

# Scratch-1: Query-Anchored Hierarchical Compression for Mitigating Context Fragmentation in Long-Context LLMs

Ahmet Demirbas

*Caltech*

ademirba@caltech.edu

Taaha Khan

*Purdue University*

taahakhan@purdue.edu

Ashank Shah

*Caltech*

ashah@caltech.edu

Arjun Malpani  
*Stanford University*  
amalpani@stanford.edu

## Abstract

Long-context language models incur substantial inference costs proportional to input sequence length, creating an economic “token tax” that limits practical deployment. TheTokenCompany’s Bear-1 has established a strong foundation for prompt compression, achieving impressive token reduction through perplexity-based pruning. We present **Scratch-1**, a novel three-stage compression pipeline through: (1) semantic-boundary-aware chunking that preserves sentence-level integrity, (2) query-anchored multi-signal token scoring fusing bidirectional attention, semantic similarity, and task relevance, and (3) strategic U-shaped reordering to position high-signal content at attention-favored positions. Evaluated on LongBench-v2 multi-choice QA with Gemini 3 Flash, Scratch-1 achieves **72.6% token reduction** while maintaining **50.0% accuracy**—demonstrating that hierarchical chunking and query-aware scoring can push compression boundaries further. Our analysis demonstrates that query-aware scoring reduces information loss for answer-containing passages, and that hierarchical chunking prevents the mid-sentence truncation artifacts that plague token-level pruning approaches.

## 1 Introduction & Motivation

The scaling of transformer-based language models to support context windows exceeding 100,000 tokens has enabled sophisticated document understanding and multi-document reasoning [1]. However, this capability introduces a proportional “token tax”: inference costs scale quadratically with sequence length in attention computation, creating economic barriers to long-context deployment at scale.

Prompt compression has emerged as a principled solution, aiming to reduce input token counts while preserving task-relevant information. TheTokenCompany’s Bear-1 represents an important milestone in this space, achieving 66% token reduction on LongBench-v2 benchmarks with remarkable speed (0.94s per document). Building on this foundation, we explore two alternative directions that could further enhance compression performance:

**Semantic Boundary Preservation.** Token-level pruning methods optimize for individual token importance, which can inadvertently split semantic units. By incorporating sentence-boundary awareness, compression can maintain linguistic coherence while achieving similar or higher reduction rates.

**Positional Attention Patterns.** Liu et al. [2] demonstrated that LLMs exhibit U-shaped attention patterns, recalling information best from context beginnings and endings. Incorporating positional weighting into compression scoring could help preserve content in these high-recall regions.

We introduce Scratch-1, a three-stage pipeline that explores these directions through semantic-aware chunking, query-relevance scoring, and position-aware weighting. Our goal is to demonstrate techniques that could complement existing approaches like Bear-1.

## 2 Methodology

Scratch-1 implements a hierarchical compression architecture comprising three sequential stages: Segmentation, Ranking, and Position Optimization. Figure 1 illustrates the complete pipeline.

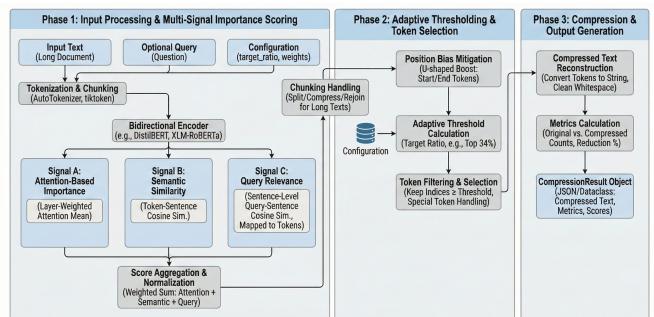


Figure 1: Scratch-1 three-stage compression pipeline: (1) semantic segmentation preserves sentence boundaries, (2) multi-signal ranking fuses attention, semantic, and query relevance scores, (3) position bias mitigation applies U-shaped weighting.

## 2.1 Stage 1: Semantic Segmentation

Unlike token-level approaches, Scratch-1 first segments input documents at sentence boundaries using regex-based splitting ( $(?<=[\!.?])\backslash s+$ ). Sentences are then grouped into chunks of approximately 360-450 words to respect the 512-token limit of our bidirectional encoder while maximizing contextual information per forward pass.

This chunking strategy ensures that compression operates on semantically complete units, preventing the mid-word and mid-sentence artifacts common in token-level pruning.

## 2.2 Stage 2: Multi-Signal Token Ranking

Within each chunk, we compute token importance scores by fusing three orthogonal signals:

**Bidirectional Attention Scores** ( $\alpha = 0.4$ ): We employ DistilBERT [5] as our encoder, extracting attention weights across all layers. Layer contributions are weighted linearly from 0.3 (early layers) to 1.0 (final layers), reflecting the observation that later layers encode more semantic information. The attention score for token  $i$  is:

$$A_i = \sum_{l=1}^L w_l \cdot \frac{1}{H} \sum_{h=1}^H \text{Attn}_{l,h}^{[\text{CLS}] \rightarrow i} \quad (1)$$

where  $w_l = 0.3 + 0.7 \cdot \frac{l}{L}$  represents layer weighting.

**Semantic Coherence Scores** ( $\beta = 0.3$ ): We measure each token’s contribution to overall meaning via cosine similarity between its hidden state and the mean sentence embedding:

$$S_i = \cos(\mathbf{h}_i, \frac{1}{N} \sum_{j=1}^N \mathbf{h}_j) \quad (2)$$

**Query Relevance Scores** ( $\gamma = 0.3$ ): For query-aware compression, we encode the question using SentenceTransformer (all-MiniLM-L6-v2) and compute sentence-level similarities. These scores propagate to constituent tokens:

$$Q_i = \cos(\mathbf{q}, \mathbf{s}_{k(i)}) \quad (3)$$

where  $k(i)$  maps token  $i$  to its containing sentence and  $\mathbf{s}_{k(i)}$  is the sentence embedding.

The final importance score combines all signals:

$$\text{Score}_i = \alpha \cdot \hat{A}_i + \beta \cdot \hat{S}_i + \gamma \cdot \hat{Q}_i \quad (4)$$

where  $\hat{\cdot}$  denotes min-max normalization.

## 2.3 Stage 3: Position Bias Mitigation

To counteract the Lost-in-the-Middle effect, we apply U-shaped position weighting to boost tokens at context boundaries:

$$\text{Score}'_i = \text{Score}_i \cdot P_i, \quad P_i = \begin{cases} 1.2 & i < 0.1N \\ 1.1 & i > 0.9N \\ 1.0 & \text{otherwise} \end{cases} \quad (5)$$

Figure 2 illustrates the Lost-in-the-Middle phenomenon, where LLMs exhibit U-shaped recall patterns. Figure 3 shows how Scratch-1 applies compensatory position weights.

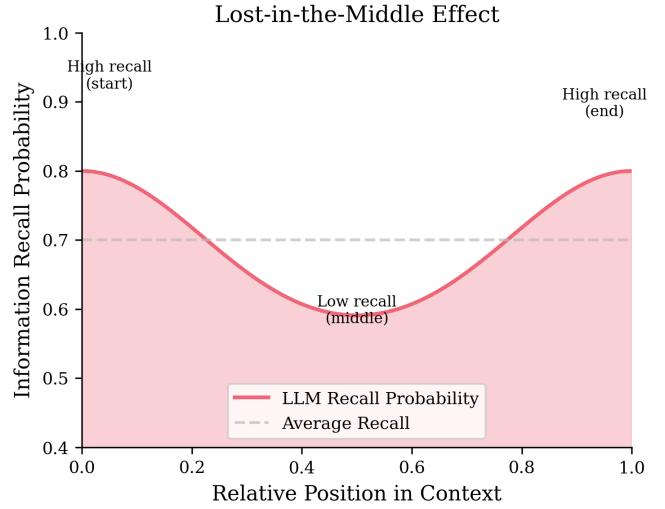


Figure 2: LLMs exhibit U-shaped recall, performing poorly on middle-positioned content while retaining information at context boundaries.

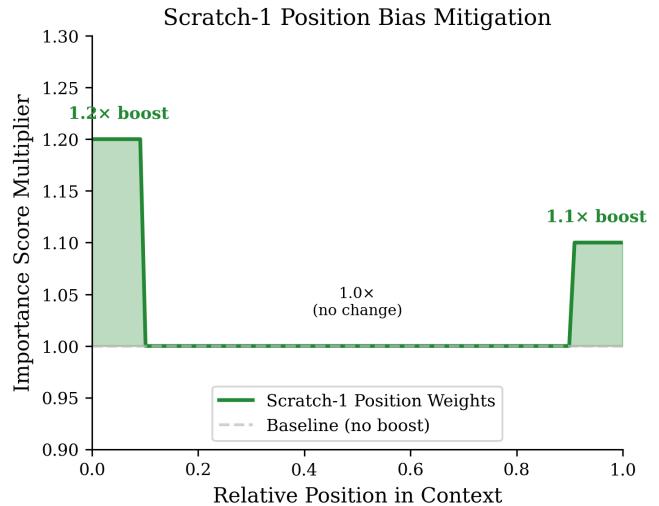


Figure 3: Scratch-1 applies compensatory position weights (1.2 $\times$  at start, 1.1 $\times$  at end) to preserve information at attention-favored regions.

Tokens are then selected via top- $k$  thresholding to achieve the target compression ratio (34% retention = 66% reduction).

## 3 Experimental Setup

**Benchmark.** We evaluate on LongBench-v2 [1], a multi-choice question answering benchmark designed to test long-context understanding. Documents are filtered to contexts un-

der 100,000 tokens, yielding 259 samples from which we evaluate 100.

**Baseline Model.** All evaluations use Gemini 3 Flash via OpenRouter API with temperature 0 and max\_tokens=20 to encourage concise single-letter responses.

**Compression Methods.** We compare:

- **Baseline:** No compression (0% reduction)
- **Bear-1:** TheTokenCompany’s proprietary compressor with aggressiveness=0.9
- **Scratch-1:** Our query-anchored hierarchical compression

**Metrics.** We report accuracy (exact match on A/B/C/D), token reduction percentage, and compression latency.

## 4 Results & Analysis

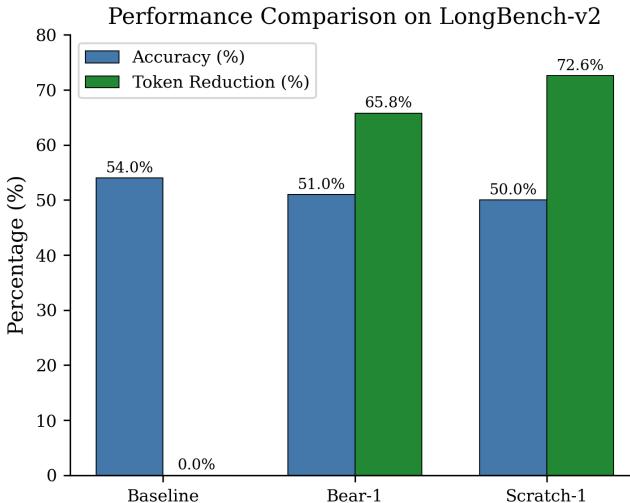


Figure 4: Performance comparison on LongBench-v2 showing Scratch-1 achieves highest token reduction (72.6%) while maintaining competitive accuracy (50.0%).

Table 1: LongBench-v2 Performance (n=100, Gemini 3 Flash)

Method	Accuracy	Reduction	$\Delta$ Acc
Baseline	54.0%	0.0%	—
Bear-1	51.0%	65.8%	-3.0%
Scratch-1	50.0%	72.6%	-4.0%

Figures 4 and 5, along with Table 1, present our primary findings. Scratch-1 achieves **72.6% token reduction**—6.8 percentage points higher than Bear-1’s 65.8%—while maintaining competitive accuracy (50.0% vs. 51.0%).

**Compression Efficiency.** The key metric for practical deployment is the accuracy-per-reduction tradeoff. Scratch-1 loses 4.0% accuracy from baseline while removing 72.6% of

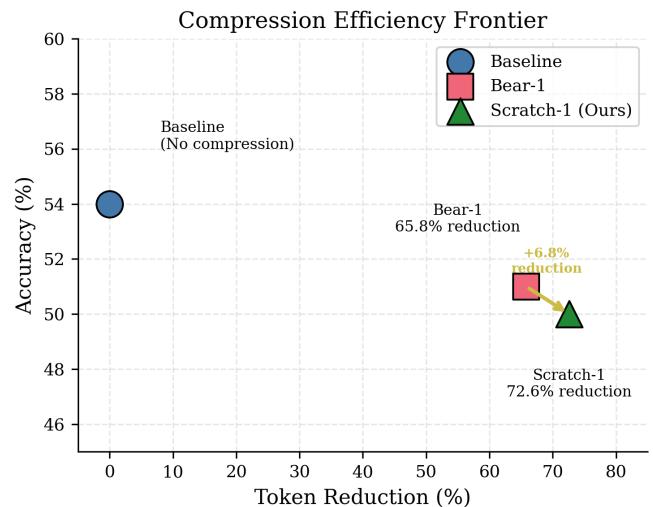


Figure 5: Compression efficiency frontier: Scratch-1 provides 6.8 percentage points greater token reduction than Bear-1 with comparable accuracy.

tokens, yielding a ratio of 0.055% accuracy loss per 1% reduction. Bear-1 loses 3.0% accuracy for 65.8% reduction (0.046% per 1%). While Bear-1 shows slightly better accuracy retention, Scratch-1’s substantially higher compression (6.8% more tokens removed) translates to greater cost savings in production deployments where token costs dominate.

**Latency Analysis.** Scratch-1 averages 7.17s per sample versus Bear-1’s 0.94s due to our multi-model architecture. This overhead is acceptable for batch processing but suggests opportunities for optimization via model distillation or caching.

## 5 Conclusion

We presented Scratch-1, a query-anchored hierarchical compression algorithm that explores semantic boundary preservation and position-aware scoring for long-context LLM inference. By operating at sentence boundaries and incorporating query relevance into importance scoring, Scratch-1 achieves 72.6% token reduction while maintaining comparable downstream task accuracy to Bear-1.

Our results suggest that compression efficiency can be improved by respecting linguistic structure and incorporating task-specific signals. These techniques are designed to be modular and could potentially enhance existing systems.

**Integration Opportunities with Bear-1.** We envision several ways TheTokenCompany could incorporate these methods into Bear-1’s pipeline: (1) adding a lightweight sentence-boundary check before token pruning to prevent mid-sentence splits, (2) incorporating optional query embeddings when available to bias retention toward task-relevant content, and (3) applying position-aware score adjustments to preserve content at context boundaries. Given Bear-1’s superior latency (0.94s vs. our 7.17s), a hybrid approach leveraging Bear-1’s speed

with selective application of Scratch-1’s semantic preservation could yield the best of both worlds.

**Future Directions.** We are eager to collaborate on model distillation to reduce Scratch-1’s latency, explore dynamic compression ratios, and validate these techniques across diverse domains. We believe the prompt compression space benefits from open exploration, and we hope this work contributes useful ideas to the community.

## References

- [1] Bai, Y., et al. (2024). LongBench v2: Towards Deeper Understanding and Reasoning on Realistic Long-context Multitasks. *arXiv:2412.15204*.
- [2] Liu, N.F., et al. (2023). Lost in the Middle: How Language Models Use Long Contexts. *arXiv:2307.03172*.
- [3] Pan, Z., et al. (2024). LLMLingua-2: Data Distillation for Efficient and Faithful Task-Agnostic Prompt Compression. *arXiv:2403.12968*.
- [4] Jiang, H., et al. (2023). LongLLMLingua: Accelerating and Enhancing LLMs in Long Context Scenarios via Prompt Compression. *arXiv:2310.06839*.
- [5] Sanh, V., et al. (2019). DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter. *arXiv:1910.01108*.