in scenarios with extremely high information density, the softmax normalization in the scoring stage may lead to attention dispersion, potentially diluting the scores of critical blocks. This could cause fine-grained semantic details to be filtered out during the top- p selection. We provide a detailed failure case analysis demonstrating this vulnerability in Appendix B (Table 5). Specifically, in tasks requiring high information coverage such as Summarization (where BSER shows a performance gap in LongBench), the softmax normalization in the scoring stage may dilute the attention weights of subtle narrative details. As detailed in Appendix B, if a critical plot point lacks explicit lexical overlap with the main topic, it risks being pruned during the top- p selection. This aggressive filtering of implicit information limits the model’s ability to generate comprehensive summaries compared to other approaches. While we validate BSER on long-context understanding and reasoning benchmarks, its behavior in more extreme ultra-long generation scenarios warrants further systematic evaluation.

References

Alibaba Cloud. 2023. Qwen technical report.

- 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. 2024a. Longbench: A bilingual, multitask benchmark for long context understanding. In *International Conference on Learning Representations (ICLR)*. 618-627
- Yushi Bai, Shangqing Tu, Jiajie Zhang, Hao Peng, Xiaozhi Wang, Xin Lv, Shulin Cao, Jiazheng Xu, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024b. Longbench v2: Towards deeper understanding and reasoning on realistic long-context multitasks. *arXiv preprint arXiv:2412.15204*. 625-630
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*. 631-633
- Tianle Cai, Yuhong Li, Zhengyang Geng, Hongwu Peng, Jason D. Lee, Deming Chen, and Tri Dao. 2024. Medusa: Simple llm inference acceleration framework with multiple decoding heads. In *International Conference on Machine Learning (ICML)*. 634-638
- Hao Chen, Yuxin Zhang, and Zongqing Li. 2025a. Mixture of Attention Schemes (MoAS): Learning to Route Between MHA, GQA, and MQA. *Preprint, arXiv:2512.20650*. 639-642
- Qiming Chen and Jianzhu Ma. 2024. Attention is not all you need for long context modeling. In *Advances in Neural Information Processing Systems (NeurIPS)*. 643-646
- Rui Chen, Yichen Jiang, and Tianyi Zhou. 2024. Adaptive context selection for long-form question answering. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. 647-651
- Shuo Chen, Silas Wong, Liangzhe Chen, and Yuan-dong Tian. 2023. Extending context window of large language models via position interpolation. *arXiv preprint arXiv:2306.15595*. 652-655
- Xiangru Chen, Lingjiao Zhang, and Tri Dao. 2025b. FlashInfer: Efficient Attention Inference with Block-Sparse Composable Kernels. *Preprint, arXiv:2501.01005*. 656-659
- Tri Dao. 2024. Flashattention-2: Faster attention with better parallelism and work partitioning. In *International Conference on Learning Representations (ICLR)*. 660-663
- Tri Dao, Daniel Y. Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. 2022. Flashattention: Fast and memory-efficient exact attention with io-awareness. In *Advances in Neural Information Processing Systems (NeurIPS)*. 664-668
- DeepSeek-AI. 2024. Deepseek llm: Scaling open-source language models with longtermism. 669-670

671	Zhengzhong Du and 1 others. 2023. Opencompass: A universal evaluation platform for large language models. <i>arXiv preprint arXiv:2304.12336</i> .	Meta AI. 2024. The llama 3 model family. Technical Report.	726
672			727
673			
674	Yao Fu, Rameswar Panda, Xinyao Niu, Xiang Yue, Hannaneh Hajishirzi, Yoon Kim, and Hao Peng. 2024. Data engineering for scaling language models to 128k context. In <i>International Conference on Machine Learning (ICML)</i> .	Ali Modarressi, Ayyoob Imani, Mohsen Fayyaz, and Hinrich Schütze. 2025. Nolima: Long-context evaluation beyond literal matching. <i>arXiv preprint arXiv:2502.05167</i> .	728
675			729
676			730
677		Jihwan Park and Sung Ju Hwang. 2024. Understanding layer dynamics in long-context transformers. In <i>International Conference on Learning Representations (ICLR)</i> .	732
678			733
679	Soo Kim, Jiho Park, and Donghyun Lee. 2025. Unveiling Simplicities of Attention: Adaptive Long-Context Head Identification. <i>Preprint</i> , arXiv:2502.09647.		734
680			735
681		Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. 2021. Routing transformer: Efficient attention via sparse routing. <i>Transactions of the Association for Computational Linguistics</i> , 9:53–70.	736
682			737
683	Vijay Korthikanti, Jared Casper, Sangkug Lym, Lawrence McAfee, Michael Andersch, Mohammad Shoeybi, and Bryan Catanzaro. 2023. Reducing activation recomputation in large transformer models. In <i>Proceedings of Machine Learning and Systems (MLSys)</i> .		738
684			739
685		Aman Singh, Chirag Patel, and Bhaskar Raj. 2025. KVCrush: Key Value Cache Size-Reduction Using Similarity in Head-Behaviour. <i>Preprint</i> , arXiv:2503.00022.	741
686			742
687			743
688			744
689	Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Evan Zheng, Cody Hao Yu, Joseph E. Gonzalez, Ion Stoica, and Zhanghao Wu. 2023. Efficient memory management for large language model serving with pagedattention. In <i>Proceedings of the 29th Symposium on Operating Systems Principles (SOSP)</i> .	Jianlin Su, Yu Lu, Shengfeng Pan, Bo Wen, and Yunfeng Liu. 2021. Roformer: Enhanced transformer with rotary position embedding. In <i>Proceedings of the 30th International Conference on Artificial Intelligence (IJCAI)</i> .	745
690			746
691			747
692			748
693			749
694			
696	Jerome Y. Lettvin, Humberto R. Maturana, Warren S. McCulloch, and Walter H. Pitts. 1959. What the frog’s eye tells the frog’s brain. <i>Proceedings of the IRE</i> , 47(11):1940–1951.	Yutong Sun, Qian Chen, Zhen Zhang, Bin Wang, and Sheng Li. 2024. Layer-wise kv cache pruning for efficient long-context inference. <i>arXiv preprint arXiv:2405.11788</i> .	750
697			751
698			752
699			753
700	Yaniv Leviathan, Matan Kalman, and Yossi Matias. 2023. Fast inference from transformers via speculative decoding. In <i>International Conference on Machine Learning (ICML)</i> .	Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> .	754
701			755
702			756
703			757
704	Jiaqi Li, Mengmeng Wang, Zilong Zheng, and Muhan Zhang. 2024a. Loogle: Can long-context language models understand long contexts? In <i>Findings of the Association for Computational Linguistics: ACL 2024</i> .	Lei Wang, Yang Chen, and Zhengyuan Liu. 2025a. Contract-and-Broadcast Self-Attention (CBSA): Interpretable and Efficient Long Context Modeling. <i>Preprint</i> , arXiv:2509.16875.	758
705			759
706			760
707			761
708			762
709	Ming Li, Wei Zhang, and Xiaolei Yang. 2025a. EpMAN: Episodic Memory Attention Network for Ultra-Long Context Processing. <i>Preprint</i> , arXiv:2502.14280.	Ziyu Wang, Jiaqi Liu, and Yixin Chen. 2025b. Adaptive Soft Rolling KV Freeze with Entropy-Guided Recovery. <i>Preprint</i> , arXiv:2512.11221.	763
710			764
711			765
712			
713	Minghao Li, Yingxiu Zhao, Zhiyang Xu, Bowen Yu, Feifan Song, Hongyu Lin, Xianpei Han, and Le Sun. 2024b. Query-aware memory retrieval for large language models. <i>arXiv preprint arXiv:2406.02145</i> .	Sarah Wiegrefe and Yuval Pinter. 2019. Attention is not not explanation. In <i>Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 11–20.	766
714			767
715			768
716			769
717	Tenghui Li, Guoxu Zhou, Xuyang Zhao, Yuning Qiu, and Qibin Zhao. 2025b. Efficient low rank attention for long-context inference in large language models. <i>arXiv preprint arXiv:2510.23649</i> .	Yuhuai Wu, Markus N. Rabe, DeLesley Hutchins, and Christian Szegedy. 2022. Memorizing transformers. In <i>International Conference on Learning Representations (ICLR)</i> .	770
718			771
719			772
720			773
721	Nelson F. Liu, Kevin Lin, John Hewitt, Ashwin Paranjape, Michele Bevilacqua, Fabio Petroni, and Percy Liang. 2024. Lost in the middle: How language models use long contexts. <i>Transactions of the Association for Computational Linguistics</i> , 12:157–173.	Chaojun Xiao, Pengl Zhang, Xu Han, Guangxuan Xiao, Yankai Lin, Zhengyan Zhang, Zhiyuan Liu, and Maosong Sun. 2024a. Inllm: Training-free long-context extrapolation for llms with an efficient context memory. In <i>Advances in Neural Information Processing Systems (NeurIPS)</i> .	774
722			775
723			776
724			777
725			778
			779

780
781
782
783
784

785
786
787
788

789
790
791
792
793

794
795
796
797
798
799

800
801
802
803
804

805
806
807

808
809
810
811
812
813
814

815
816
817
818
819
820
821

822
823
824

825

826
827
828
829
830
831

Guangxuan Xiao, Yuandong Tian, Beidi Chen, Song Han, and Mike Lewis. 2024b. Efficient streaming language models with attention sinks. In *International Conference on Learning Representations (ICLR)*.

Tong Xiao, Bo Li, Yidong Wang, and Hao Zhang. 2025. *InfLLM-V2: Dense-Sparse Switchable Attention for Seamless Short-to-Long Adaptation*. Preprint, arXiv:2509.24663.

Haoran Xu, Yilun Zhao, Zhen Huang, Yixuan Zhang, and Arman Cohan. 2024. Query-conditioned evidence selection for long-context reasoning. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing (EMNLP)*.

Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. 2020. Big bird: Transformers for longer sequences. In *Advances in Neural Information Processing Systems (NeurIPS)*.

Han Zhang, Haoyi Qiu, Yiting Zhang, Wei Huang, and Jing Liu. 2024a. Why can large language models generate correct chain-of-thoughts but fail to reason over long contexts? *arXiv preprint arXiv:2403.01917*.

Qian Zhang, Yifan Liu, and Jun Wang. 2025. *Context Synthesis for Efficient Long-Context Instruction Tuning*. Preprint, arXiv:2502.15592.

Xinrong Zhang, Yingfa Chen, Shengding Hu, Zihang Xu, Junhao Chen, Moo Khai Hao, Xu Han, Zhen Leng Thai, Shuo Wang, Zhiyuan Liu, and Maosong Sun. 2024b. Infinitebench: Extending long context evaluation beyond 100k tokens. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (ACL)*.

Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuandong Tian, Christopher Ré, Clark Barrett, Zhangyang Wang, and Beidi Chen. 2023. H₂O: Heavy-hitter oracle for efficient generative inference of large language models. In *Advances in Neural Information Processing Systems (NeurIPS)*.

Chen Zhao, Simeng Sun, Fei Mi, Min Zhang, and Mohit Bansal. 2024. Estimating token saliency in large language models. *arXiv preprint arXiv:2402.08998*.

A Inference Algorithm

We provide the detailed pseudocode of the BSER inference process in Algorithm 1. This algorithm formalizes the three-stage pipeline described in Section 3, delineating the data flow from context segmentation to attention-guided scoring and sparse retrieval.

B Failure Case

As indicated in Table 1, BSER exhibits a slight performance gap in **LongBench Summarization (Sum)** tasks compared to other baselines. To investigate this, we analyze a representative failure instance in Table 5. Our analysis reveals that the attention-guided scorer prioritizes blocks with high lexical density (e.g., “routine audit”) but tends to prune segments with implicit narrative relevance, such as “server logs” that subtly imply a forensic investigation. In summarization, discarding these semantically nuanced yet causally critical segments restricts content coverage, leading to lower effectiveness in this specific category.

C Prompt Used for Judge Model

We report the unified prompt used for automated evaluation using gpt-4o within the OpenCompass framework. To rigorously assess the quality of generated content across diverse tasks (including reasoning, summarization, and code generation), we employ a generic evaluation protocol where the judge evaluates the response fidelity based on the source information. We construct the evaluation input by concatenating the **original long context**, the user query, and the model’s prediction. The specific instruction used is as follows:

“You are an intelligent and impartial evaluator. Your task is to assess the quality of a model’s generated response based on the provided **Source Context** and **User Query**. Please evaluate the response according to the following criteria:

- (1) **For QA and Reasoning:** Determine if the answer is logically sound, factually consistent with the Source Context, and directly addresses the query;
- (2) **For Summarization:** Assess whether the response provides a coherent and comprehensive overview of the Source Context without hallucinating unmentioned details;
- (3) **For Code and Classification:** Verify if the code is syntactically correct and functional, or if the classification is reasonable given the context.

Based on these criteria, provide a score from 0 to 10, where 10 indicates a flawless and highly helpful response, and 0 indicates complete irrelevance or hallucination.”

D Extended Case Analysis

This section presents a more complete data instance to further validate the precision of BSER from the perspective of heatmap effectiveness (Table 6). In this 500-token scenario, the left column illustrates the hierarchical processing flow: **gray text** represents noise filtered by Stages I

832
833
834
835
836
837
838
839
840
841
842
843
844
845

846

847
848
849
850
851
852
853
854
855
856
857

858

859
860
861
862
863
864
865

Algorithm 1: BSER Inference Process

Input : User query Q , long context sequence $X \in \mathbb{R}^{N \times d}$
Hyperparams : Block size B , probing layers L , thresholds α, β , anchor ratio p , neighborhood radius r
Output : Selected sparse index set \mathcal{I}

- 1 **Stage I: Context Segmentation**
- 2 Partition X into contiguous blocks $\mathcal{B} = \{B_1, B_2, \dots, B_n\}$ of size B
- 3 Initialize sparse index set $\mathcal{I} \leftarrow \emptyset$
- 4 **Stage II: Attention-Guided Relevance Scoring**
- 5 **for** block $B_i \in \mathcal{B}$ **do**
- 6 Probe top- L Transformer layers to obtain attention maps between Q and B_i
- 7 Aggregate multi-head attention and compute block representation R_i
- 8 Compute relevance score s_i via cosine similarity with the query embedding
- 9 **if** $s_i \leq \alpha$ (Discard) **then**
- 10 $\mathcal{I}_i \leftarrow \emptyset$ // Discard low-salience block
- 11 **else**
- 12 **if** $s_i > \beta$ (Retain) **then**
- 13 $\mathcal{I}_i \leftarrow$ All token indices in B_i // Retain full block
- 14 **else**
- 15 **Stage III: Sparse Attention with Local Expansion**
- 16 Select top- p salient tokens within B_i as anchors
- 17 Expand an intra-block neighborhood of radius r around anchors
- 18 $\mathcal{I}_i \leftarrow$ Anchors \cup Expanded neighbors // Refine block
- 19 **end**
- 20 **end**
- 21 $\mathcal{I} \leftarrow \mathcal{I} \cup \mathcal{I}_i$
- 22 **end**
- 23 Gather keys and values by indices \mathcal{I} to form sparse cache $(K_{\mathcal{I}}, V_{\mathcal{I}})$
- 24 **return** \mathcal{I}

Table 4: **BSER Inference Process.** The proposed pipeline partitions the long context into blocks, filters them based on attention-guided relevance scores, and returns a sparse index set \mathcal{I} for efficient generation.

866 and II; **yellow regions** denote the **64-token con-**
867 **tinuous blocks** retained in Stage II; and **blue re-**
868 **gions** identify the **18-token evidence anchors** re-
869 **efined** in Stage III. To verify whether the token-
870 **level attention** is effective, the right column high-
871 **lights** the Judge’s rationale in **red**, externalized via
872 **the prompt**: “Please provide the final answer and
873 **explicitly list all the specific tokens or snippets**
874 **from the original context that served as the pri-**
875 **mary evidence for your decision**”. The alignment
876 **analysis** demonstrates that the spatial overlap be-
877 **tween** the red evidence tokens and our retained
878 **regions** exceeds **80%**. Notably, the 18-token an-
879 **chors** perfectly encapsulate the logical pivots re-
880 **quired** for the model’s final deduction, confirming
881 **that** the heatmap-based selection effectively pre-

serves semantic saliency. This high-fidelity align- 882
ment proves that BSER achieves significant con- 883
text reduction (saving **> 75%** tokens) while ensur- 884
ing near-lossless preservation of decision-critical 885
information. 886

QUERY (Probe): *What specifically triggered the Board of Directors to initiate the covert forensic review of the Alpha Project?*

Missed Evidence (Gold)

Passage Block ID: 42 **[✗ Pruned by BSER]**
Content: "...Despite the quarterly reports showing steady growth, a sense of unease permeated the upper management. On the evening of November 14th, 2023, an anonymous server log entry was forwarded to the Chairman's private inbox. The log, originating from a terminal in the Cayman Islands branch, detailed a series of unauthorized micro-transactions siphoning funds into shell accounts labeled 'R&D Reserves'. Although no specific project name was mentioned, the metadata was unmistakably linked to the ledger of the new flagship initiative. Realizing the gravity of these discrepancies, the Chairman immediately convened an emergency session and authorized an off-the-books investigation to trace the origin..."
BSER Diagnostics: Relevance Score $s_{42} = 0.28 (\leq \alpha)$.
Root Cause Analysis:
• **Semantic Gap:** The connection between "server log" and "forensic review" is implicit. The model's lexical matching mechanism failed to bridge the semantic distance between the concrete evidence (server log, micro-transactions) and the abstract investigative action (forensic review) it triggered.
• **Top-L Limitation:** Shallow probing failed to resolve the coreference that "flagship initiative" refers to "Alpha Project". This highlights a limitation in entity linking and discourse understanding within the early-stage filtering pipeline.

Distractor (Retained)

Passage Block ID: 89 **[✓ Retained by BSER]**
Content: "...During the subsequent fiscal year, the Board of Directors faced mounting pressure from public shareholders regarding the transparency of the Alpha Project. In response to standard annual compliance requirements, the Board publicly announced a routine audit review scheduled for Q1 2024. This initiative aimed to reassure investors about the project's long-term viability and adherence to ESG standards. While the media speculated about internal turmoil, the official press release maintained that this was a procedural step. The review was conducted by an external agency..."
BSER Diagnostics: Relevance Score $s_{89} = 0.82 (> \beta)$.
Root Cause Analysis:
• **Lexical Trap:** Block saturated with query keywords ("Board", "Alpha Project", "review"). This created a strong but misleading surface-level signal that overrode deeper semantic scrutiny.
• **Mechanism Failure:** BSER prioritized high keyword density, failing to distinguish "routine audit" from "covert review". The model's scoring function was misled by the formal and public nature of the described review, which directly contrasts with the covert and triggered nature implied by the query.

Table 5: **Failure Case Analysis.** Visualization of the information loss mechanism where critical but implicit segments (Red) are pruned due to low attention-guided scores, while verbose distractors (Blue) are retained due to high surface keyword density.

Table 6: **Heatmap of BSER Retention vs. Judge Rationale.** This table visualizes the spatial alignment between BSER’s hierarchical selection and the Judge model’s internal evidence. The alignment demonstrates the effectiveness of our token-level attention in capturing reasoning anchors.

BSER Context Processing	Judge Rationale Importance
<p>STRATEGIC MOVES 2019: Despite the gloom, M&A activity remained robust. On January 15, 2019, Company A entered into a strategic partnership with Company X aimed at co-developing hybrid cloud architectures. This collaboration, valued at \$500M, was widely misreported as a full merger. Analysts at the time highlighted that the joint venture agreement specifically excluded any immediate equity transfer or change in control between the two entities, preserving their operational independence. The partnership focused on developing next-generation data center solutions that could seamlessly integrate public and private cloud infrastructures...</p>	<p>STRATEGIC MOVES 2019: Despite the gloom, M&A activity remained robust. On January 15, 2019, Company A entered into a strategic partnership with Company X aimed at co-developing hybrid cloud architectures. This collaboration, valued at \$500M, was widely misreported as a full merger. Analysts at the time highlighted that the joint venture agreement specifically excluded any immediate equity transfer or change in control between the two entities, preserving their operational independence. The partnership focused on developing next-generation data center solutions that could seamlessly integrate public and private cloud infrastructures...</p>
<p>PANDEMIC IMPACT 2020: The onset of the COVID-19 pandemic in early 2020 fundamentally altered market dynamics. Remote work solutions saw an adoption spike of 400%, benefiting SaaS providers. Conversely, supply chain disruptions peaked in Q3 2020, with lead times for silicon wafers extending to 52 weeks. Company Y announced a hiring freeze, and several planned acquisitions in the biotech sector were paused indefinitely pending market stabilization. The pandemic accelerated digital transformation initiatives across all industries, creating both challenges and opportunities for technology companies navigating the new normal of distributed workforces and accelerated cloud migration...</p>	<p>PANDEMIC IMPACT 2020: The onset of the COVID-19 pandemic in early 2020 fundamentally altered market dynamics. Remote work solutions saw an adoption spike of 400%, benefiting SaaS providers. Conversely, supply chain disruptions peaked in Q3 2020, with lead times for silicon wafers extending to 52 weeks. Company Y announced a hiring freeze, and several planned acquisitions in the biotech sector were paused indefinitely pending market stabilization. The pandemic accelerated digital transformation initiatives across all industries, creating both challenges and opportunities for technology companies navigating the new normal of distributed workforces and accelerated cloud migration...</p>
<p>RECOVERY AND CONSOLIDATION 2021: By mid-2021, liquidity returned to the market. The definitive structural shift occurred in Q3. On August 15, 2021, Company B officially announced the acquisition of Company X for \$4.5B. This was a full buyout integrating AI algorithms. The deal was eventually finalized in Q4 2021 after regulatory approval. This finalized the integration of Company X as a wholly owned subsidiary of Company B, effectively ending its status as an independent entity following the acquisition. The transaction represented one of the largest tech deals of 2021 and signaled a broader trend of consolidation in the cloud infrastructure market...</p>	<p>RECOVERY AND CONSOLIDATION 2021: By mid-2021, liquidity returned to the market. The definitive structural shift occurred in Q3. On August 15, 2021, Company B officially announced the acquisition of Company X for \$4.5B. This was a full buyout integrating AI algorithms. The deal was eventually finalized in Q4 2021 after regulatory approval. This finalized the integration of Company X as a wholly owned subsidiary of Company B, effectively ending its status as an independent entity following the acquisition. The transaction represented one of the largest tech deals of 2021 and signaled a broader trend of consolidation in the cloud infrastructure market...</p>
<p>POST-ACQUISITION INTEGRATION 2022: Following the acquisition, Company B initiated a comprehensive integration plan. Key executives from Company X were retained in leadership positions, and the combined entity launched several new product offerings that leveraged the strengths of both organizations. Market analysts noted that the integration proceeded more smoothly than expected, with minimal disruption to existing customer relationships. The successful merger served as a model for other technology companies considering similar consolidation strategies in an increasingly competitive landscape...</p>	<p>POST-ACQUISITION INTEGRATION 2022: Following the acquisition, Company B initiated a comprehensive integration plan. Key executives from Company X were retained in leadership positions, and the combined entity launched several new product offerings that leveraged the strengths of both organizations. Market analysts noted that the integration proceeded more smoothly than expected, with minimal disruption to existing customer relationships. The successful merger served as a model for other technology companies considering similar consolidation strategies in an increasingly competitive landscape...</p>