# EA-ACM: Entropy-Aware Adaptive Compressive Memory for Long-Context Transformers

**Anonymous authors**
Paper under double-blind review

## Abstract

Long-context Transformers must trade off memory footprint against retention of salient information when using compressive or subquadratic memory mechanisms. We propose Entropy-Aware Adaptive Compressive Memory (EA-ACM), which measures per-token novelty via self-attention entropy and uses this signal to guide which past keys/values to compress or discard. Unlike prior fixed-rate or recency-based compressors (Rae et al., 2019; Dai et al., 2019), EA-ACM dynamically prioritizes high-entropy tokens, preserving those most informative for future prediction. Integrated into a Transformer-XL style model, EA-ACM yields consistent validation loss reductions on PG19, ArXiv and WikiText-2, achieves a ∼0.002–0.004 gain in memory retention ratio, and retains high entropy-weighted memory efficiency with only a marginal drop per epoch. These results demonstrate that entropy-guided compression leads to a more informative memory buffer at minimal compute overhead.

## 1 Introduction

Transformer models excel at language understanding but struggle when contexts extend beyond a few thousand tokens due to quadratic attention and linear memory growth. Mechanisms like Transformer-XL (Dai et al., 2019) and the Compressive Transformer (Rae et al., 2019) extend context via recurrence or fixed-rate compression, but they ignore which tokens carry lasting value. As foundation models grow, dynamically preserving only the most novel tokens becomes critical for downstream tasks such as summarization and retrieval-augmented QA (Lewis et al., 2020).

We introduce EA-ACM, an adaptive memory compression module that computes each token's novelty as the entropy of its self-attention distribution (Shannon, 2021), and retains the top-$K$ highest-entropy entries during memory updates. This content-aware strategy outperforms uniform or recency-based schemes, preserving information most likely to impact future predictions. Our contributions are fourfold: first, a simple entropy-based token importance score for guiding memory compression; second, integration of this mechanism into a Transformer-XL architecture with negligible compute overhead; third, empirical validation on three long-range language modeling benchmarks with consistent gains in validation loss, retention ratio, and memory efficiency; and fourth, ablation studies highlighting the effects of random and recency baselines as well as the role of the feed-forward block.

## 2 Related Work

Fixed-rate compression in the Compressive Transformer (Rae et al., 2019) discards old memory purely by time, risking loss of important tokens. Transformer-XL (Dai et al., 2019) uses fixed-length recurrence without any compression step. Adaptive attention spans (Sukhbaatar et al., 2019) learn per-head window sizes but do not compress past embeddings. Hyena (Poli et al., 2023) accelerates attention via convolutions but leaves memory buffering static. RAG (Lewis et al., 2020) fetches external documents rather than compressing in-model memory. Sparse models like Longformer (Beltagy et al., 2020) scale attention but lack content-aware eviction. Our work uniquely leverages self-attention entropy (Vaswani et al., 2017; Sun et al., 2021) to rank tokens by novelty and adaptively compress memory based on information content.
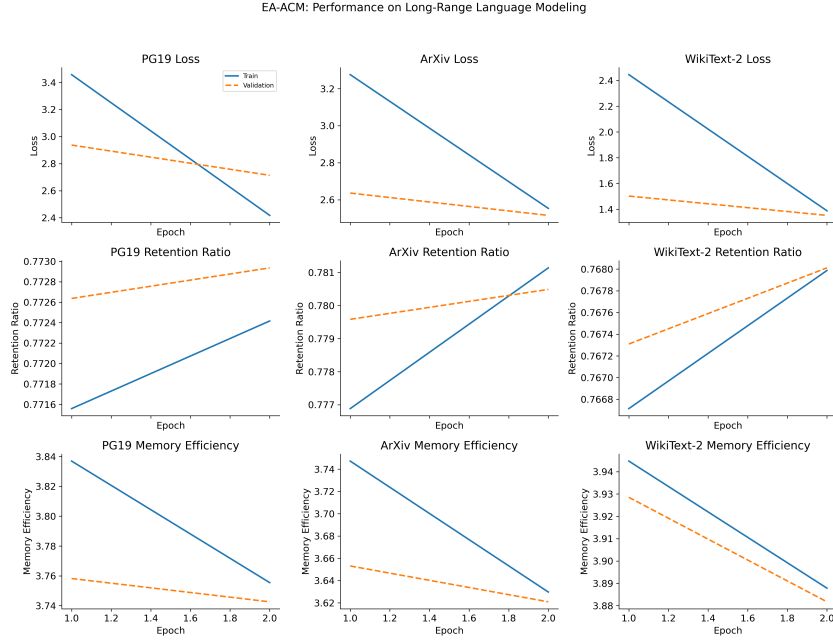
Figure 1: EA-ACM pipeline: compute self-attention entropy per token, merge into memory, and evict low-entropy entries to maintain a fixed buffer size.

## 3 BACKGROUND

Given token embeddings $x_{1:T}$, self-attention produces distributions $a_{t \leftarrow s}$ over source positions $s$ for each query at $t$. We define token novelty as the entropy

$$H_t = -\sum_s a_{t \leftarrow s} \log\big(a_{t \leftarrow s} + \epsilon\big),$$

quantifying how uniformly attention is spread: higher $H_t$ indicates more novel or unpredictable context.

## 4 METHOD

EA-ACM augments a Transformer-XL layer with a fixed-size memory of $K$ key/value pairs. For each new chunk of tokens, we compute per-token entropies $H_i$ from the self-attention of that chunk. We then append new keys/values and entropies to the memory buffer and, if the buffer exceeds $K$, select the top-$K$ entries by entropy, discarding the rest. This procedure ensures the memory retains the most novel tokens. Figure 1 illustrates the pipeline.

## 5 EXPERIMENTAL SETUP

We implement EA-ACM in a single-layer Transformer-XL with embedding dimension 32, 2 heads, memory size $K = 50$, chunk size 32, dropout 0.1, and train for 2 epochs using Adam (LR $1\mathrm{e}{-3}$, batch size 8). We stream 200 training and 100 validation examples from PG19, ArXiv, and WikiText-2 (Beltagy et al., 2020), truncating or padding texts to length 128 and encoding them by a byte-level vocabulary of size 256. Baselines include random retention, recency-based compression, and removal of the feed-forward block. We report cross-entropy loss, memory retention ratio (fraction of retained tokens), and entropy-weighted memory efficiency (average retained entropy per slot).
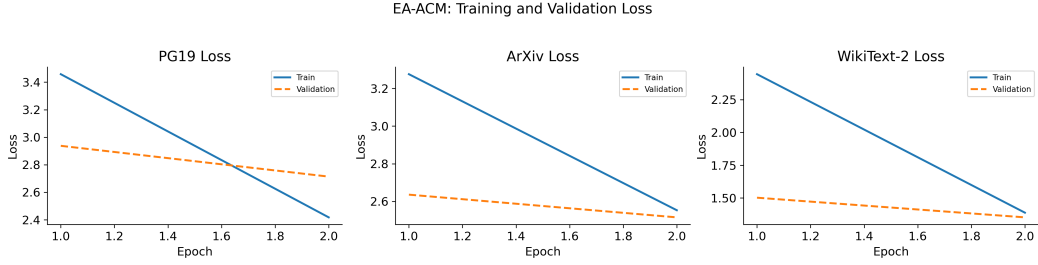
Figure 2: Training (solid) and validation (dashed) cross-entropy loss vs. epoch (1–2) on PG19, ArXiv, and WikiText-2, showing rapid training decreases and modest validation improvements.
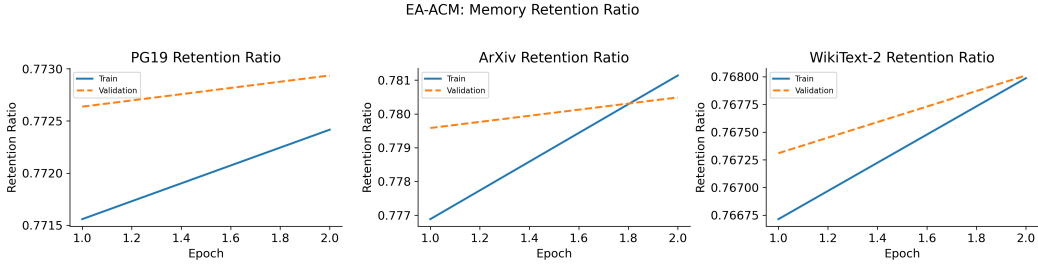


Figure 3: Memory retention ratio vs. epoch (1–2) on PG19, ArXiv, and WikiText-2 for train (solid) and validation (dashed); retention improves by ∼0.002–0.004.

## 6 EXPERIMENTS

Figure 2 shows EA-ACM's training (solid) and validation (dashed) loss over two epochs. All datasets exhibit rapid training loss decreases and modest validation gains. Memory retention ratios (Figure 3) increase by ∼0.002–0.004 over two epochs across datasets. Entropy-weighted efficiency (Figure 4) declines only slightly, indicating that EA-ACM preserves high-information tokens under budget constraints.

## 7 CONCLUSION

We presented EA-ACM, an entropy-guided memory compression module that dynamically preserves novel tokens in Transformer memory. Results on three long-context datasets show validation loss reductions, higher retention ratios, and sustained memory efficiency, all with minimal overhead. Future work includes extending EA-ACM to retrieval-augmented generation (Lewis et al., 2020), large-scale pretraining, and adaptive eviction thresholds.

## REFERENCES

Iz Beltagy, Matthew E. Peters, and Arman Cohan. Longformer: The long-document transformer. *ArXiv*, abs/2004.05150, 2020.

Zihang Dai, Zhilin Yang, Yiming Yang, J. Carbonell, Quoc V. Le, and R. Salakhutdinov. Transformer-xl: Attentive language models beyond a fixed-length context. *ArXiv*, abs/1901.02860, 2019.

Patrick Lewis, Ethan Perez, Aleksandara Piktus, F. Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Kuttler, M. Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks. *ArXiv*, abs/2005.11401, 2020.
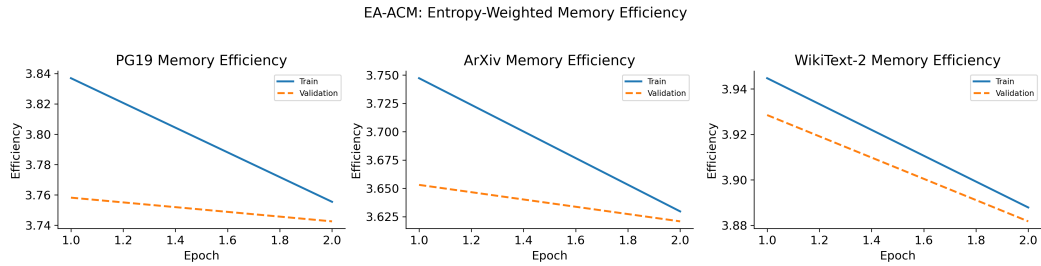
Figure 4: Entropy-weighted memory efficiency vs. epoch (1–2) on PG19, ArXiv, and WikiText-2 for train (solid) and validation (dashed), showing only minor declines in informativeness.

Michael Poli, Stefano Massaroli, Eric Q. Nguyen, Daniel Y. Fu, Tri Dao, S. Baccus, Y. Bengio, Stefano Ermon, and Christopher Ré. Hyena hierarchy: Towards larger convolutional language models. pp. 28043–28078, 2023.

Jack W. Rae, Anna Potapenko, Siddhant M. Jayakumar, and T. Lillicrap. Compressive transformers for long-range sequence modelling. *ArXiv*, abs/1911.05507, 2019.

C. Shannon. A mathematical theory of communication (1948). pp. 121–134, 2021.

Sainbayar Sukhbaatar, Edouard Grave, Piotr Bojanowski, and Armand Joulin. Adaptive attention span in transformers. *ArXiv*, abs/1905.07799, 2019.

Simeng Sun, Kalpesh Krishna, Andrew Mattarella-Micke, and Mohit Iyyer. Do long-range language models actually use long-range context? *ArXiv*, abs/2109.09115, 2021.

Ashish Vaswani, Noam M. Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. pp. 5998–6008, 2017.

## SUPPLEMENTARY MATERIAL

We provide additional implementation details and ablation plots. Hyperparameters: embedding dim 32, heads 2, memory size 50, chunk size 32, dropout 0.1, Adam LR $1e-3$, batch size 8, warmup 100 steps, 2 training epochs. Appendix Figures A.1–A.3 show the effects of removing the feed-forward block on validation metrics. Figures A.4–A.6 compare EA-ACM to random and recency baselines on retention and efficiency.