

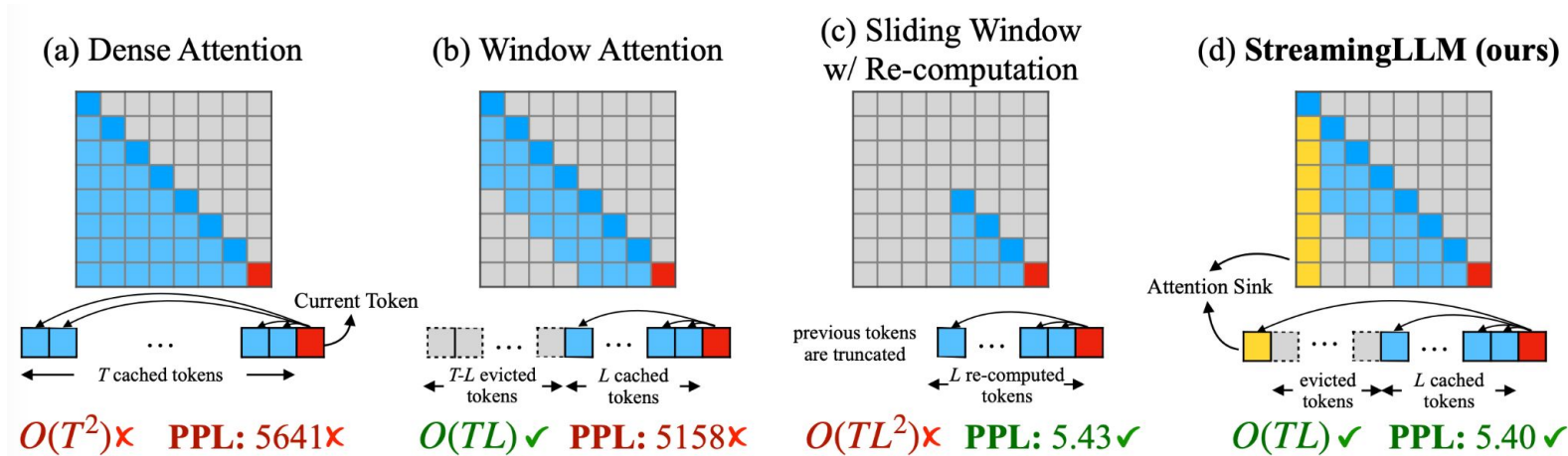
# Training Free Exponential Context Extension via Cascading KV Cache

Jeffrey Willette<sup>1</sup>, Heejun Lee<sup>1</sup>, Youngwan Lee<sup>12</sup>, Myeongjae Jeon<sup>3</sup>, Sung Ju Hwang<sup>14</sup>

<sup>1</sup>KAIST <sup>2</sup>ETRI <sup>3</sup>POSTECH <sup>4</sup>DeepAuto.ai

# Sliding Window Attention (Streaming LLM)

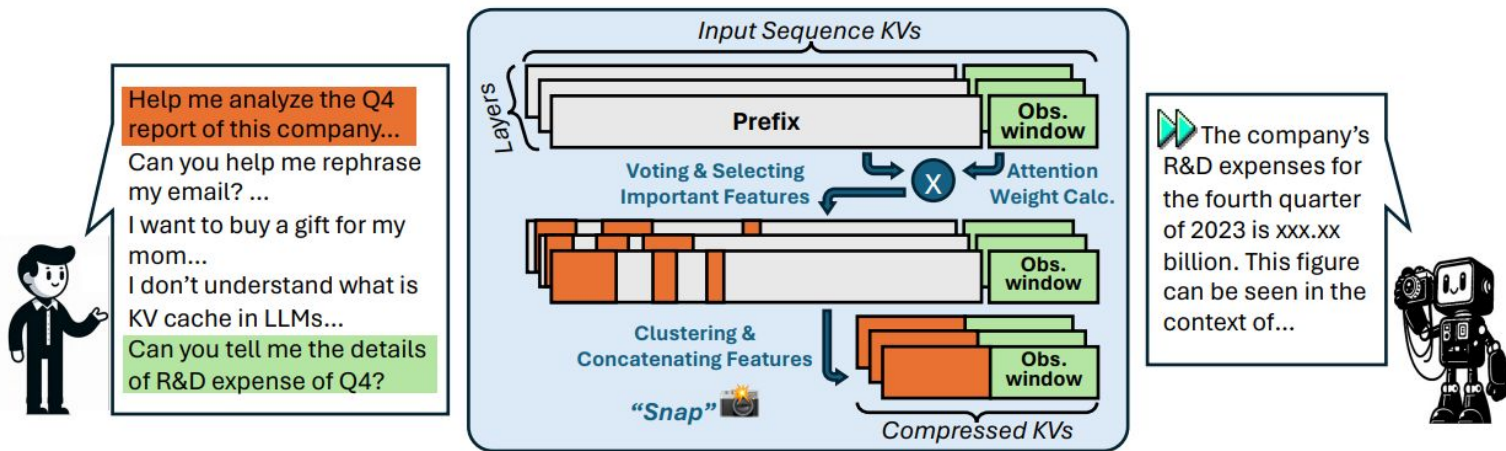
Recent works have tackled the problem of stabilizing sliding window attention during inference.



However, the token eviction strategy is naive, which limits the effectiveness of the model.

# KV Caching (Recent Works)

Recent works (such as H2O [1] and SnapKV [2]) have proposed KV cache compression.



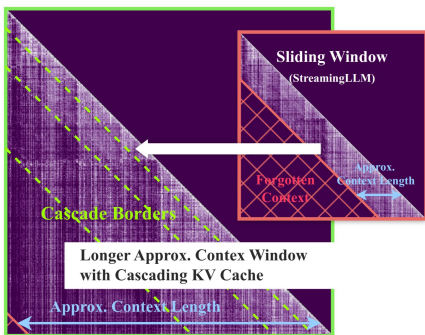
However, these methods still require quadratic complexity during the prefill stage and only prune tokens after prefill and before decoding.

[1] Zhang, Z., Sheng, Y., Zhou, T., Chen, T., Zheng, L., Cai, R., ... & Chen, B. (2023). H2o: Heavy-hitter oracle for efficient generative inference of large language models. *Advances in Neural Information Processing Systems*, 36, 34661-34710.

[2] Li, Y., Huang, Y., Yang, B., Venkitesh, B., Locatelli, A., Ye, H., ... & Chen, D. (2024). Snapkv: Llm knows what you are looking for before generation. *Advances in Neural Information Processing Systems*, 37, 22947-22970.

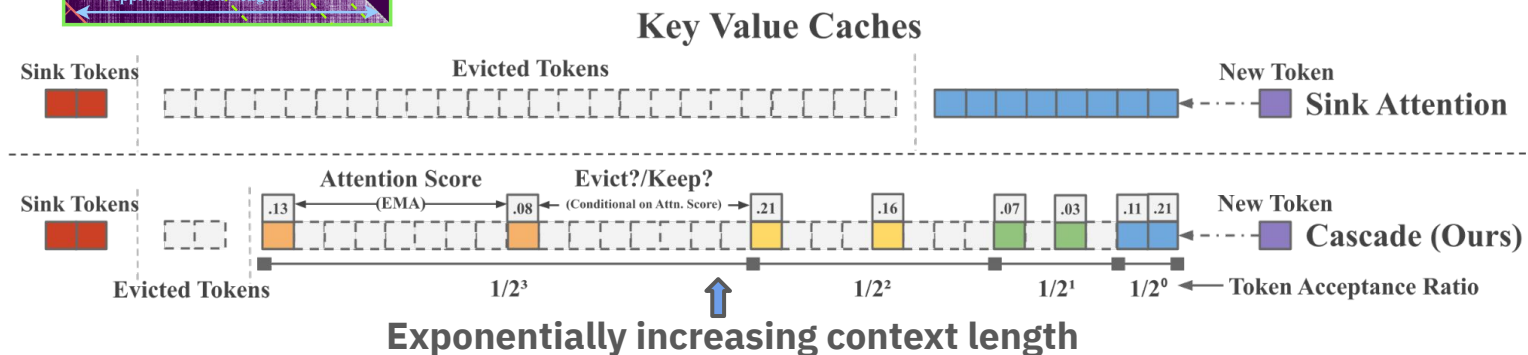
# Cascading KV Cache: Exponential Window Extension

We proposed a method called Cascading KV caching to further extend the context window size of any pretrained transformer by retaining more important tokens while evicting less important ones.



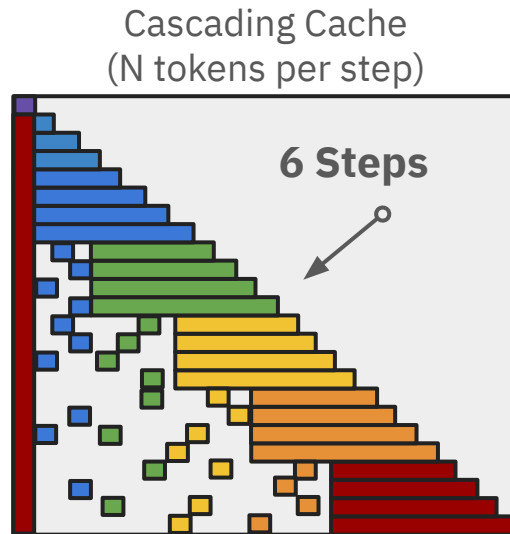
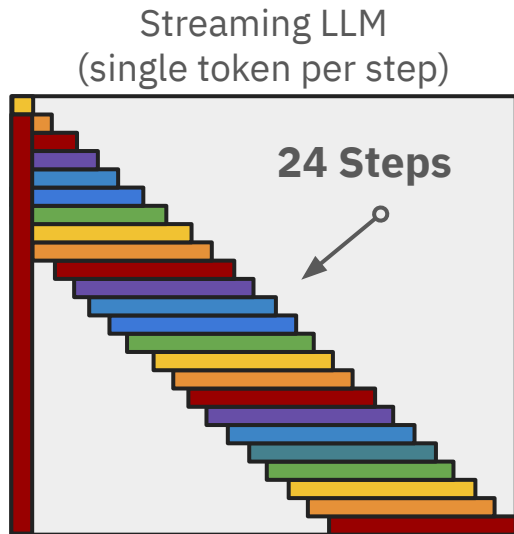
The dynamic cascading KV cache:

- Evicts less relevant tokens
- Keeps important tokens for as long as possible
- Maintains important patterns in the attention matrix
- Maintains linear inference complexity
- Uses the same memory as a fixed sliding window



# Cascading KV Cache: Batch Processing Prefill Tokens

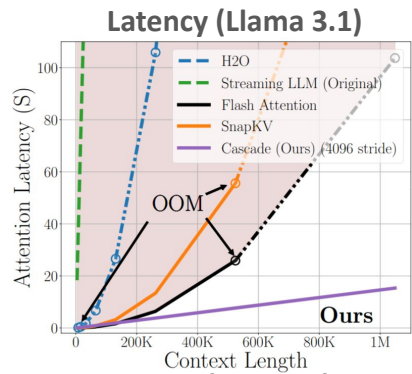
The key to realizing actual latency improvements over other linear sliding window methods lies in the way our Cascading KV Cache processes batches of tokens at a time instead of single tokens.



Using this method, Cascading KV Cache delivers a more than 2 order of magnitude improvement in latency over Streaming LLM.

# Cascading KV Cache: Latency and Long Context Tasks

Given the same cache size, our method shows stronger performance on long context tasks, while **scaling latency linearly**, and **keeping memory size constant**.



1. Outperforms static window on long context tasks

2. Linearly Scaling Latency

3. Better than random at 16x cache size

**1M Token Passkey Retrieval (Llama 3.1)**  
Total Cache Size = 65K, 8 Cascades

Ours

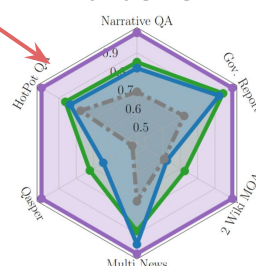
Insert Loc (%)	0-20	20-40	40-60	60-80	80-100
32K	100	100	100	100	100
65K	100	100	100	100	100
131K	100	100	96	97	100
262K	90	77	87	87	100
524K	59	38	61	87	100
1M	23	21	33	45	88

Streaming LLM

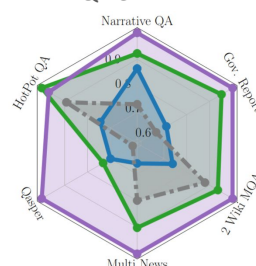
Insert Loc (%)	0-20	20-40	40-60	60-80	80-100
32K	100	100	100	100	100
65K	100	100	100	100	100
131K	14	16	61	100	100
262K	13	9	10	45	100
524K	12	13	14	12	59
1M	11	11	8	11	48

**LongBench (Llama 3.1)**

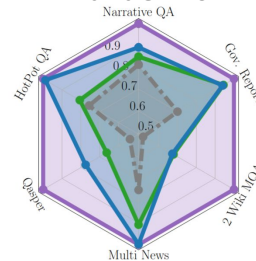
Llama 3.1 8B



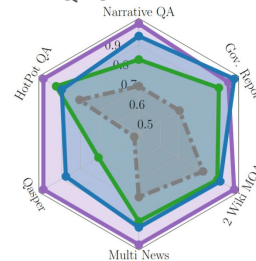
Qwen 2 7B



Llama 3.1 70B



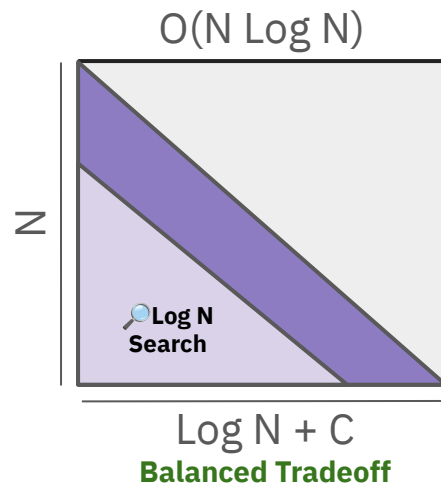
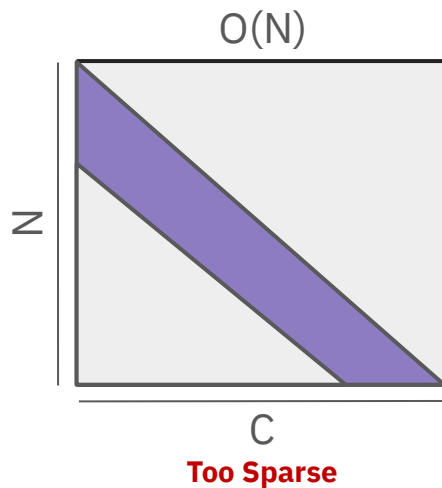
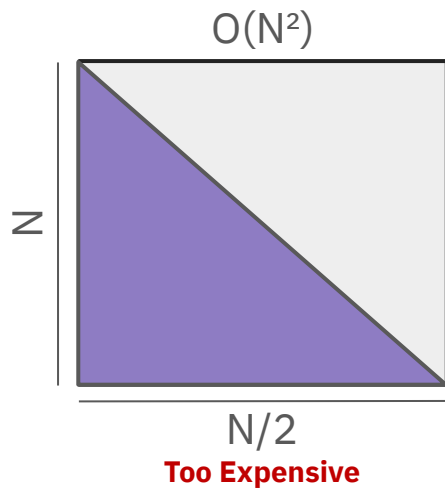
Qwen 2 72B



— Streaming LLM — Big Bird - - - Flash Attention (truncated) — Cascade (Ours)

# Future Research: Logarithmic Token Search

Cascading KV cache must eventually evict tokens. However, if all tokens can be stored and relevant tokens can be retrieved in logarithmic time, then an overall  $O(N \log N)$  transformer would result.



Devising an effective logarithmic token search would allow for practically infinite context lengths while maintaining all tokens in memory so they may retain influence over future predictions.

**Thank you! We look forward to seeing you at  
ICLR 2025!**