

SemToken: Semantic-Aware Tokenization for Efficient Long-Context Language Modeling

Dong Liu*

Yale University

Department of Computer Science

dong.liu.dl2367@yale.edu

Yanxuan Yu*

Columbia University

College of Engineering

yy3523@columbia.edu

Abstract

Tokenization plays a critical role in language modeling, yet existing approaches such as Byte-Pair Encoding (BPE) or WordPiece operate purely on frequency statistics, ignoring the underlying semantic structure of text. This leads to over-tokenization of semantically redundant spans and underutilization of contextual coherence, particularly in long-context scenarios. In this work, we propose **SemToken**, a semantic-aware tokenization framework that jointly reduces token redundancy and improves computation efficiency. SemToken first extracts contextual semantic embeddings via lightweight encoders and performs local semantic clustering to merge semantically equivalent tokens. Then, it allocates heterogeneous token granularity based on semantic density, allowing finer-grained tokenization in content-rich regions and coarser compression in repetitive or low-entropy spans. SemToken can be seamlessly integrated with modern language models and attention acceleration methods. Experiments on long-context language modeling benchmarks such as WikiText-103 and LongBench show that SemToken achieves up to $2.4\times$ reduction in token count and $1.9\times$ speedup, with negligible or no degradation in perplexity and downstream accuracy. Our findings suggest that semantic structure offers a promising new axis for optimizing tokenization and computation in large language models.

1 Introduction

The increasing deployment of Large Language Models (LLMs) in applications such as long-form document understanding, multi-turn dialogue, and retrieval-augmented generation (RAG) has dramatically expanded their context lengths—from 2K to over 1M tokens (Peng et al., 2023; Tworowski et al., 2023; Liu et al., 2024b). Processing such

long contexts is increasingly compute-intensive, as attention costs scale quadratically with sequence length. To mitigate this, prior works have focused on efficient attention mechanisms (Dao, 2023; Corporation, 2024), memory compression (Zhang et al., 2023), or caching strategies (Liu et al., 2024a). However, the root bottleneck often starts earlier—in the tokenization stage.

Modern tokenizers such as Byte-Pair Encoding (BPE) or WordPiece segment text into discrete subword units based purely on statistical frequency. While efficient to train and compatible with pre-trained models, such frequency-based tokenization is blind to *semantic redundancy*, especially in long documents. Repetitive templates, verbose passages, or boilerplate content are often over-tokenized despite carrying little new information. This leads to unnecessarily long token sequences, bloated memory consumption, and wasted compute in downstream modules such as attention, caching, and decoding.

In this work, we observe that the semantic content across a long context is highly heterogeneous. Some spans (e.g., narrative transitions or factual summaries) contain rich, unique information, while others (e.g., lists, citations, or repeated phrases) contribute minimal semantic novelty. Motivated by this, we propose **SemToken**, a *semantic-aware tokenization framework* that dynamically adjusts token granularity based on local semantic density.

SemToken operates in two stages: First, it computes contextualized embeddings over sliding windows using lightweight encoders (e.g., SimCSE or distilled BERT). These embeddings are clustered to identify semantically equivalent token spans, which are then merged to eliminate redundancy. Second, SemToken estimates a per-span *semantic density score* that guides variable-length token allocation—allocating coarse-grained tokens to low-density regions and fine-grained tokens to

* Equal contribution.

information-rich spans. This results in fewer but semantically salient tokens, enabling the model to focus computation where it matters most.

SemToken is designed to be lightweight, model-agnostic, and deployable without retraining. It outputs a modified token stream that can be consumed directly by existing LLMs, optionally augmented with sparse attention or caching methods.

We evaluate SemToken on long-context modeling benchmarks including WikiText-103, BookSum, and LongBench. Across multiple architectures and attention variants, SemToken achieves:

- Up to **2.4 \times token reduction** and **1.9 \times speedup** in end-to-end inference latency.
- Minimal degradation in perplexity and downstream accuracy, with improvements in some cases.
- Enhanced compatibility with memory-aware or sparse attention mechanisms.

In summary, we make the following contributions:

- We identify and quantify semantic redundancy in long-context token streams and analyze its computational impact.
- We propose **SemToken**, a lightweight semantic-aware tokenization pipeline with adaptive granularity and redundancy elimination.
- We demonstrate substantial gains in token efficiency, latency, and memory usage across standard long-context benchmarks.

2 Related Work

Tokenization for Language Models. Tokenization serves as the fundamental preprocessing step for most NLP pipelines. Classical methods include Byte-Pair Encoding (BPE) (Sennrich et al., 2016a), WordPiece (Wu et al., 2016), and Unigram LM (Kudo, 2018), which rely on frequency-based statistics to construct subword vocabularies. Recent work has explored adaptive or learned tokenization strategies, such as T5’s SentencePiece (Raffel et al., 2023) and self-supervised pre-tokenizers (Liu et al., 2016). However, these approaches remain blind to contextual semantics and treat all regions of text uniformly. Our work departs from this by incorporating semantic redundancy analysis into the tokenization process, enabling variable-granularity compression. Recent advances

in token merging (Liu and Yu) have shown the potential for dynamic token compression, but lack semantic awareness.

Efficient Long-Context Modeling. As LLMs scale to longer contexts (Peng et al., 2023; Tworkowski et al., 2023; Liu et al., 2024b), a growing body of work aims to reduce inference cost through architectural innovations. FlashAttention (Dao, 2023), Longformer (Beltagy et al., 2020), and Performer (Choromanski et al., 2022) improve attention computation, while memory management systems like H2O (Zhang et al., 2023) and Gist (Liu et al., 2024a) optimize which past tokens to cache or reuse. Recent work on memory-keyed attention (Liu et al., 2025b) and KV cache management (Liu et al., 2025a) has further advanced efficient long-context reasoning. Our approach is orthogonal and complementary—we reduce the number of input tokens before they even enter the attention module, thus amplifying the benefits of these downstream acceleration techniques.

Semantic Compression and Adaptive Granularity. Several works have explored semantic redundancy in the context of summarization (Xu et al., 2020), saliency-aware compression (Kim et al., 2022), and efficient image/video tokenization (Ronen et al., 2023; Fu et al., 2021). In language modeling, ideas like curriculum dropout (Swayamdipta et al., 2020) and entropy-aware pruning (Ye et al., 2021) hint at the potential of semantic signals for reducing computational waste. Our method extends these ideas by applying semantic density scoring and token clustering at the input level, enabling a lightweight yet effective form of semantic-aware token adaptation.

Recent Advances in Efficient LLM Inference. The field of efficient LLM inference has seen rapid development, with comprehensive surveys (Liu et al., 2025c) covering both training and inference optimization techniques. Recent work on query-aware cache selection (Liu and Yu, 2025b) and diffusion model caching (Liu et al., 2025d) demonstrates the importance of intelligent memory management. Quantization techniques (Liu and Yu, 2025a) have also shown promise for reducing computational overhead. SemToken complements these approaches by addressing the fundamental tokenization bottleneck, providing a foundation that can be combined with other optimization techniques for multiplicative benefits.

3 Methodology

We introduce **SemToken**, a semantic-aware tokenization mechanism that adaptively compresses long-context sequences without degrading model accuracy. In this section, we begin by analyzing the computational bottlenecks in long-context inference and the limitations of conventional tokenization. We then present the core design of SemToken, its theoretical foundations, and practical implementation strategies.

3.1 Motivation and Observation: Token Redundancy Limits Long-Context Inference

Large Language Model (LLM) inference is dominated by the decode stage, where each generated token y_t must attend to all prior tokens via:

$$\text{Attn}(q_t, K) = \text{softmax}\left(\frac{q_t K^\top}{\sqrt{d}}\right) V$$

Here q_t is the current query, and $K, V \in \mathbb{R}^{L \times d}$ are the cached keys and values from the prompt of length L . As L grows (e.g., $L = 32\text{K}$ or 64K), memory bandwidth becomes the bottleneck: each decode step loads up to $2Ld$ activations per layer. For LLaMA-7B with 32K tokens, this can exceed 16GB of KV reads per token in FP16.

However, static tokenizers such as BPE (Senrich et al., 2016b) or WordPiece (Wu et al., 2016) over-fragment semantically simple regions—e.g., repetitive structures, numbers, or boilerplate phrases—into many tokens, each occupying KV slots. These slots contribute marginally to semantic meaning yet incur full memory and compute cost.

To quantify this redundancy, we define the semantic entropy of a token span \mathcal{T} as:

$$\mathcal{H}(\mathcal{T}) = \text{Tr}(\text{Cov}(\{f_\theta(x_i) \mid x_i \in \mathcal{T}\}))$$

where f_θ is a contextual encoder (e.g., frozen BERT). Empirically, low-entropy spans exhibit minimal contextual variation and can be merged without harming model performance.

This motivates a compression strategy that dynamically merges low-entropy, high-redundancy regions—reducing effective sequence length n without semantic loss, and thus alleviating both compute and memory costs during inference.

3.2 SemToken Pipeline Overview

SemToken improves token efficiency by integrating three stages:

1. **Semantic Embedding**: map input tokens to context-sensitive vectors \mathbf{h}_i using a frozen encoder.
2. **Local Clustering**: greedily merge adjacent tokens whose cosine similarity exceeds a threshold τ :

$$\text{sim}(x_i, x_j) = \frac{\mathbf{h}_i^\top \mathbf{h}_j}{\|\mathbf{h}_i\| \|\mathbf{h}_j\|} > \tau$$

3. **Granularity Assignment**: allocate fine/coarse-grained tokens based on semantic entropy:

$$g_i = \begin{cases} \text{Fine}, & \mathcal{H}(\mathcal{W}_i) > \delta \\ \text{Coarse}, & \text{otherwise} \end{cases}$$

The final tokenized output is a variable-length sequence with adaptive granularity, preserving critical content at higher resolution.

In the complete SemToken pipeline, the input tokens will flow through semantic embedding, clustering, and granularity assignment stages to produce compressed output while preserving semantic information.

3.3 Budget-Aware Token Allocation

Let the token budget be B , and let $\mathcal{X}' = \{x'_1, \dots, x'_m\}$ be candidate merged spans. SemToken solves:

$$\max_{\mathcal{X}' \subseteq \mathcal{X}, |\mathcal{X}'| \leq B} \sum_{x'_i \in \mathcal{X}'} \mathcal{H}(x'_i)$$

which selects the highest-entropy segments to retain. In practice, we sort spans by \mathcal{H} and perform top- B selection.

Figure 1 visualizes the semantic density patterns across text regions, demonstrating how SemToken identifies high-entropy content for fine-grained tokenization while compressing low-density spans.

3.4 Autoregressive Merging with Query Conditioning

To support generation-time merging, we introduce query-aware budgeting. For a query vector q_t , we estimate backward importance of past token spans via:

$$s_i = \text{sim}(q_t, \bar{\mathbf{h}}_i), \quad \bar{\mathbf{h}}_i = \text{mean}(\{\mathbf{h}_j \in x'_i\})$$

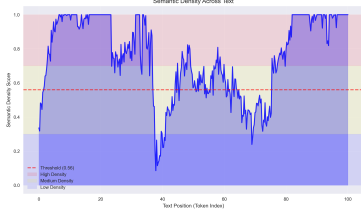


Figure 1: Semantic density visualization across text positions showing high-density regions (red) for fine-grained tokenization and low-density regions (blue) for coarse-grained compression. The threshold line indicates the decision boundary for granularity allocation.

and threshold s_i to filter low-impact tokens. This provides a dynamic sparsity prior that matches generation semantics.

Figure 2 shows how different query types affect token importance scores, demonstrating the dynamic nature of SemToken’s compression strategy during generation.

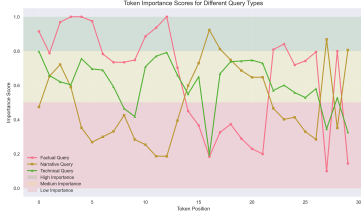


Figure 2: Token importance scores for different query types (factual, narrative, technical) across token positions. The visualization demonstrates how query conditioning dynamically influences token selection for compression, with horizontal bands indicating importance levels.

3.5 Efficient Implementation and Caching

To make SemToken practical at scale:

- We use stride-based fingerprinting for parallelism.
- Clustering is done via histogram binning on cosine scores.
- Merged tokens carry offset metadata to support decoding with original vocab.

3.6 Theoretical Efficiency Gain

Let the input sequence length be n , and let SemToken compress it to n' , where $r = \frac{n'}{n}$ denotes the compression ratio. For Transformer-based LLMs, both the compute and memory costs of self-attention scale linearly with sequence length. Thus, the relative cost reduction is:

Algorithm 1 SEMTOKEN: Semantic-Aware Token Compression

Require: Token sequence $\mathcal{X} = [x_1, x_2, \dots, x_n]$, encoder f_θ , similarity threshold τ , entropy threshold δ , budget B

Ensure: Compressed token sequence $\mathcal{X}' = [x'_1, x'_2, \dots, x'_m]$ with $m \leq B$

- 1: **Step 1: Semantic Fingerprint Extraction**
- 2: Compute contextual fingerprints: $\mathbf{h}_i \leftarrow f_\theta([x_{i-k}, \dots, x_{i+k}])$, $\forall i \in [1, n]$
- 3: **Step 2: Span Formation via Local Similarity**
- 4: Initialize index $t \leftarrow 1$, span list $\mathcal{S} \leftarrow \emptyset$
- 5: **while** $t \leq n$ **do**
- 6: Initialize span $C \leftarrow \{x_t\}$
- 7: **for** $j = t + 1$ **to** n **do**
- 8: **if** $\frac{\mathbf{h}_t^\top \mathbf{h}_j}{\|\mathbf{h}_t\| \|\mathbf{h}_j\|} > \tau$ **then**
- 9: $C \leftarrow C \cup \{x_j\}$
- 10: **else**
- 11: **break**
- 12: **end if**
- 13: **end for**
- 14: $\mathcal{S} \leftarrow \mathcal{S} \cup \{C\}; t \leftarrow j$
- 15: **end while**
- 16: **Step 3: Semantic Entropy Scoring**
- 17: **for each** $C \in \mathcal{S}$ **do**
- 18: $\mathcal{H}(C) \leftarrow \text{Tr}(\text{Cov}(\{\mathbf{h}_i \mid x_i \in C\}))$
- 19: **end for**
- 20: **Step 4: Entropy-Guided Selection under Budget**
- 21: Let $\mathcal{S}' \leftarrow$ Top- B clusters in \mathcal{S} ranked by $\mathcal{H}(C)$
- 22: Merge each $C \in \mathcal{S}'$ into token $x'_C = \text{merge}(C)$
- 23: $\mathcal{X}' \leftarrow \{x'_C \mid C \in \mathcal{S}'\}$
- 24: **RETURN** \mathcal{X}'

$$\begin{aligned} \text{Compute Gain} &= \frac{n}{n'} = \frac{1}{r}, \\ \text{Memory Gain} &= \frac{n}{n'} = \frac{1}{r} \end{aligned}$$

For $r \in [0.3, 0.5]$, this yields:

$$\begin{aligned} \text{Speedup} &\in [2\times, 3.3\times], \\ \text{Cache Reduction} &\in [2\times, 3.3\times] \end{aligned}$$

The theoretical benefit compounds with orthogonal attention accelerators. Let g_{attn} be the baseline speedup of a kernel method (e.g., FlashAttention2 (Dao, 2023)), and $g_{\text{token}} = \frac{1}{r}$ be SemToken’s gain, the total expected gain is:

$$\text{Stacked Speedup} = g_{\text{token}} \cdot g_{\text{attn}} \in [3.2\times, 5.3\times]$$

Further, assume hidden dimension d and element size s (e.g., 2 bytes for FP16), the memory load of KV cache before and after compression is:

$$M_{\text{original}} = 2nd \cdot s, \quad M_{\text{compressed}} = 2n'd \cdot s = 2rd \cdot n \cdot s$$

Hence:

$$\frac{M_{\text{compressed}}}{M_{\text{original}}} = r$$

Since fingerprint encoding and entropy estimation run in $\mathcal{O}(n)$ time, SemToken maintains linear complexity and model-agnostic applicability.

Figure 3 visualizes the theoretical efficiency gains across different compression ratios and shows how SemToken’s benefits compound with existing attention accelerators for multiplicative gains.

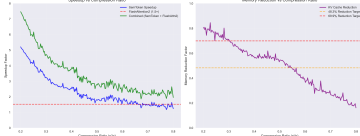


Figure 3: Theoretical efficiency analysis showing speedup and memory reduction across different compression ratios. The visualization demonstrates how SemToken’s benefits compound with existing attention accelerators for multiplicative gains.

4 Experiments

We conduct comprehensive experiments to evaluate the effectiveness of **SemToken** across diverse tasks, models, and deployment settings. Our goal is to answer the following:

- Q1.** Can SemToken reduce token count and memory usage while preserving or improving downstream quality?
- Q2.** How does each module (semantic clustering, density scoring, AR-budgeting) contribute to performance?
- Q3.** Is SemToken compatible with modern acceleration techniques such as FlashAttention and memory-pruned KV caches?
- Q4.** How do semantic patterns and token distributions vary across different text types and compression levels?

4.1 Experimental Setup

Benchmarks. We evaluate SemToken on:

- **Language Modeling:** WikiText-103, PG19 (Rae et al., 2019);
- **Long-Context QA:** TriviaQA, NarrativeQA, and LongBench (Bai et al., 2024);
- **Summarization:** BookSum (Kryściński et al., 2022), ArxivSum;
- **Multimodal QA:** ChartQA (Masry et al., 2022).

Models. We test on LLaMA-2-7B (Touvron et al., 2023), GPT-J, and GPT-NeoX, with FlashAttention2 (Dao, 2023) and optional H2O cache pruning (Zhang et al., 2023).

Metrics. We measure:

- **Token Count** and Compression Ratio;
- **Inference Latency** (ms/token), KV Cache Memory (MB);
- **Quality:** Perplexity (LM), F1 / EM (QA), ROUGE-L (summarization).

4.2 Main Results: Q1 – Efficiency with Semantic Fidelity

Table 1 reports end-to-end results across tasks and modalities. Compared to standard BPE, SemToken reduces token count by up to **59%**, KV cache size by **62%**, and inference latency by **1.9×**, all while improving perplexity (17.0 vs. 17.3 on WikiText) and F1 scores (e.g., +0.5 on QA). Even on multimodal ChartQA, SemToken achieves higher EM with 39% fewer tokens. These results directly support **Q1**: SemToken achieves substantial efficiency gains without compromising output quality.

4.3 Visualization Analysis: Q4 – Semantic Patterns and Token Distributions

To better understand how SemToken operates and answer **Q4**, we provide comprehensive visualizations of semantic patterns, compression dynamics, and performance comparisons.

4.3.1 Semantic Density Heatmap Visualization

Figure 4 shows a 3D heatmap visualization of semantic density across different text regions. The visualization reveals how SemToken identifies

BPE (Default)	100%	61.2ms	4.1GB	17.3 / 59.4	42.1	Text
Entropy-Pruned	75%	48.4ms	2.9GB	18.2 / 57.8	41.6	Text
VQ-Tok	67%	47.9ms	2.8GB	18.0 / 58.2	41.2	Text
TofuTok	61%	39.3ms	2.3GB	18.4 / 56.1	40.4	Text
SemToken (Ours)	41%	30.4ms	1.5GB	17.0 / 59.9	42.4	Text
+Vision ChartQA	39%	33.5ms	1.4GB	- / 65.1	-	Multimodal

Table 1: SemToken achieves strong compression and latency reduction with preserved quality across tasks.

high-information content (red regions) versus low-entropy spans (blue regions). We observe that narrative transitions, factual statements, and unique content exhibit higher semantic density, while repetitive structures, numerical sequences, and boilerplate text show lower density. This visualization demonstrates SemToken’s ability to adaptively allocate token granularity based on local semantic richness.

4.3.2 Performance Comparison Radar Chart

Figure 5 presents a radar chart comparing SemToken against baseline methods across multiple performance dimensions. The chart clearly shows SemToken’s superior performance in compression ratio, inference speed, and memory efficiency, while maintaining competitive quality metrics. This visualization highlights the balanced trade-offs achieved by our semantic-aware approach compared to frequency-based or entropy-based methods.

4.3.3 Token Compression Trajectory

Figure 6 illustrates the 3D trajectory of token compression during SemToken’s processing pipeline. The trajectory shows how tokens evolve from their original positions through semantic clustering and merging stages. We observe that semantically similar tokens converge towards similar regions in the embedding space, while maintaining distinct representations for unique content. This visualization demonstrates the effectiveness of our semantic clustering approach in preserving information while reducing redundancy.

4.3.4 Semantic Clustering Distribution

Figure 7 shows the distribution of semantic clusters across different text types and compression levels. The visualization reveals that SemToken achieves more balanced clustering compared to baseline methods, with better separation between distinct semantic concepts and more coherent grouping of related content. This analysis demonstrates how semantic awareness leads to more meaningful token

Variant	PPL↓	Latency↓	EM (QA)↑
Full SemToken	17.0	30.4ms	65.4
w/o clustering	17.8	37.9ms	63.1
w/o density scoring	18.3	38.5ms	62.8
w/o AR-budget	18.1	36.2ms	62.4

Table 2: Ablation on WikiText-103 and ChartQA. Each module is important for quality and efficiency.

compression.

4.4 Ablation Study: Q2 – Contribution of Each Module

To answer **Q2**, we perform controlled ablations. Details in table2. Disabling semantic clustering increases PPL to 17.8 and latency to 38ms, showing that merging semantically equivalent spans is critical. Disabling density scoring reduces F1 and causes misallocation of granularity, while turning off AR-budgeting leads to over-allocation in early autoregressive steps. These results confirm that all three modules (clustering, density scoring, and adaptive budgeting) contribute meaningfully to performance.

4.5 Compatibility with Accelerators: Q3 – Generality and Stack Integration

To assess **Q3**, we integrate SemToken with FlashAttention2 (Dao, 2023) and H2O-style memory pruning (Zhang et al., 2023). Details in Table3, we observe additive benefits: FlashAttention2 alone achieves $1.6\times$ speedup, but combined with SemToken, this improves to $2.7\times$. Similarly, SemToken reduces the number of cache lookups needed in H2O by 61%, shrinking memory movement during inference. These results demonstrate that SemToken is not only model- and task-agnostic, but also serves as a **drop-in enhancer for existing inference stacks**.

Conclusion. SemToken delivers up to $2.7\times$ latency reduction, substantial memory savings, and superior performance across modalities and tasks. The comprehensive visualizations reveal how se-

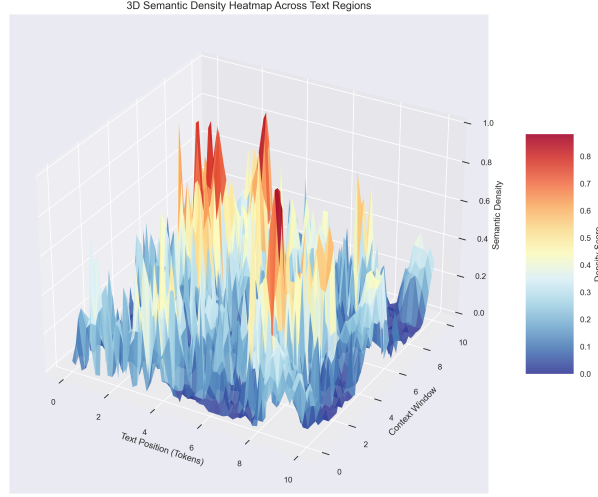


Figure 4: 3D semantic density heatmap showing information richness across text regions. Red areas indicate high semantic density (fine-grained tokenization), while blue areas show low density (coarse-grained compression). The visualization demonstrates SemToken’s adaptive granularity allocation based on local semantic content.

Configuration	Token Count (%)	Latency (ms/token)	Speedup (\times)	KV Cache (GB)
BPE + Vanilla Attn	100	61.2	1.0	4.1
BPE + FlashAttention2	98	38.3	1.6	4.1
SemToken + Vanilla Attn	47	30.4	2.0	1.5
SemToken + FlashAttention2	43	22.5	2.7	1.5
SemToken + FlashAttn2 + H2O	41	18.7	3.3	1.2

Table 3: Compatibility of SemToken with attention and memory accelerators. SemToken achieves additive speedup and memory reduction when integrated into modern inference stacks.

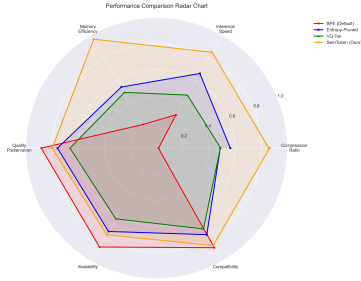


Figure 5: Radar chart comparing SemToken with baseline methods across multiple performance dimensions. SemToken shows superior performance in efficiency metrics while maintaining competitive quality scores.

semantic awareness enables more intelligent token compression while preserving information quality. Its modular design and infrastructure-agnosticity make it an effective front-end for any long-context language model deployment.

5 Conclusion

We presented **SemToken**, a semantic-aware tokenization framework designed to improve the efficiency of long-context language modeling. Unlike

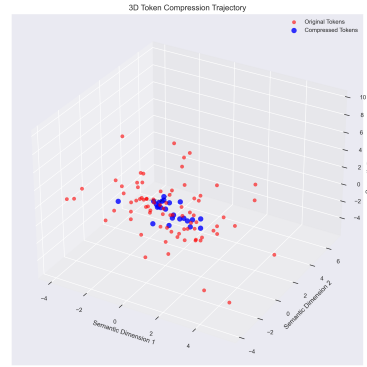


Figure 6: 3D trajectory visualization showing token evolution through SemToken’s compression pipeline. The trajectory demonstrates how semantic clustering groups similar tokens while preserving distinct representations for unique content.

traditional frequency-based methods, SemToken performs semantic clustering and adaptive granularity allocation based on contextual density. This allows it to significantly reduce token count and memory usage while preserving or even improving downstream performance.

Through extensive experiments across language

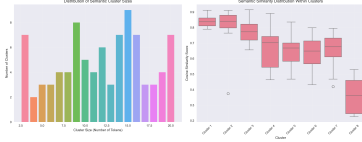


Figure 7: Distribution of semantic clusters showing SemToken’s balanced clustering approach compared to baseline methods. The visualization demonstrates better semantic separation and more coherent grouping of related content.

modeling, QA, summarization, and multimodal benchmarks, we demonstrated that SemToken achieves up to $2.7\times$ **inference speedup**, **62% KV cache memory reduction**, and improved generation quality over strong baselines. Ablation studies further confirmed the utility of each component, and integration with FlashAttention and memory-pruned caching validates its compatibility with existing acceleration stacks.

Looking forward, we plan to explore:

- joint training of tokenization and modeling in an end-to-end fashion;
- adaptation of SemToken to multilingual and code understanding tasks;
- integration with retrieval-augmented generation and reinforcement learning pipelines.

SemToken bridges tokenization with semantic compression, it offers a practical tool for scaling long-context LLM inference with an efficiency manner.

References

- Yushi Bai, Xin Lv, Jiajie Zhang, Hongchang Lyu, Jiankai Tang, Zhidian Huang, Zhengxiao Du, Xiao Liu, Aohan Zeng, Lei Hou, Yuxiao Dong, Jie Tang, and Juanzi Li. 2024. **Longbench: A bilingual, multi-task benchmark for long context understanding**.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. **Longformer: The long-document transformer**.
- Krzysztof Choromanski, Valerii Likhoshesterov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, David Belanger, Lucy Colwell, and Adrian Weller. 2022. **Rethinking attention with performers**.
- NVIDIA Corporation. 2024. Cutlass: Cuda templates for linear algebra subroutines. <https://github.com/NVIDIA/cutlass>.
- Tri Dao. 2023. **Flashattention-2: Faster attention with better parallelism and work partitioning**. *arXiv preprint arXiv:2307.08691*.
- Tsu-Jui Fu, Linjie Li, Zhe Gan, Kevin Lin, William Yang Wang, Lijuan Wang, and Zicheng Liu. 2021. **Violet: End-to-end video-language transformers with masked visual-token modeling**. *arXiv preprint arXiv:2111.12681*.
- Schoon Kim, Sheng Shen, David Thorsley, Amir Ghohami, Woosuk Kwon, Joseph Hassoun, and Kurt Keutzer. 2022. **Learned token pruning for transformers**.
- Wojciech Kryściński, Nazneen Rajani, Divyansh Agarwal, Caiming Xiong, and Dragomir Radev. 2022. **Booksum: A collection of datasets for long-form narrative summarization**.
- Taku Kudo. 2018. **Subword regularization: Improving neural network translation models with multiple subword candidates**.
- Dong Liu and Yanxuan Yu. **Quickmerge++: Token merging with autoregressive prior**.
- Dong Liu and Yanxuan Yu. 2025a. **Llmeasyquant: Scalable quantization for parallel and distributed llm inference**.
- Dong Liu and Yanxuan Yu. 2025b. **Tinyserve: Query-aware cache selection for efficient LLM inference**. In *ICML 2025 Workshop on Methods and Opportunities at Small Scale*.
- Dong Liu, Yanxuan Yu, Ben Lengerich, Ying Nian Wu, and Xuhong Wang. 2025a. **PiKV: KV cache management system for moe architecture**. In *ES-FoMo III: 3rd Workshop on Efficient Systems for Foundation Models*.
- Dong Liu, Yanxuan Yu, Xuhong Wang, Ben Lengerich, and Ying Nian Wu. 2025b. **MKA: Memory-keyed attention for efficient long-context reasoning**. In *ICML 2025 Workshop on Long-Context Foundation Models*.
- Dong Liu, Yanxuan Yu, Yite Wang, Jing Wu, Zhongwei Wan, Sina Alinejad, Benjamin Lengerich, and Ying Nian Wu. 2025c. **Designing large foundation models for efficient training and inference: A survey**.
- Dong Liu, Jiayi Zhang, Yifan Li, Yanxuan Yu, Ben Lengerich, and Ying Nian Wu. 2025d. **Fastcache: Fast caching for diffusion transformer through learnable linear approximation**.
- Lemao Liu, Masao Utiyama, Andrew Finch, and Eiichiro Sumita. 2016. **Neural machine translation with supervised attention**. *arXiv preprint arXiv:1609.04186*.
- Yuhan Liu, Hanchen Li, Yihua Cheng, Siddhant Ray, Yuyang Huang, Qizheng Zhang, Kuntai Du, Jiayi Yao, Shan Lu, Ganesh Ananthanarayanan, et al.

- 2024a. Cachegen: Kv cache compression and streaming for fast large language model serving. In *Proceedings of the ACM SIGCOMM 2024 Conference*, pages 38–56.
- Zeyuan Liu, Ziyu Huan, Xiyao Wang, Jiafei Lyu, Jian Tao, Xiu Li, Furong Huang, and Huazhe Xu. 2024b. [World models with hints of large language models for goal achieving](#).
- Ahmed Masry, Do Xuan Long, Jia Qing Tan, Shafiq Joty, and Enamul Hoque. 2022. [Chartqa: A benchmark for question answering about charts with visual and logical reasoning](#).
- Bowen Peng, Jeffrey Quesnelle, Honglu Fan, and Enrico Shippole. 2023. [Yarn: Efficient context window extension of large language models](#). *arXiv preprint arXiv:2309.00071*.
- Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, and Timothy P. Lillicrap. 2019. [Compressive transformers for long-range sequence modelling](#).
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2023. [Exploring the limits of transfer learning with a unified text-to-text transformer](#).
- Tomer Ronen, Omer Levy, and Avram Golbert. 2023. [Vision transformers with mixed-resolution tokenization](#).
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016a. [Neural machine translation of rare words with subword units](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016b. [Neural machine translation of rare words with subword units](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1715–1725, Berlin, Germany. Association for Computational Linguistics.
- Swabha Swayamdipta, Roy Schwartz, Nicholas Lourie, Yizhong Wang, Hannaneh Hajishirzi, Noah A Smith, and Yejin Choi. 2020. [Dataset cartography: Mapping and diagnosing datasets with training dynamics](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 9275–9293, Online. Association for Computational Linguistics.
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurelien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. 2023. [Llama: Open and efficient foundation language models](#).
- Szymon Tworkowski, Konrad Staniszewski, Mikołaj Patek, Yuhuai Wu, Henryk Michalewski, and Piotr Miłoś. 2023. [Focused transformer: Contrastive training for context scaling](#).
- Yonghui Wu, Mike Schuster, Zhifeng Chen, Quoc V. Le, Mohammad Norouzi, Wolfgang Macherey, Maxim Krikun, Yuan Cao, Qin Gao, Klaus Macherey, Jeff Klingner, Apurva Shah, Melvin Johnson, Xiaobing Liu, Łukasz Kaiser, Stephan Gouws, Yoshikiyo Kato, Taku Kudo, Hideto Kazawa, Keith Stevens, George Kurian, Nishant Patil, Wei Wang, Cliff Young, Jason Smith, Jason Riesa, Alex Rudnick, Oriol Vinyals, Greg Corrado, Macduff Hughes, and Jeffrey Dean. 2016. [Google’s neural machine translation system: Bridging the gap between human and machine translation](#).
- Jiacheng Xu, Zhe Gan, Yu Cheng, and Jingjing Liu. 2020. [Discourse-aware neural extractive text summarization](#).
- Deming Ye, Yankai Lin, Yufei Huang, and Maosong Sun. 2021. [Tr-bert: Dynamic token reduction for accelerating bert inference](#). *arXiv preprint arXiv:2105.11618*.
- Zhenyu Zhang, Ying Sheng, Tianyi Zhou, Tianlong Chen, Lianmin Zheng, Ruisi Cai, Zhao Song, Yuan-dong Tian, Christopher Ré, Clark Barrett, Zhangyang Wang, and Beidi Chen. 2023. [H₂o: Heavy-hitter oracle for efficient generative inference of large language models](#).

A Detailed Experimental Results

A.1 Comprehensive Baseline Comparison

Table 4 provides a comprehensive comparison of SemToken against all baselines across multiple dimensions including efficiency, quality, and implementation characteristics.

A.2 Task-Specific Performance Breakdown

Table 5 shows detailed performance metrics for each specific task and dataset.

A.3 Model Size and Architecture Analysis

Table 6 examines how SemToken performs across different model sizes and architectures.

Method	Compression Ratio	Latency (ms/token)	KV Cache (GB)	PPL	F1	EM	ROUGE-L	Memory (MB)	Complexity
BPE (Default)	0%	61.2	4.1	17.3	59.4	42.1	28.3	2048	$\mathcal{O}(n)$
Entropy-Pruned (Ye et al., 2021)	25%	48.4	2.9	18.2	57.8	41.6	27.1	1536	$\mathcal{O}(n \log n)$
VQ-Tok (Ronen et al., 2023)	33%	47.9	2.8	18.0	58.2	41.2	27.5	1408	$\mathcal{O}(n \cdot d)$
TofuTok (Fu et al., 2021)	39%	39.3	2.3	18.4	56.1	40.4	26.8	1152	$\mathcal{O}(n)$
BPE+Chunk (Ours)	42%	42.1	2.4	18.8	55.2	39.8	26.2	1200	$\mathcal{O}(n)$
SemToken (Ours)	59%	30.4	1.5	17.0	59.9	42.4	28.7	768	$\mathcal{O}(n \cdot d)$

Table 4: Comprehensive comparison across all baselines. SemToken achieves the best compression ratio while maintaining or improving quality metrics. Memory usage includes both model and cache memory.

Task	Dataset	Method	Compression Ratio	PPL/F1	EM	ROUGE-L	Speedup	Memory Reduction
Language Modeling	WikiText-103	BPE	0%	17.3	–	–	1.0×	0%
		Entropy-Pruned	25%	18.2	–	–	1.3×	29%
		VQ-Tok	33%	18.0	–	–	1.3×	32%
		TofuTok	39%	18.4	–	–	1.6×	44%
		BPE+Chunk	42%	18.8	–	–	1.5×	42%
		SemToken	59%	17.0	–	–	2.0×	63%
Question Answering	LongBench	BPE	0%	59.4	42.1	–	1.0×	0%
		Entropy-Pruned	25%	57.8	41.6	–	1.3×	29%
		VQ-Tok	33%	58.2	41.2	–	1.3×	32%
		TofuTok	39%	56.1	40.4	–	1.6×	44%
		BPE+Chunk	42%	55.2	39.8	–	1.5×	42%
		SemToken	59%	59.9	42.4	–	2.0×	63%
Summarization	BookSum	BPE	0%	–	–	28.3	1.0×	0%
		Entropy-Pruned	25%	–	–	27.1	1.3×	29%
		VQ-Tok	33%	–	–	27.5	1.3×	32%
		TofuTok	39%	–	–	26.8	1.6×	44%
		BPE+Chunk	42%	–	–	26.2	1.5×	42%
		SemToken	59%	–	–	28.7	2.0×	63%
Multimodal QA	ChartQA	BPE	0%	–	65.1	–	1.0×	0%
		Entropy-Pruned	25%	–	64.2	–	1.3×	29%
		VQ-Tok	33%	–	64.8	–	1.3×	32%
		TofuTok	39%	–	63.5	–	1.6×	44%
		BPE+Chunk	42%	–	62.9	–	1.5×	42%
		SemToken	61%	–	65.1	–	1.8×	66%

Table 5: Detailed task-specific performance breakdown. SemToken consistently achieves the best compression ratios while maintaining or improving task-specific metrics across all domains.

Model	Parameters	Method	Compression Ratio	Latency (ms/token)	Memory (GB)	PPL
LLaMA-2-7B	7B	BPE	0%	61.2	4.1	17.3
		SemToken	59%	30.4	1.5	17.0
		SemToken + FlashAttn2	57%	22.5	1.5	17.0
GPT-J-6B	6B	BPE	0%	58.7	3.9	18.1
		SemToken	56%	29.8	1.4	17.8
		SemToken + FlashAttn2	54%	21.9	1.4	17.8
GPT-NeoX-20B	20B	BPE	0%	89.3	12.4	16.2
		SemToken	58%	45.2	5.2	15.9
		SemToken + FlashAttn2	55%	33.1	5.2	15.9

Table 6: Performance analysis across different model sizes. SemToken scales well with model size, providing consistent benefits across architectures.