

# Core Context Aware Transformers for Long Context Language Modeling

ICML 2025

Yaofo Chen<sup>1</sup>, Zeng You<sup>1,2</sup>, Shuhai Zhang<sup>1,3</sup>, Haokun Li<sup>1,2</sup>, Yiru Li<sup>1</sup>,  
Yaowei Wang<sup>2,4</sup>, Mingkui Tan<sup>1,3,5</sup>

South China University of Technology<sup>1</sup>, Peng Cheng Laboratory<sup>2</sup>, Pazhou Laboratory<sup>3</sup>,  
Harbin Institute of Technology<sup>4</sup>, Key Laboratory of Big Data and Intelligent Robot, Ministry of Education<sup>5</sup>

June 10, 2025

# Outline

1. Background
2. Core Context Aware Attention
3. Experimental Results
4. Conclusion

# Outline

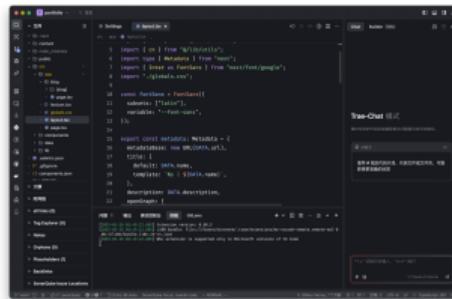
1. Background
2. Core Context Aware Attention
3. Experimental Results
4. Conclusion

# Large Language Model

Large language models (LLMs) have demonstrated exceptional proficiency across various applications, particularly in natural language understanding and generation.



chatbot



coding



education

Figure: Applications of LLMs: empowering chatbots, coding assistance, and educational tools

# Importance of Long-context Modeling

The ability to process longer contexts has been a key factor in improving the performance of LLMs across a wide range of tasks.



Legal document analysis

why long-context modeling ability is important for LLM?

Okay, the user is asking why long-context modeling is important for LLMs. Let me start by breaking down the key reasons. First, I recall that real-world tasks often involve long documents or conversations. For example, legal contracts or research papers can be hundreds of pages. If an LLM can't handle that, it misses crucial context.

Hmm, the user might be a developer or researcher working with LLMs. They probably encountered limitations with standard models and need to understand the value of longer contexts. Their deeper need could be justifying the use of more advanced models or optimizing their current setup.

I should explain coherence in extended interactions. Like in customer service chats, past messages matter. Without context, responses feel disjointed. Also, information retrieval: finding answers in long texts requires understanding the whole document. Missing sections could lead to incomplete answers.

Then there's complex reasoning. Tasks like summarizing a book or analyzing codebases need connecting ideas across many pages. Short contexts force chunking, which loses connections. The user might be dealing with such tasks and facing accuracy issues.

Oh, and domain-specific applications! Medical records or financial reports are lengthy. Professionals in these fields need LLMs that grasp the entire document. Maybe the user is in such a field and needs models to handle their data better.

Chain-of-thought reasoning

**Figure:** Long context matters: enabling advanced long text analysis and chain-of-thought reasoning with large language models.

# Challenges of Processing Long Context

When the context length becomes very large (*i.e.*, 128K), the redundant information in the context tends to increase, which not only **incurs unnecessary computational and storage overhead** but also **may hamper the modeling representation performance**.

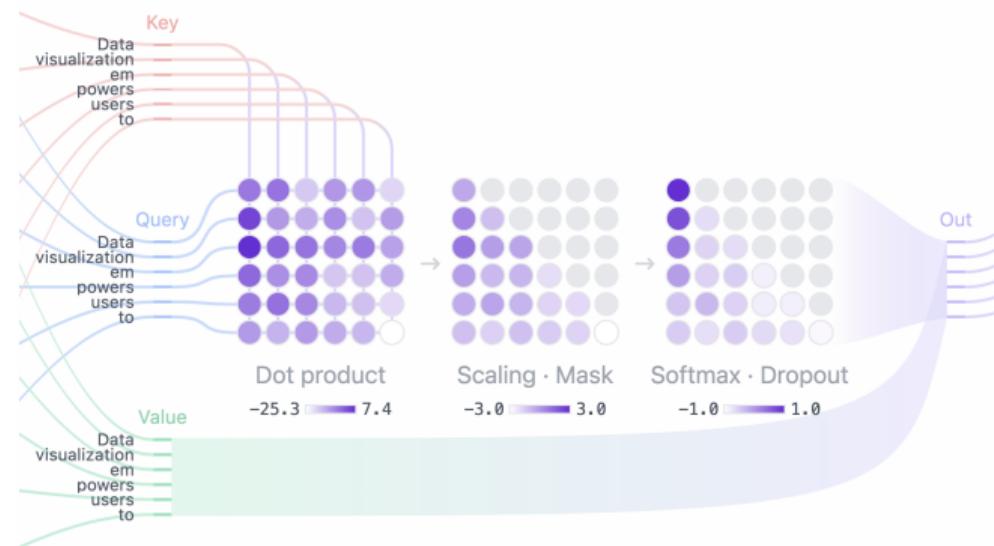


Figure: The computational complexity and KV cache of self-attention grows quadratically and linearly with the length of the context. This figure is from <https://poloclub.github.io/transformer-explainer>.

## Preliminaries

Given a sequence of tokens  $\mathbf{X} = [\mathbf{X}_1; \mathbf{X}_2; \dots; \mathbf{X}_L] \in \mathbb{R}^{L \times d}$ , where each token  $\mathbf{X}_i \in \mathbb{R}^{1 \times d}$ , the self-attention begins by transforming  $\mathbf{X}$  into a query  $\mathbf{Q} = \mathbf{X}\mathbf{W}^Q$ , a key  $\mathbf{K} = \mathbf{X}\mathbf{W}^K$ , and a value  $\mathbf{V} = \mathbf{X}\mathbf{W}^V$  through three linear projections  $\mathbf{W}^{Q,K,V} \in \mathbb{R}^{d \times d}$ . The self-attention is calculated by

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left( \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d}} \right) \mathbf{V}. \quad (1)$$

The computational and storage complexity:

- The computational complexity of self-attention grows **quadratically** with the length of the context.
- During inference, we cache  $\mathbf{K}$  and  $\mathbf{V}$  for accelerating attention computation. The size and KV cache grows **linearly** with the length of the context.

# Outline

1. Background
2. Core Context Aware Attention
3. Experimental Results
4. Conclusion

# Motivation

As the context length scales to extremely large magnitudes, it becomes impractical for a token to maintain significant semantic connections with all tokens in the context.

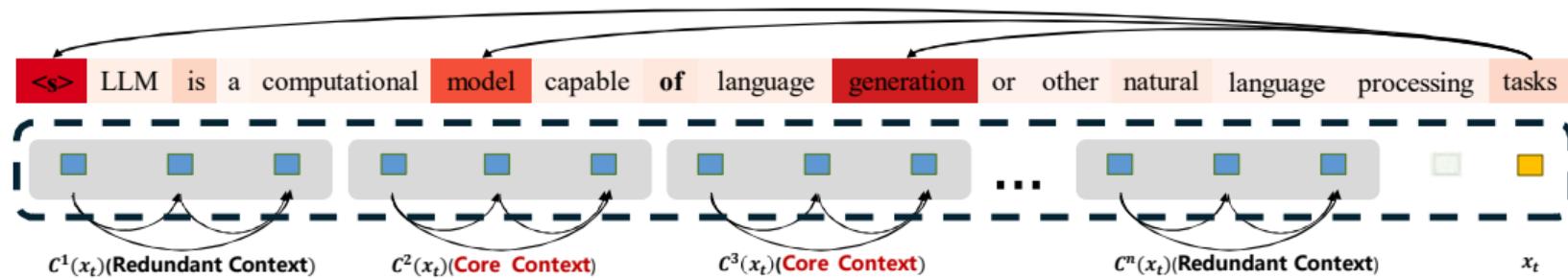


Figure: Illustration of core contexts and redundant contexts. We show attention scores of the last token relative to the other tokens in LLaMA2-7B (darker shadows indicate higher attention scores). The last token exhibits high attention scores towards core contexts. The remains are considered as redundant contexts, introducing unnecessary computational overhead for attention.

# Motivation

The core context (e.g.,  $C^2(\mathbf{x}_t)$  w.r.t.  $\mathbf{x}_t$  in Figure) refers to the contextual information that is highly relevant to the token  $\mathbf{x}_t$ , which is crucial for the token's representation. Therefore, the model should **prioritize the core context over redundant parts**.

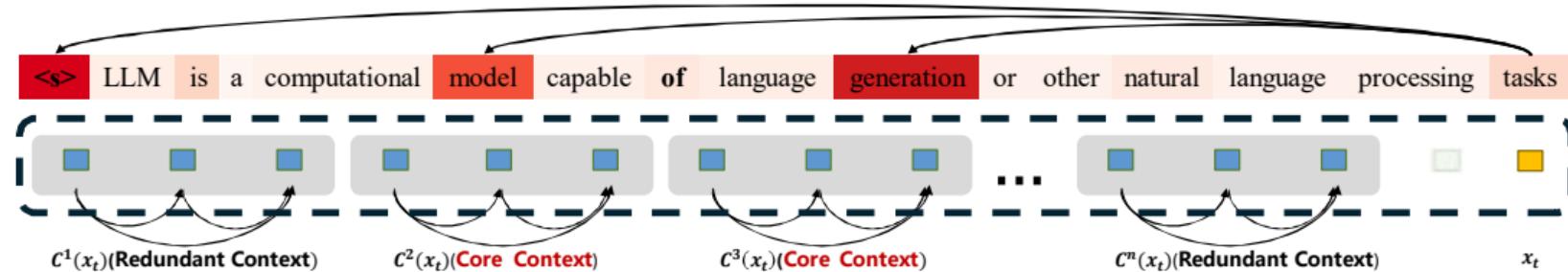


Figure: Illustration of core contexts and redundant contexts. We show attention scores of the last token relative to the other tokens in LLaMA2-7B (darker shadows indicate higher attention scores). The last token exhibits high attention scores towards core contexts. The remains are considered as redundant contexts, introducing unnecessary computational overhead for attention.

## Overview

We propose a **Core Context Aware Attention** that employs globality-aware pooling and locality-preserving modules to capture both global and local context dependencies.

- The globality-aware pooling module operates by generating representative core tokens from segmented groups of the input sequence.
- The locality-preserving module is responsible for capturing the local information of  $s$  neighborhood tokens to ensure comprehensive coverage.
- We devise a differentiable fusion strategy to combine the insights from global and local modules to generate the final outputs.

# Overview

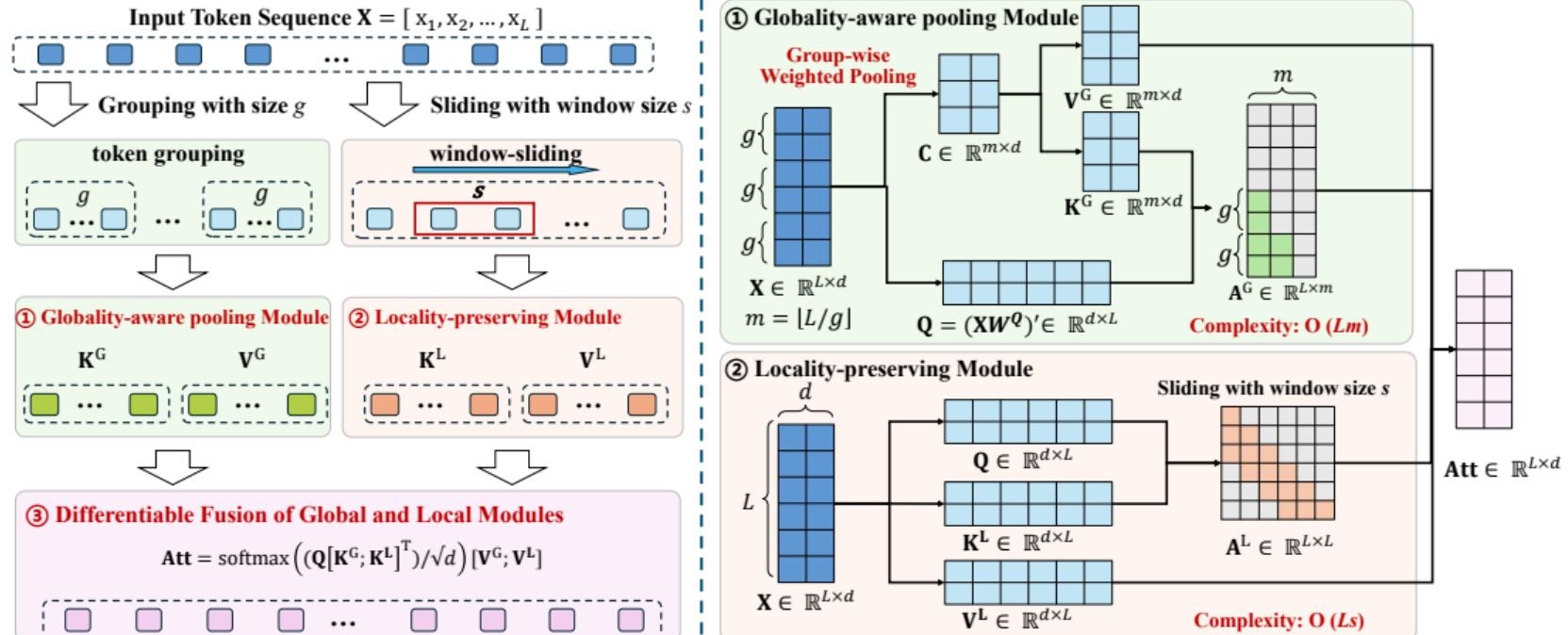


Figure: Illustration of CCA-Attention, which includes two components: Globality-aware pooling module and Locality-preserving module. We produce the final output  $\text{Att}$  by fusing these two modules.

## Globality-aware Pooling Module: Core Tokens Generation

**Motivation:** The context redundancy indicates that computational resources can be dynamically allocated to core contexts while reducing emphasis on the remaining ones.

**Solution:** Given input tokens  $\mathbf{X} = [\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_L] \in \mathbb{R}^{L \times d}$ , we segment it into  $m = \lfloor L/g \rfloor$  groups, each group containing  $g$  tokens (denoting the  $i$ -th group by  $\mathbf{X}_{\mathcal{I}_i}^G \in \mathbb{R}^{g \times d}$ ).

To identify prominent tokens in the  $i$ -th group, we devise a group-wise weighted pooling strategy that employs the last token  $\mathbf{x}_{ig}$  to evaluate the importance. Formally, we derive one core token  $\mathbf{c}_i \in \mathbb{R}^{1 \times d}$  from each group by

$$\mathbf{c}_i = \text{softmax} \left( \frac{\mathbf{Q}_{ig} \mathbf{K}_{\mathcal{I}_i}'}{\sqrt{d}} \right) \mathbf{X}_{\mathcal{I}_i}^G \in \mathbb{R}^{1 \times d}, i = 1, \dots, m, \quad (2)$$

where  $\mathbf{K}'_{\mathcal{I}_i} = \mathbf{X}_{\mathcal{I}_i}^G \mathbf{W}^K$ . The core token  $\mathbf{c}_i$  encapsulates crucial information of the group. With  $m$  groups, we derive  $m$  core tokens  $\mathbf{C} = [\mathbf{c}_1; \mathbf{c}_2; \dots; \mathbf{c}_m] \in \mathbb{R}^{m \times d}$ .

## Globality-aware Pooling Module: Core Tokens as Efficient Proxy

To reduce the redundancy, we use the sequence of core tokens  $\mathbf{C}$  instead of the original tokens  $\mathbf{X}$  for attention computation. Formally, we adopt core tokens to calculate the key and value matrices for a specific query  $\mathbf{Q}_i$  as follows

$$\begin{aligned}\tilde{\mathbf{K}}_{T_i}^G &= [\mathbf{K}_1^G; \dots; \mathbf{K}_j^G] \in \mathbb{R}^{j \times d}, \quad \tilde{\mathbf{V}}_{T_i}^G = [\mathbf{V}_1^G; \dots; \mathbf{V}_j^G] \in \mathbb{R}^{j \times d}, \\ \mathbf{K}^G &= \mathbf{C}\mathbf{W}^K \in \mathbb{R}^{m \times d}, \quad \mathbf{V}^G = \mathbf{C}\mathbf{W}^V \in \mathbb{R}^{m \times d},\end{aligned}\tag{3}$$

The index  $j$  in the above Eqn. can be calculated as  $j = \max(0, \lfloor (i-s)/g \rfloor)$ . When the context is short ( $i < (g+s)$ ), the key and value  $\mathbf{K}^G$ ,  $\mathbf{V}^G$  would be excluded from attention calculation since the redundancy in the context is negligible.

During inference, as tokens are generated sequentially, we derive a new core token via the aforementioned Eqn., once the number of generated tokens reaches  $g$ . Different from the full self-attention, we cache  $\tilde{\mathbf{K}}^G$  and  $\tilde{\mathbf{V}}^G$  for inference.

## Locality-preserving Module

**Challenge:** The globality-aware pooling module focuses on coarse-grained global information, potentially overlooking fine-grained local context.

**Solution:** We introduce a locality-preserving module that complements the globality-aware pooling module by focusing on neighboring tokens to capture detailed local dependencies. The key and value matrices for a specific query  $\mathbf{Q}_i$  in the locality-preserving module are defined as follows:

$$\begin{aligned}\tilde{\mathbf{K}}_{\mathcal{U}_i}^L &= [\mathbf{K}_k^L; \dots; \mathbf{K}_i^L], \quad \tilde{\mathbf{V}}_{\mathcal{U}_i}^L = [\mathbf{V}_k^L; \dots; \mathbf{V}_i^L], \\ \mathbf{K}^L &= \mathbf{XW}^K, \mathbf{V}^L = \mathbf{XW}^V\end{aligned}\tag{4}$$

During generation, it is challenging to maintain the number of tokens as a multiple of the group size  $g$ . To address this, we set the local window size to  $s + ((i - s) \bmod g)$ .

# Differentiable Fusion of Global and Local Modules

**Challenge:** Both globality-aware pooling and locality-preserving modules involve only a portion of tokens, leading to a limited comprehensive understanding of the context.

**Solution:** We seek to combine the involved tokens of these two attentions by concatenating the key and value matrices from both attentions, *i.e.*,  $[\tilde{\mathbf{K}}_{\mathcal{T}_i}^G; \tilde{\mathbf{K}}_{\mathcal{U}_i}^L]$  and  $[\tilde{\mathbf{V}}_{\mathcal{T}_i}^G; \tilde{\mathbf{V}}_{\mathcal{U}_i}^L]$ , to leverage the combined information. Formally, the proposed CCA-Attention is computed as follows:

$$\mathbf{Att}_i = \text{softmax}\left(\frac{\mathbf{Q}_i [\tilde{\mathbf{K}}_{\mathcal{T}_i}^G; \tilde{\mathbf{K}}_{\mathcal{U}_i}^L]^\top}{\sqrt{d}}\right) [\tilde{\mathbf{V}}_{\mathcal{T}_i}^G; \tilde{\mathbf{V}}_{\mathcal{U}_i}^L]. \quad (5)$$

We represent the final output of our CCA-Attention as  $\mathbf{Att} = [\mathbf{Att}_1; \mathbf{Att}_2; \dots; \mathbf{Att}_L]$ .

# Algorithm

**Require:** Inputs  $\mathbf{X} = [\mathbf{x}_1; \mathbf{x}_2; \dots; \mathbf{x}_L]$ , parameters  $\mathbf{W}^{Q,K,V}$ , group size  $g$ , local window size  $s$ .

- 1: Calculate the query  $\mathbf{Q} = \mathbf{X}\mathbf{W}^Q$ , #groups  $m = \lfloor L/g \rfloor$
- 2: **for**  $i$  in  $\{1, 2, \dots, m\}$  **do**
- 3:    $\mathbf{X}_{\mathcal{I}_i}^G = [\mathbf{x}_{(i-1)g+1}; \mathbf{x}_{(i-1)g+2}; \dots; \mathbf{x}_{ig}]$
- 4:    $\mathbf{c}_i = \text{softmax} \left( \frac{\mathbf{Q}_{ig} \mathbf{K}_{\mathcal{I}_i}^{\prime \top}}{\sqrt{d}} \right) \mathbf{X}_{\mathcal{I}_i}^G$ , where  $\mathbf{K}_{\mathcal{I}_i}' = \mathbf{X}_{\mathcal{I}_i}^G \mathbf{W}^K$
- 5: **end for**
- 6: Let  $\mathbf{C} = [\mathbf{c}_1; \mathbf{c}_2; \dots; \mathbf{c}_m]$
- 7:  $\mathbf{K}^G = \mathbf{C}\mathbf{W}^K$ ,  $\mathbf{V}^G = \mathbf{C}\mathbf{W}^V$  // *Globality-aware Pooling Module*
- 8:  $\mathbf{K}^L = \mathbf{X}\mathbf{W}^K$ ,  $\mathbf{V}^L = \mathbf{X}\mathbf{W}^V$  // *Locality-preserving Module*
- 9: **for**  $i$  in  $\{1, 2, \dots, L\}$  **do**
- 10:    $\tilde{\mathbf{K}}_{\mathcal{T}_i}^G = \mathbf{K}_{1:j}^G$ ,  $\tilde{\mathbf{V}}_{\mathcal{T}_i}^G = \mathbf{V}_{1:j}^G$ ,  $j = \max(0, \lfloor (i-s)/g \rfloor)$
- 11:    $\tilde{\mathbf{K}}_{\mathcal{U}_i}^L = \mathbf{K}_{k:i}^L$ ,  $\tilde{\mathbf{V}}_{\mathcal{U}_i}^L = \mathbf{V}_{k:i}^L$ ,  $k = \max(1, i-s - ((i-s) \bmod g))$
- 12:    $\mathbf{Att}_i = \text{softmax}(\mathbf{Q}_i [\tilde{\mathbf{K}}_{\mathcal{T}_i}^G; \tilde{\mathbf{K}}_{\mathcal{U}_i}^L]^{\top})/\sqrt{d} [\tilde{\mathbf{V}}_{\mathcal{T}_i}^G; \tilde{\mathbf{V}}_{\mathcal{U}_i}^L]$
- 13: **end for**
- 14: **Return:** Representations of tokens  $\mathbf{Att} = [\mathbf{Att}_1; \dots; \mathbf{Att}_L]$

# Efficient and Parallel Implementation

We implement our CCA-Attention by leveraging Triton([github.com/triton-lang/triton](https://github.com/triton-lang/triton)) to perform low-level operator fusion between our globality-aware pooling and locality-preserving modules. This enables us to integrate our CCA-Attention as a standalone, cache-friendly operator, effectively eliminating redundant computations.

```

1 @triton.jit
2 def _global_pooling_attn_fwd_inner(
3     global_acc, l_i, m_i, q, #
4     K_block_ptr, V_block_ptr, #
5     start_m, qk_scale, #
6     BLOCK_M: tl.constexpr, HEAD_DIM: tl.constexpr, BLOCK_N: tl.constexpr, #
7     STAGE: tl.constexpr, offs_m: tl.constexpr, offs_n: tl.constexpr, #
8     N_CTX: tl.constexpr, group_size: tl.constexpr, fp8_v: tl.constexpr
9 ):
10    N_GROUP = N_CTX//group_size
11    lo, hi = 0, (start_m + 1) * BLOCK_M//group_size
12    K_block_ptr = tl.advance(K_block_ptr, (0, lo))
13    V_block_ptr = tl.advance(V_block_ptr, (lo, 0))
14    # loop over k, v and update accumulator
15    for start_n in range(lo, hi, BLOCK_N):
16        # -- compute qk ----
17        k = tl.load(K_block_ptr, boundary_check=(0,1))
18        qk = tl.dot(q, k)
19        offs_group = offs_m//group_size
20        mask = offs_group[:, None] > (start_n + offs_n[None, :])
21        qk = qk * qk_scale + tl.where(mask, 0, -1.0e6)
22        # print(f'from:{start_m},mask:{mask},qk:{qk[:]}')
23        m_ij = tl.maximum(m_i, tl.max(qk, 1))
24        qk -= m_ij[:, None]
25        p = tl.math.exp2(qk)
26        l_ij = tl.sum(p, 1)
27
28        alpha = tl.math.exp2(m_i - m_ij)
29        l_i = l_i * alpha + l_ij
30        # -- update output accumulator --
31        global_acc = global_acc * alpha[:, None]
32        # update acc

```

# Outline

1. Background

2. Core Context Aware Attention

3. Experimental Results

4. Conclusion

# Experimental Results

We apply our CCA-Attention to pretrained LLMs and compare it with state-of-the-art efficient attention methods in terms of long-context modeling and efficiency.

## Dataset and Evaluation Metrics

- **LongBench** is a benchmark for the bilingual, multi-task, and comprehensive assessment of LLMs' long context understanding capabilities.
- **Exact Match Score (EM Score)** in multidocument QA is a metric for measuring the model's ability to find the key information within a long context.

## Models and Compared Methods

- **Models:** We apply our proposed CCAAttention to LLaMA2-7B-32K, LLaMA2-7B-80K, LLaMA3.1-8B-128K, and Qwen2.5-7B-128K.
- **Compared Methods:** We compare our proposed CCA-Attention with vanilla attention, StreamingLLM, LM-infinite, and MInference.

# Comparisons on Long Context Modeling

## Comparisons on Longbench-E:

- CCA-LLM outperforms all efficient attention baselines across varying context lengths, attributed to its global-aware pooling for enhanced core-context focus.
- CCA-LLM attains the lowest inference latency and memory usage among competitors, significantly reducing storage requirements while maintaining speed.

Methods	S. QA	M. QA	Sum.	FS. Learning	Synthetic	Code	Avg.	FTL (s)	Mem. (GB)
<i>LLaMA2-7B-32K (Vanilla Self-Attention)</i>	2.75	1.85	<b>12.43</b>	<b>66.28</b>	0.34	<b>48.99</b>	<b>22.11</b>	9.15	35.58
• StreamingLLM	<b>4.75</b>	2.94	2.97	48.20	<u>0.66</u>	30.16	14.95	5.75 (1.6×)	22.94 (35%↓)
• LM-Infinite	2.04	2.33	1.98	57.45	0.3	48.46	18.76	4.72 (1.9×)	26.35 (26%↓)
• MInference	<u>3.68</u>	<u>3.05</u>	<u>10.97</u>	<u>66.26</u>	0.61	42.30	21.14	4.20 (2.2×)	33.52 (6%↓)
• CCA-LLM (Ours)	<b>3.63</b>	<b>3.98</b>	<b>7.79</b>	61.79	<b>2.64</b>	<b>51.36</b>	<u>21.86</u>	<b>2.59 (3.5×)</b>	<b>19.12 (46%↓)</b>
<i>LLaMA2-7B-80K (Vanilla Self-Attention)</i>	<u>3.22</u>	2.71	3.90	<b>64.98</b>	0.56	<b>59.16</b>	<b>22.42</b>	32.43	60.03
• StreamingLLM	2.07	2.32	0.37	45.03	<b>2.67</b>	37.17	14.94	9.04 (3.6×)	37.45 (37%↓)
• LM-Infinite	2.54	1.53	2.22	61.29	<u>1.08</u>	<u>58.54</u>	21.20	8.27 (3.9×)	41.54 (31%↓)
• MInference	2.44	<u>3.49</u>	<u>4.41</u>	<u>64.26</u>	0.28	57.60	22.08	8.14 (4.0×)	54.09 (10%↓)
• CCA-LLM (Ours)	<b>5.62</b>	<b>4.34</b>	<b>8.99</b>	59.60	0.48	54.40	<u>22.24</u>	<b>6.42 (5.7×)</b>	<b>33.86 (44%↓)</b>

\* Refer our paper for results of LLaMA3.1-8B-Instruct-128K and Qwen2.5-7B-128K

# Comparisons on Long Context Modeling

## Comparisons on Long-document QA:

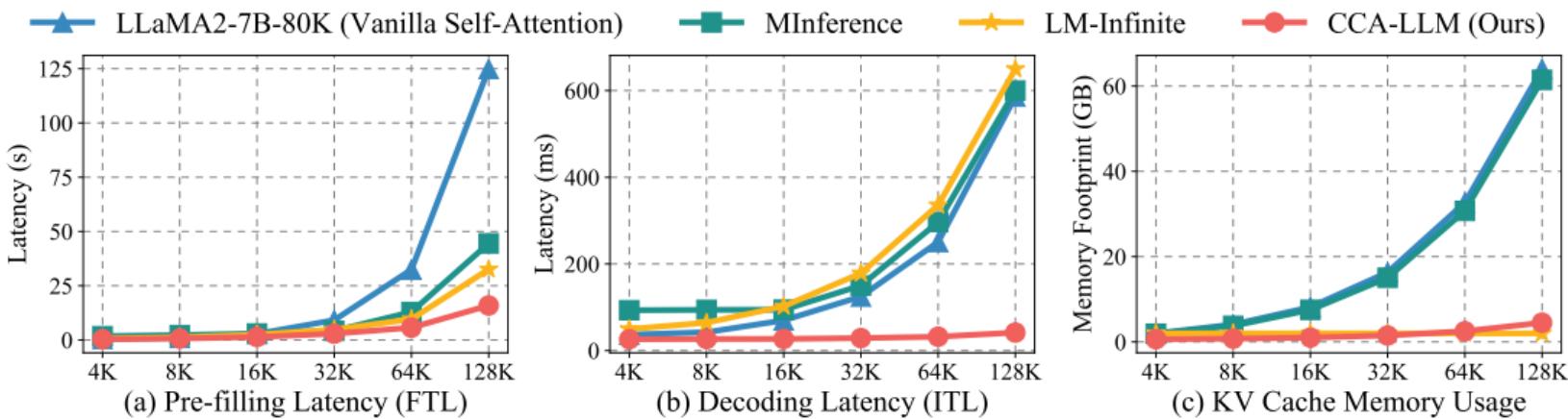
- CCA-LLM outperforms baselines on short sequences (e.g., 4K) by preserving both local and global dependencies without token discard.
- For extreme lengths (64K–128K), CCA-LLM delivers up to 7.9× faster inference and higher accuracy than vanilla self-attention, thanks to adaptive context compression.
- Our method consistently excels across diverse context scales (4K–128K), avoiding performance degradation by dynamically balancing core/redundant context.

Methods	4K	8K	16K	32K	64K	128K	Avg.	FTL (s)
<i>LLaMA2-7B-80K (Vanilla Self-Attention)</i>	<b>39.4</b>	<b>37.8</b>	<b>37.6</b>	<b>36.2</b>	<u>34.6</u>	<u>30.3</u>	<b>36.0</b>	124.85
• StreamingLLM	33.6	26.0	32.2	30.6	27.4	25.1	29.2	34.74 (3.6×)
• LM-Infinite	31.6	25.6	32.4	32.2	28.2	26.3	29.4	32.57 (3.8×)
• MInference	39.0	32.4	<u>37.4</u>	<u>36.0</u>	32.3	28.9	34.3	20.18 (6.2×)
• CCA-LLM (Ours)	<u>39.3</u>	<u>33.2</u>	35.4	31.4	<b>35.3</b>	<b>32.0</b>	<u>34.4</u>	<b>15.89 (7.9×)</b>

# Comparisons on Computational Efficiency

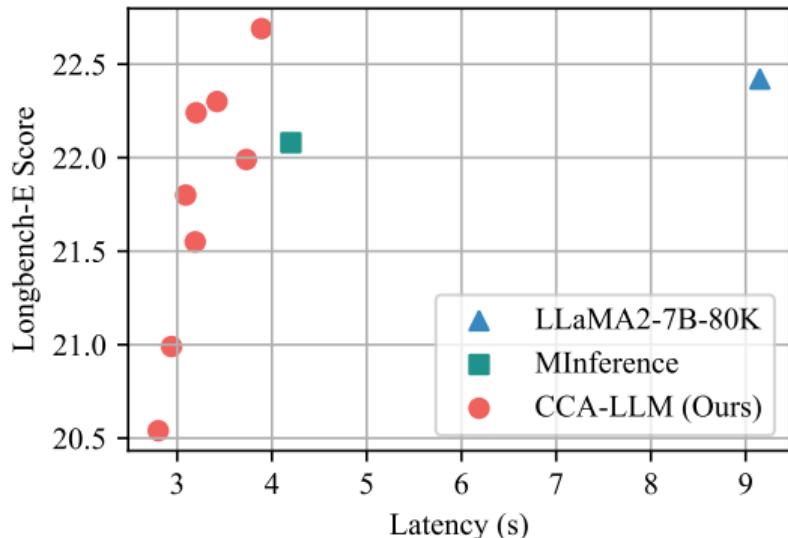
We compare our CCA-LLM with existing methods in terms of inference latency and memory footprint during inference on a single NVIDIA A800 GPU.

- CCA-LLM achieves **7.9 $\times$**  faster pre-filling speed than full self-attention (15.89s vs. 124.85s at 128K) and accelerates both pre-filling and decoding phrase.
- With **93%** KV cache usage reducing (4.5GB vs. 64GB at 128K), CCA-LLM slashes GPU memory demands while maintaining end-to-end efficiency.



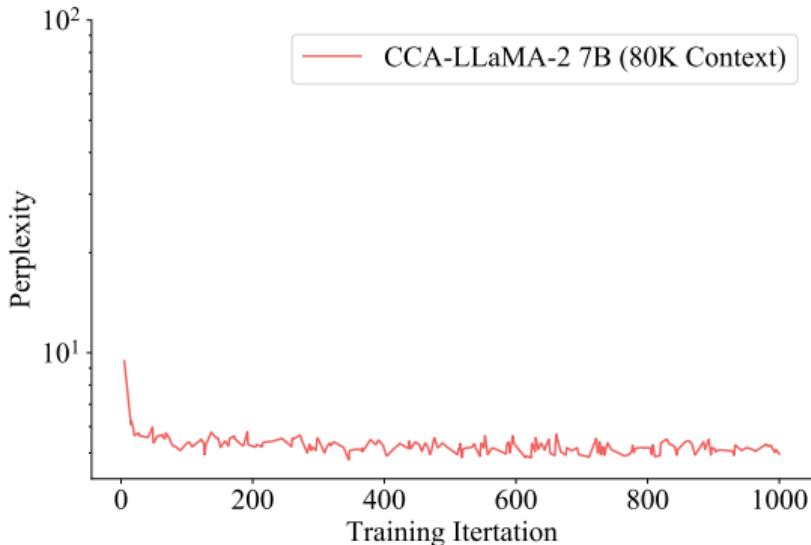
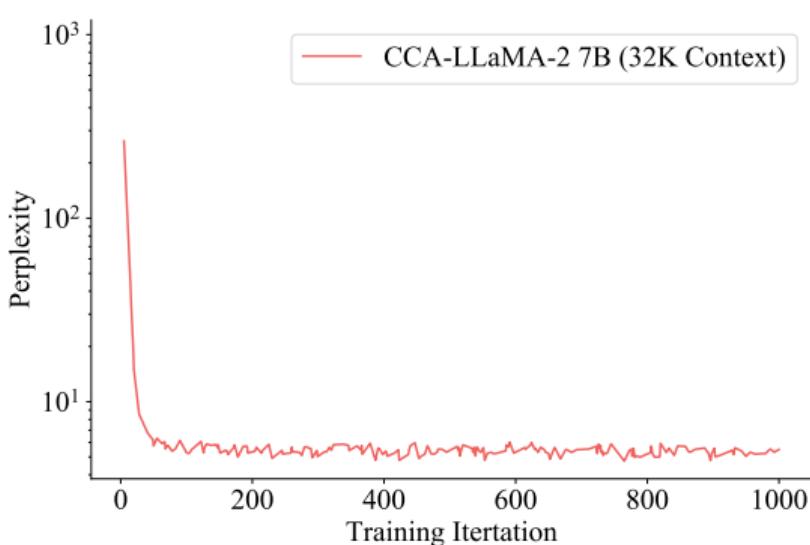
# Demonstration of Inference Flexibility

- With **no retraining needed**, CCA offers a tunable tradeoff between accuracy and speed, enabled by its context-aware compression and scale-invariant local attention.
- We can **dynamically adapt group/window sizes** to maintain high throughput during peak traffic, eliminating server bottlenecks while retaining performance.



# Training Convergence Curve

- Our method enables both LLaMA2-7B-32K and 80K models to achieve rapid convergence within 100 steps while maintaining training stability over 1,000 steps.
- The consistent convergence performance demonstrates the effectiveness of CCA-Attention as a readily integrable component for existing LLMs, regardless of their initial context window sizes.



# Effect of Different Updating Strategies

We can fine-tune the models by two updating strategies: 1) updating all the parameters (full finetuning) and 2) only updating  $\mathbf{W}^Q$ ,  $\mathbf{W}^K$ ,  $\mathbf{W}^V$  (partial finetuning).

- Our method achieves optimal performance when fully fine-tuning all parameters, enabling complete adaptation to the proposed attention mechanism.
- Partial fine-tuning offers a resource-efficient solution with competitive accuracy, ideal for scenarios with computational constraints.

Strategies	Single-Doc. QA	Multi-Doc. QA	Sum.	FS. Learning	Synthetic	Code	Avg.
Partial Finetuning	5.39	3.62	9.21	60.41	1.34	51.77	21.96
Full Finetuning	5.62	4.34	8.99	59.60	0.48	54.40	<b>22.24</b>

# Outline

1. Background

2. Core Context Aware Attention

3. Experimental Results

4. Conclusion

# Conclusion

- Extremely long contexts in LLMs introduce redundant computations while sparse attention methods compromise token connectivity. Our Core Context Aware Attention addresses this by dynamically compress context into core tokens.
- CCA combines (1) a globality-aware pooling module to compress context into adaptive core tokens and (2) a locality-preserving module for fine-grained context – achieving linear complexity without retraining.
- CCA enables  $7.9\times$  faster inference on 128K contexts while outperforming baselines, offering a plug-and-play solution for efficient long-context modeling in LLMs.

**Thanks for Your Attention**

**Q&A**