

# ATTNCACHE: ACCELERATING SELF-ATTENTION INFERENCE FOR LLM PREFILL VIA ATTENTION CACHE

Dinghong Song<sup>1</sup> Yuan Feng<sup>1</sup> Yiwei Wang<sup>1</sup> Shangye Chen<sup>1</sup> Cyril Guyot<sup>2</sup> Filip Blagojevic<sup>2</sup> Hyeran Jeon<sup>1</sup>  
Pengfei Su<sup>1</sup> Dong Li<sup>1</sup>

## ABSTRACT

Large Language Models (LLMs) are widely used in generative applications such as chatting, code generation, and reasoning. However, many real-world workloads—such as classification, question answering, recommendation, and text embedding—rely solely on the prefill stage of inference, where the model encodes input sequences without performing autoregressive decoding. In these prefill-only scenarios, the self-attention computation becomes the primary performance bottleneck due to its quadratic complexity with respect to sequence length. In this paper, we observe that semantically different sentences often produce similar attention maps across layers and heads. Building on this insight, we propose AttnCache, a framework that accelerates the prefill stage of LLM inference by retrieving and reusing similar attention maps. Based on an attention map memoization database, AttnCache employs efficient caching and similarity search techniques to identify and reuse pre-cached attention maps during inference, thereby reducing the computational overhead of self-attention. Experimental results show that AttnCache achieves an average of  $1.2\times$  end-to-end and  $2\times$  attention speedup on CPU, and  $1.6\times$  end-to-end and  $3\times$  attention speedup on GPU, with negligible accuracy degradation. AttnCache is available at: <https://github.com/dinghongsong/AttnCache>.

## 1 INTRODUCTION

Large language models (LLMs) are extensively used in generative tasks, including chatting (e.g., ChatGPT (OpenAI, 2023), Deepseek (DeepSeek, 2025), Claude (Anthropic, 2025)), code generation (e.g., GitHub Copilot (GitHub, 2025), Trae (ByteDance, 2025), Cursor (Cursor, 2025)). Each input prompt first undergoes a prefill phase to encode the context and generate the initial output token, followed by a decoding stage that autoregressively generates subsequent tokens step by step.

Nevertheless, there are also many **prefill-only applications** of LLMs (Du et al., 2025), such as classification (Wang et al., 2018; Gholamian et al., 2024; Vajjala & Shimangaud, 2025), question answering (Talmor et al., 2019; Hendrycks et al., 2020; Phan et al., 2025), recommendation (Wang et al., 2023b; Wu et al., 2024a; Firooz et al., 2025), and data labeling (He et al., 2023; Lan et al., 2024; Zhang et al., 2023a). These workloads use only the embeddings from the final layer, either to produce a single token or as input for downstream tasks. This process does not involve the decoding stage or require generating multiple tokens, so the KV cache is not reused for extended decoding. Therefore, storing the KV cache is unnecessary, and only the prefill stage of LLM inference needs to be executed. For example, in a question answering application, the input prompt could be “What is the capital of France? A. Berlin, B. London, C.

Paris, D. Rome. Your answer is:”, and the LLM only needs to generate a single answer token (i.e., A, B, C, or D).

Furthermore, LLMs can also serve as text encoders to extract general-purpose sentence embeddings, excelling in text representation tasks (Lee et al., 2024b; BehnamGhader et al., 2024; Lee et al., 2024a; Li et al., 2024). In such tasks, only the prefill stage is used to encode the input sentences, with either the hidden states of the last token (Wang et al., 2023a; Lei et al., 2024) or the pooling of all token representations (Li & Zhou, 2024; Lei et al., 2025) as the sentence embedding, without involving decoding stage of LLM inference.

Central to the prefill stage of LLM inference is the self-attention mechanism, which enables LLMs to capture dependencies and relationships across different positions within a sequence. Attention maps, computed as the product of Query (Q) and the transpose of Key (K), encode the relevance of each position to others. However, the quadratic time complexity of this computation with respect to sequence length poses a significant performance bottleneck.

In this paper, we observe that semantically different input sentences can exhibit high similarity in their attention maps across layers or heads during inference. By pre-storing (or caching) these similar attention maps in a vector database (called attention map database) or other memory systems

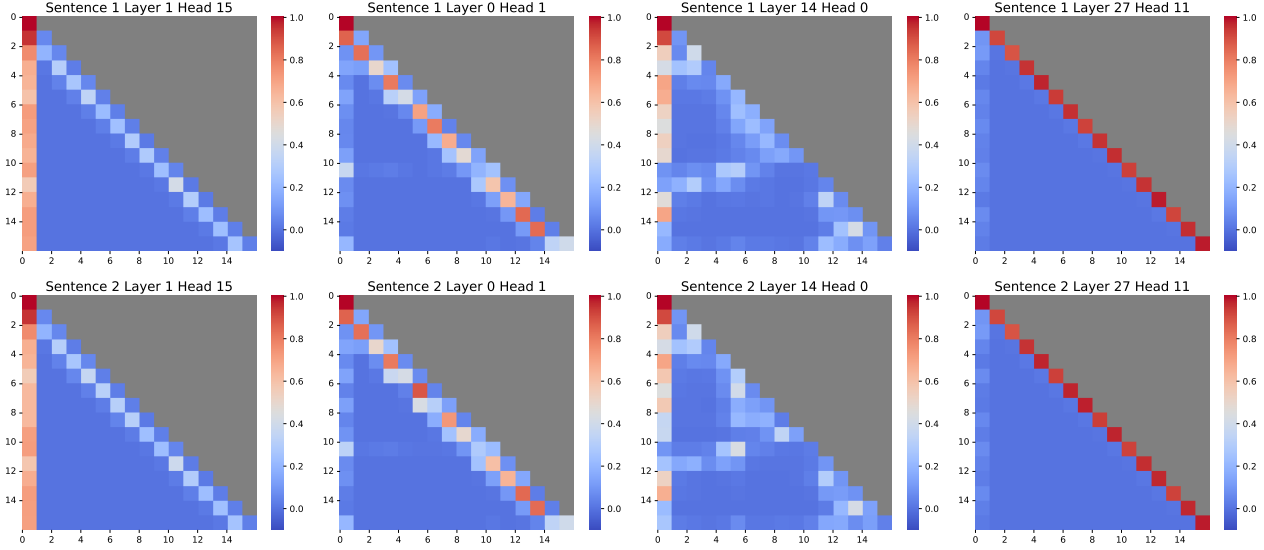


Figure 1. Visualization of the attention maps in Llama-3.2-3B over two sentences, each with a length of 32. Sentence 1 is “This sentence: ‘you should never do it.’ means in one word:”. Sentence 2 is “This sentence: ‘how do you do that?’ means in one word:”. The plots reveals that although Sentence 1 and Sentence 2 have different meanings, their attention maps at different layers and different heads are similar.

such as SRAM, DRAM or HBM, we can retrieve and reuse them to reduce self-attention computation. For example, consider the two sentences shown in Figure 1. Sentence 1 is “This sentence: ‘you should never do it.’ means in one word:”. Sentence 2 is “This sentence: ‘how do you do that?’ means in one word:”. Although the two sentences have different semantics, their attention maps at different layers and different heads are similar, which indicates they have similar relevance at each token position. Consequently, the attention maps computed for Sentence 1 can be reused for Sentence 2. Building on this interesting insight, we introduce *AttnCache*, a framework designed to accelerate self-attention computation. Given that reusing attention maps eliminates the need for storing and computing key and query states in the KV cache, *AttnCache* primarily focuses on accelerating the prefill stage of LLM inference, rather than the decoding stage, which necessitates the storage of past key and value states.

Implementing *AttnCache* presents two key challenges. The first challenge lies in finding an effective data representation. Both the representations of input sentences and attention maps in LLMs are high-dimensional tensors. Therefore, it is practically infeasible to find similar attention maps by directly comparing the representations of input sentences. Instead, we design a lightweight embedding neural network to represent attention maps efficiently. This network must be computationally inexpensive so that its overhead, combined with the search in the attention map database, remains lower than the cost of self-attention computation.

The second challenge is the high cost of memory accesses when storing and fetching pre-populated attention maps. A large attention map database improves search hit rates but leads to sparse memory accesses, as accesses to attention maps exhibit poor spatial and temporal locality. Additionally, modern deep learning frameworks like PyTorch require tensors to be placed in consecutive memory addresses to enable vectorized data accesses for Single Instruction Multiple Data (SIMD) operations. Therefore, once a tensor is fetched from the pre-populated database, it must be copied to a consecutive memory buffer before being loaded to the processor, incurring two memory reads and one write per fetch. To reduce memory access overhead, we first store all attention maps of a layer as a single file object, and arrange the attention maps of neighboring layers continuously in the database, enhancing spatial and temporal locality. Then, *AttnCache* eliminates expensive tensor copying through memory mapping between a consecutive virtual-memory space and scattered physical addresses of individual tensors.

When the memory footprint exceeds GPU capacity, we leverage CPU to demonstrate our approach in *AttnCache*. Generating LLM embeddings for large-scale text corpora often surpasses GPU memory limits. In applications such as recommendation systems, which involve processing billions of text chunks, throughput across many concurrent instances is more critical than per-instance latency. For such workloads, CPUs offer better efficiency in time, energy and cost. Therefore, it is meaningful and valuable to evaluate the performance of *AttnCache* on CPUs. Our evaluation shows that *AttnCache* achieves an average of 1.2× end-to-end and

2× attention speedup on CPU, and 1.6× end-to-end and 3× attention speedup on GPU, with negligible accuracy loss.

## 2 RELATED WORK

**Sentence Embedding.** Sentence embeddings encode the semantic information of sentences into high-dimensional vector representations. Prior works (Li & Zhou, 2024; Muennighoff et al., 2024; Ni et al., 2021) have demonstrated the capability of LLMs to generate high-quality sentence embeddings. Recent studies (Zhuang et al., 2024; Qin et al., 2023; Zhang et al., 2024a) have explored converting LLMs into sentence encoders without additional training. To enhance embedding quality, prompt-based techniques have gained traction. MetaEOL (Lei et al., 2024) uses multitask prompts to generate general-purpose embeddings. The research by (Jiang et al., 2023) illustrates how to extract a sentence embedding by prompting LLMs with the instruction “This sentence: ‘[text]’ means in one word:”. In this work, we leverage LLMs to generate sentence embeddings without fine-tuning.

**LLM inference acceleration.** Most KV cache optimization approaches (Beltagy et al., 2020; Zhang et al., 2023b; Oren et al., 2024) focus on accelerating the LLM decoding phase by reducing redundancy in Key and Value matrices. StreamingLLM (Xiao et al., 2023) identifies “attention sinks” and keeps initial and recent tokens’ KV to anchor attention computation, while FastGen (Ge et al., 2023) prunes tokens during decoding by profiling attention heads. However, these approaches do not reduce the prefill costs. In contrast, several recent efforts (Wu et al., 2024b; Jiang et al., 2024; Tang et al., 2024) center on optimizing the LLM prefill phase, benefiting tasks like sentence embedding generation. PromptCache (Gim et al., 2024) and ChunkAttention (Ye et al., 2024) reduce time-to-first-token latency by sharing KV tensors of common prompt prefixes. Other acceleration approaches (He et al., 2024; Men et al., 2024; Song et al., 2024; Zhang et al., 2024b) focus on removing the redundant attention or transformer layer in LLMs. Token pruning (Ham et al., 2020; Wang et al., 2021) reduces computation by excluding less important tokens from the input, while layer-wise reuse (Ying et al., 2021; Xiao et al., 2019; Bhojanapalli et al., 2021) reduces computation by sharing attention maps calculated in prior layers in multiple subsequent layers. These approaches complement AttnCache and can further enhance memory efficiency.

**Reuse mechanism in Neural Networks.** The reuse mechanism exploits the inherent redundancy in neural networks to enhance efficiency. Prior works (Ning et al., 2019; Ning & Shen, 2019; Wu et al., 2022; Köpüklü et al., 2019) have explored reusing similar computation results to improve performance. Silfa et al. (2019) accelerate RNN training by

reusing neuron outputs. Studies (Bhojanapalli et al., 2021; Xiao et al., 2019) have shown that transformer attention maps (Vaswani et al., 2017) exhibit similar distributions across adjacent layers. Many prior efforts (Hunter et al., 2023; Xiao et al., 2019; Bhojanapalli et al., 2021; Ying et al., 2021; Liao & Vargas, 2024) focus on sharing computed attention weights across multiple layers for the same input sequence. However, this approach may introduce dissimilar attention maps, which can degrade performance. In contrast, our work efficiently reuses similar attention maps across different sequences, overcoming the limitations of intra-sequence reuse.

## 3 METHODOLOGY

As shown in Figure 2, given an input sentence, AttnCache embeds it into a feature vector using a lightweight neural network (feature projector). The feature vector is used to retrieve the index of the attention maps that have the highest similarity to the input sentence. Then, the search engine uses the index to fetch the corresponding attention maps from the attention map database. The fetched attention maps are used in the self-attention computation during online inference, while the prefill stage in LLM inference is utilized to generate the sentence embedding.

### 3.1 Search Engine

As illustrated in Figure 1, two sentences with completely different semantics can produce highly similar attention maps. Because input sentences are represented as high-dimensional hidden states, directly comparing these representations provides little insight into the similarity of their corresponding attention patterns. To overcome this limitation, AttnCache uses the feature vector of input hidden states, which is embedded by the feature projector. By searching for similar feature vectors, AttnCache can efficiently retrieve input embeddings that yield similar attention maps.

**Feature Projector.** We use two layers of Multi-Layer Perceptron (MLP) as the feature projector, which maps the input embedding to a feature vector with lower dimension size. The network structure of Feature Projector is important to the accuracy and efficiency of the search process. Compared with other embedding models, such as convolutional neural network or transformer, MLP is lightweight with less computational complexity and shorter inference time. Training the feature projector is challenging due to a lack of labeled data. Deciding the similarity between input embeddings and labeling them as similar or not is prohibitively expensive. We use the Siamese network (Koch et al., 2015), which contains two identical feature projectors and shares the same weights, as shown in Figure 3.

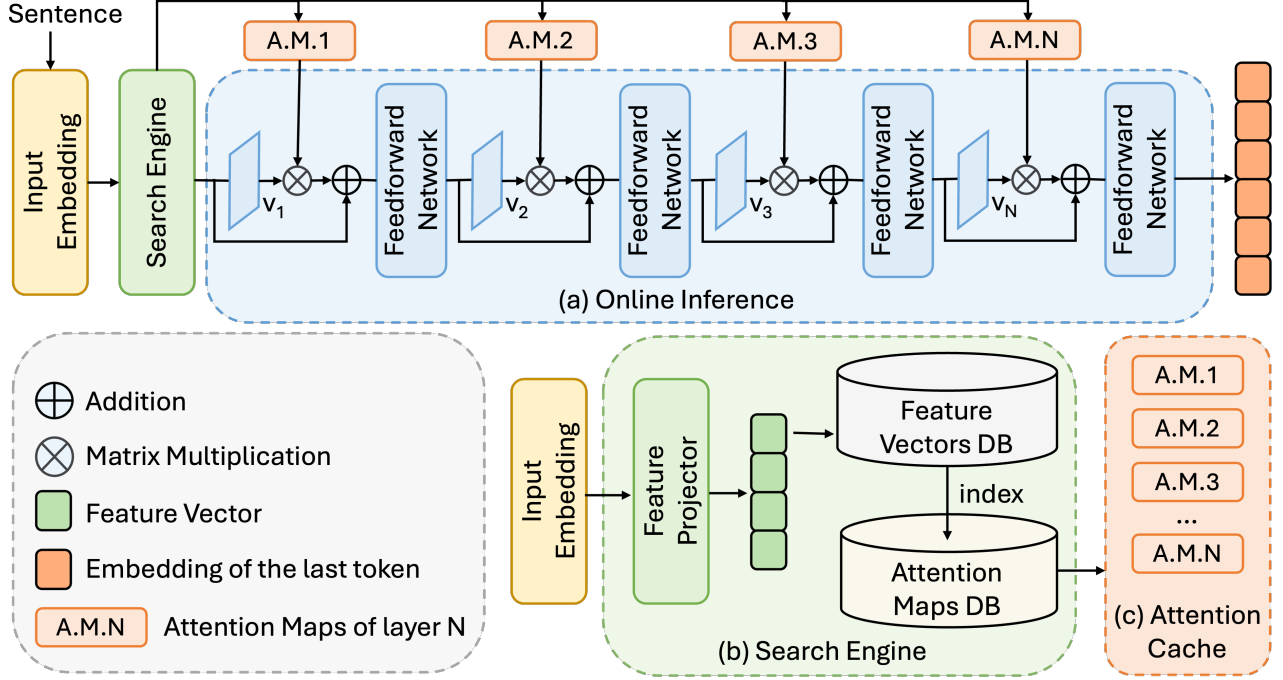


Figure 2. AttnCache overview. The search engine will identify the index of the sentence that produces the most similar attention maps based on the feature vector of the current input sentence and prefetch attention maps for each layer from the attention map database using the index. These fetched attention maps are stored in the attention cache and reused for the matrix multiplication calculation with value projection during the self-attention computation.

#### Algorithm 1 Search Engine

```

1: Input: Sentence  $S$ , Threshold  $\theta$ ;
2: Output: Attention Cache  $attn\_cache$ ,
   Input embedding  $h$ ;
3: Function SEARCH_ENGINE( $S, \theta$ )
4:    $h \leftarrow \text{encode}(S)$ 
5:    $f \leftarrow \text{feature\_projector}(h)$ 
6:    $(idx, sims) \leftarrow \text{VecDB.search}(f)$ 
7:    $attn\_cache \leftarrow []$ 
8:   if  $sims \geq \theta$  then
9:      $n \leftarrow \text{num\_layers}$ 
10:     $ams \leftarrow \text{AttnMapsDB.get}(idx, n)$ 
11:     $attn\_cache.append(ams)$ 
12:   end if
13: return ( $attn\_cache, h$ )

```

During each training iteration, two input embeddings are used as input to the two identical feature projectors in the Siamese network. After getting the feature vectors, the Euclidean distance (i.e. L2-norm) is calculated as follows.

$$\hat{y} = \|f_{\mathbf{W}}(\mathbf{X}_1) - f_{\mathbf{W}}(\mathbf{X}_2)\|_2 \quad (1)$$

where  $\mathbf{X}$  is the input embedding,  $f_{\mathbf{W}}$  is the feature projector, and  $\|\cdot\|_2$  is the L2 norm. Besides, we measure the similarity

score using the attention maps and the sequence length of tokens, which associate with the two input embeddings. We use the metric as the labels for training the feature projector based on the average distance of heads, which is defined as follows.

$$y = \frac{1}{n} \times \alpha \sum_{p=1}^n \frac{1}{2} \|\mathbf{A}_1[p, :] - \mathbf{A}_2[p, :]\|_2 + \|s_1 - s_2\|_1 \quad (2)$$

where  $\mathbf{A}$  denotes the attention map,  $n$  indicates the number of head,  $\mathbf{A}[p, :]$  is the  $p^{th}$  row of the attention map,  $\|\cdot\|_1$  is the L1 norm,  $s$  denotes the length of input token sequence, and  $\alpha$  is the hyperparameter to control the relative importance of the similarity of the attention maps and the token length. In addition to the inherent similarity of the attention maps, the token sequence also plays an important role in determining whether two attention maps are similar. When the token sequences of two attention maps are very different in length, even if the attention maps are similar, they cannot be used directly in AttnCache, otherwise it may cause a large inference error. The final loss function of the feature projector is defined as follows.

$$L = \begin{cases} 0.5(\hat{y} - y)^2 & \text{if } |\hat{y} - y| < 1 \\ |\hat{y} - y| - 0.5 & \text{if } |\hat{y} - y| \geq 1 \end{cases} \quad (3)$$

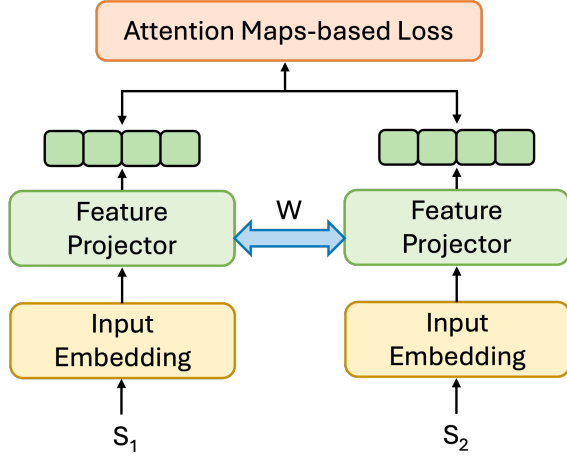


Figure 3. The training of the feature projector. The feature projector maps input embedding of a sentence  $S$  into a feature vector. Then we train the feature projector using the attention maps-based loss function.

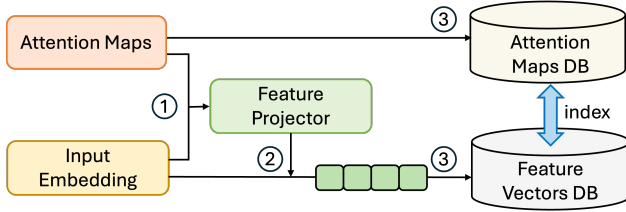


Figure 4. Databases building include three steps. 1. Train the feature projector with input embeddings and attention maps; 2. Embed the input embeddings to feature vectors; 3. Store the feature vectors and attention maps to their respective databases. Both databases share the same index.

We use Smooth L1 Loss (Girshick, 2015) as the loss function, which is able to balance the effects of outliers. The training process iteratively updates the parameters of the feature projector to minimize the loss function.

**Databases.** To minimize the costly search for attention maps, we construct an indexed database, where feature vectors are stored and indexed for fast search. In essence, the feature vector database is a key-value store where the key and value are the feature vector and its index. Figure 4 illustrates the process of building the databases. The attention maps associated with the feature vectors are stored in the attention map database, where the key is the index and value is the attention map. Both databases have the same index.

Algorithm 1 illustrates the process of finding the most similar attention maps, referred to as `search_engine`. The input sentence is embedded by input embedding (Line 2). The input embedding includes tokenization of the sentence, position encoding, and layer normalization. Then the result

is mapped into a feature vector with lower dimension (Line 3). The feature vector is used for querying in the feature vector database. After the query, the indices that have the highest similarity to the feature vector are returned (Line 4). When the similarity is not less than the threshold  $\theta$ , the corresponding index  $idx$  is used to fetch attention maps from the attention map database.

The retrieved index  $idx$  corresponds to a sentence  $S$  whose attention maps are similar to those of the current input sentence. The corresponding attention maps are fetched for all layers of the LLM and stored in a contiguous memory region, referred to as the attention cache. Specifically, these attention maps are used in the matrix multiplication calculation with value projection in `online_inference`. All layers of the attention maps are fetched for the computation of self-attention before the LLM inference starts.

### 3.2 Online Inference

Algorithm 2 illustrates online inference with AttnCache. In the attention block of each layer, the value projection is computed. If similar attention maps are found, the attention output can be obtained by multiplying the attention maps by the value  $v$ . Thus, finding similar attention maps and reusing them in the self-attention calculation leads to performance benefits.

However, AttnCache cannot always find similar attention maps. For those hidden states with low similarities, the attention maps must be calculated at each layer during the inference, which means the query, key, rotary positional encoding, and softmax normalization must be computed. In this regard, AttnCache does not bring benefit in inference speed, and instead degrades performance due to its search overhead. However, given a batch of inferences, as long as the success rate of retrieving for all inferences is high, the overall inference is still accelerated.

## 4 APPLICABILITY OF ATTNCACHE

AttnCache is well-suited for tasks that rely solely on the prefill stage inference of LLMs. By reusing similar attention maps, AttnCache effectively reduces three time-consuming matrix multiplications at each layer, i.e.,

$$\begin{aligned} \mathbf{Q} &= \mathbf{X}\mathbf{W}_Q, & \mathbf{Q} &\in \mathbb{R}^{L \times d_k} \\ \mathbf{K} &= \mathbf{X}\mathbf{W}_K, & \mathbf{K} &\in \mathbb{R}^{L \times d_k} \\ \text{AttnMaps} &= \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right), & \text{AttnMaps} &\in \mathbb{R}^{L \times L} \end{aligned} \quad (4)$$

where  $L$  denotes the input sentence length in tokens,  $\mathbf{X} \in \mathbb{R}^{L \times d}$ ,  $\mathbf{W}_Q, \mathbf{W}_K \in \mathbb{R}^{d \times d_k}$ . For simplicity, we omit the notation for batch size and the number of attention heads. When reusing attention maps,  $\mathbf{Q}$  and  $\mathbf{K}$  are neither com-



---

**Algorithm 2** Online Inference

---

```
1: Input: Attention Cache  $attn\_cache$ ,  
   Input embedding  $h$ ;  
2: Output: Hidden states of last layer  $h$ ;  
3: Function ONLINE_INFERENCE( $attn\_cache, h$ )  
4:   for  $l$  in range( $num\_layers$ ) do  
5:      $residual \leftarrow h$   
6:      $v \leftarrow v\_projection(h)$   
7:     if  $attn\_cache$  is not NULL then  
8:        $attn\_map \leftarrow attn\_cache[l]$   
9:        $h \leftarrow mat\_mul(attn\_map, v)$   
10:    else  
11:       $q \leftarrow q\_projection(h)$   
12:       $k \leftarrow k\_projection(h)$   
13:       $(q, k) \leftarrow rotary\_pos\_emb(q, k)$   
14:       $attn\_map \leftarrow softmax(q, k)$   
15:       $h \leftarrow mat\_mul(attn\_map, v)$   
16:    end if  
17:     $h \leftarrow residual + h$   
18:     $h \leftarrow h + feed\_forward(h)$   
19:  end for  
20: return  $h$ 
```

---

puted nor stored. However, during LLM decoding inference, the keys stored in the KV cache are required for each step of autoregressive token generation. Therefore, although AttnCache can accelerate inference in the prefill stage, its inability to compute and store  $\mathbf{K}$  makes it unsuitable for decoding scenarios. Another limitation is the mismatch in attention map dimensions between the prefill and decoding stages. In the decoding phase, since only one token is generated at a time, the corresponding attention map has a shape of  $1 \times (L + t)$  rather than  $L \times L$ , i.e.,

$$\begin{aligned} \mathbf{Q}_t &= \mathbf{X}_t \mathbf{W}_Q, & \mathbf{Q}_t &\in \mathbb{R}^{1 \times d_k} \\ \mathbf{K}_{\leq t} &= \text{cached Keys} & \mathbf{K}_{\leq t} &\in \mathbb{R}^{(L+t) \times d_k} \\ \text{AttnMaps} &= \text{softmax} \left( \frac{\mathbf{Q}_t \mathbf{K}_{\leq t}^\top}{\sqrt{d_k}} \right), & \text{AttnMaps} &\in \mathbb{R}^{1 \times (L+t)} \end{aligned} \quad (5)$$

AttnCache currently can only accelerate computation during the prefill stage, as the reused AttnMaps have a shape of  $L \times L$  and therefore cannot be applied in the decoding stage. Investigating how to reuse AttnMaps with a shape of  $1 \times (L + t)$  could be a promising research direction to reduce KV cache storage and the computational cost of self-attention during LLM decoding.

## 5 EXPERIMENTS

### 5.1 Datasets

We take three representative datasets, including Semantic Textual Similarity (STS) (Muennighoff et al., 2022), Stan-

ford Sentiment Treebank v2 (SST-2) (Wang et al., 2018), and Massive Multitask Language Understanding (MMLU) (Hendrycks et al., 2020). STS datasets contain STS 12-16, STS-B and SICK-R (Lei et al., 2024). The semantic similarity of each sentence pair is annotated with a score of 0-5. We use the Spearman correlation score (Lei et al., 2024) between the ground-truth similarity scores and the predicted similarity scores as the evaluation metric. SST-2 is a binary sentiment classification dataset derived from movie reviews. Each sentence is labeled as either positive or negative. MMLU is a multi-choice dataset containing questions from 57 diverse subjects, including math, computer science, engineering, physics, and more. For SST-2 and MMLU, we use accuracy as the evaluation metric. These datasets and their outputs and labels are summarized in the Table 1.

Table 1. Summary of STS-B, SST-2, and MMLU with task types and labels.

Task	Type	Output	Label
STS	Semantic similarity estimation	Continuous score	Similarity score (0–5)
SST-2	Sentiment classification	Categorical label	Positive / Negative (1 / 0)
MMLU	Multitask multiple-choice QA	Categorical label	Multiple choice (A/B/C/D)

### 5.2 Models

We conduct experiments using three representative open-source models, including Llama-2-7B (Touvron et al., 2023), Llama-3-8B (Touvron et al., 2024), and Mistral-7B (Jiang et al., 2023). All these models are run with weights stored in full-precision (fp32) floating-point format, and evaluation is conducted using the SentEval toolkit (Conneau & Kiela, 2018), measuring performance on CPU. For single-GPU scenarios, we evaluate the full-precision Llama-3-3B (Touvron et al., 2024) and the 4-bit quantized Llama-3-8B, Deepseek-MoE-16B (DeepSeek, 2024), and Qwen1.5-MoE-A2.7B (Team, 2024) on MMLU datasets. To assess performance under varying context lengths, we evaluate Llama-3-3B on MMLU under different input n-shot settings. To demonstrate the generality of AttnCache, in addition to Transformer decoder-based LLMs, we also evaluate its performance on Transformer encoder-based models, such as BERT base (Devlin et al., 2018), RoBERTa (Liu et al., 2019), and DeBERTa (He et al., 2020).

### 5.3 Experimental Setting

We evaluate AttnCache on a server equipped with two sockets, each with 24-core Intel(R) Xeon(R) Silver 4410Y processors. The platform provides 512 GB DRAM and a 14

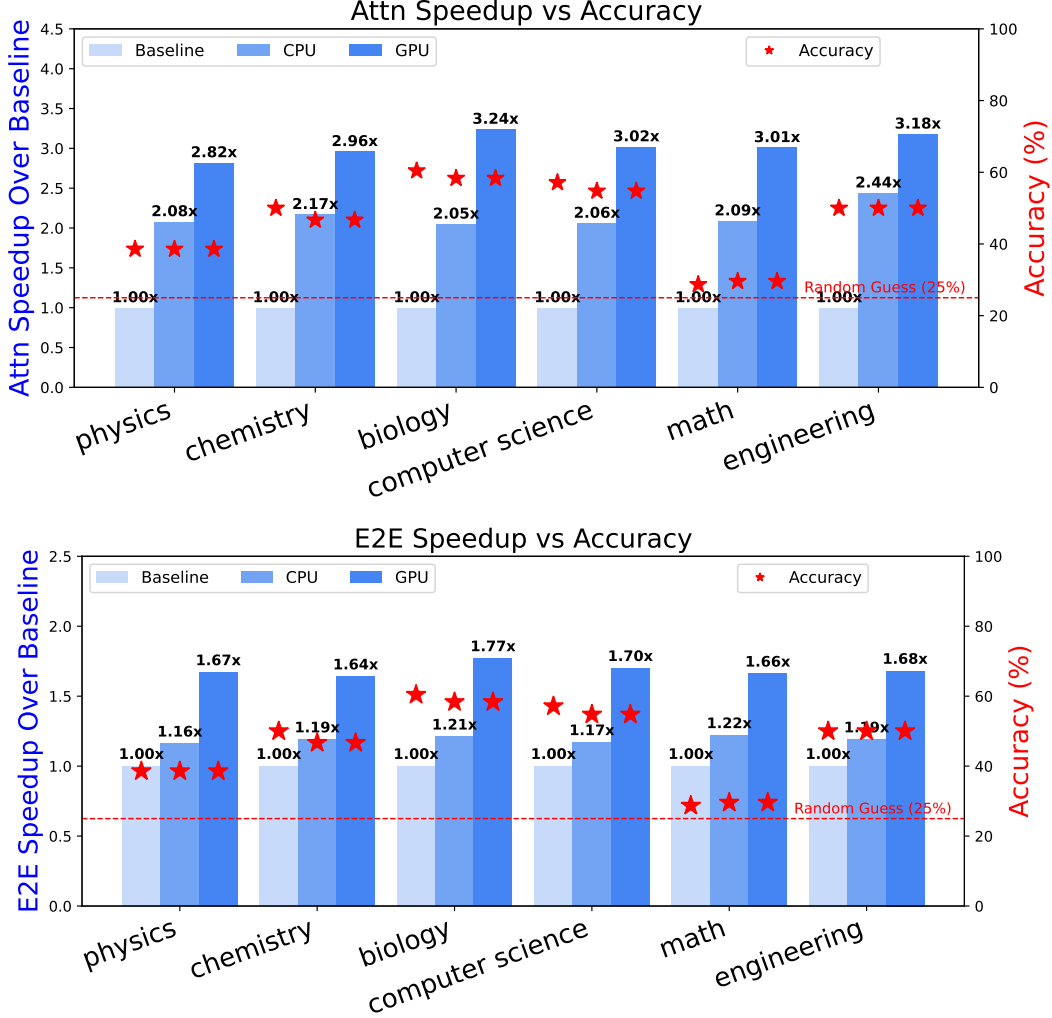


Figure 5. Attention and end-to-end speedup of Llama-3-3B with AttnCache on MMLU STEM. By selecting an appropriate threshold, AttnCache achieves an average of 1.2 $\times$  end-to-end and 2 $\times$  attention speedup on CPU, and 1.6 $\times$  end-to-end and 3 $\times$  attention speedup on GPU, while incurring only a negligible drop in accuracy.

TB hard disk drive (HDD), with the DRAM used to store the attention map database and feature vector database. In addition, the platform includes an NVIDIA A100 GPU with 80 GB of HBM. Since AttnCache performs inference only once and generates sentence embeddings solely during the prefill phase—without creating or storing KV cache—a 80 GB GPU is sufficient for small models, such as Llama-3-3B. If the GPU runs out of memory, online inference can only be performed on the CPU. To build the feature vector database, we use Faiss (Johnson et al., 2019), a vector database enabling efficient similarity search by the Hierarchical Navigable Small Worlds algorithm (Malkov & Yashunin, 2018). We use the standard LLM inference as the baseline, named *full model*.

#### 5.4 Details of Implementation

For each task, we collect the input hidden states and their corresponding attention maps at each layer, which are used for training the feature projector and building databases; then we randomly select 1K samples that are not involved in training to measure AttnCache. The dimensions of the feature vector and batch size are set to 128 and 64, respectively. To maintain high inference accuracy, we set the similarity threshold  $\theta$  to 0.99, and set  $\alpha$ , which is used to train the feature projectors (see Equation 2), to 0.2. For efficient similarity search, we construct the feature vector database using Faiss (Johnson et al., 2019). Faiss is highly efficient for similarity search. For example, our evaluation shows that with Faiss, searching 100K vectors with a vector-dimension size of 128 takes less than 0.5 ms, which yields 360 $\times$  and 10 $\times$  speedups over self-attention computation and embedding

Table 2. Spearman correlation score (in %) across 7 STS tasks.

Llama-2-7B										
Method	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	Avg. (↑)	E2E Speedup (↑)	$\gamma$ (↓)
Full Model	60.88	73.93	58.30	70.27	75.46	73.89	67.44	68.60	1.00×	–
SAN	5.02	42.63	19.84	43.49	44.70	18.01	38.71	30.34	1.45×	0.85
LazyFormer	23.79	34.88	27.80	35.93	44.04	32.50	42.45	34.48	1.39×	0.87
AttnCache-f	22.06	67.75	31.52	61.15	53.89	53.97	62.40	50.39	1.14×	1.30
<b>AttnCache</b>	60.59	73.46	57.97	69.01	75.38	72.02	65.85	67.75	1.19×	0.04
Llama-3-8B										
Method	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	Avg. (↑)	E2E Speedup (↑)	$\gamma$ (↓)
Full Model	61.57	76.41	63.23	75.27	80.41	75.84	70.45	71.88	1.00×	–
SAN	27.61	53.81	37.18	57.20	57.43	39.46	54.98	46.81	1.49×	0.51
LazyFormer	27.25	60.37	36.21	53.85	59.21	40.30	48.24	46.49	1.42×	0.60
AttnCache-f	24.89	51.15	36.19	67.81	61.39	48.05	63.77	50.46	1.16×	1.34
<b>AttnCache</b>	60.82	72.49	60.59	74.67	79.52	72.61	66.68	69.63	1.21×	0.11
Mistral-7B										
Method	STS12	STS13	STS14	STS15	STS16	STS-B	SICK-R	Avg. (↑)	E2E Speedup (↑)	$\gamma$ (↓)
Full Model	63.28	74.89	61.57	75.64	81.89	78.26	69.39	72.13	1.00×	–
SAN	25.04	54.66	35.30	53.11	61.55	39.59	55.45	46.39	1.44×	0.58
LazyFormer	38.90	54.41	38.71	37.18	57.61	42.23	50.66	45.67	1.38×	0.70
AttnCache-f	35.03	55.07	40.28	54.51	50.22	54.75	64.52	50.63	1.15×	1.43
<b>AttnCache</b>	62.66	72.23	61.85	73.32	81.59	74.66	65.89	70.31	1.20×	0.09

generation, respectively. As a result, the search process does not create a performance bottleneck for AttnCache. In addition, we store each layer’s attention maps as a file object in memory. When retrieving attention maps as a batch, the file objects are mapped into a contiguous virtual memory space as a tensor without a memory copy. After self-attention calculation, the file objects are unmapped. When the combined size of the attention map and feature vector databases exceeds the available DRAM on our platform, we evaluate model performance using a hybrid storage setup that spans DRAM and HDD. While this configuration introduces I/O latency during access, it does not affect the correctness or quality of the model outputs. In this case, to evaluate the model inference time on a “virtual” big DRAM system with enough capacity to store attention maps, we use our limited DRAM assuming that the needed attention maps are in the DRAM for measuring time.

## 5.5 Baselines

We use three baselines for evaluation.

**LazyFormer** (Ying et al., 2021) divides all layers of the transformer to multiple subblocks. In each subblock, the attention maps are only computed in the first layer and then used by the remaining layers in the same subblock. Like LazyFormer, we set the number of layers in each sub-block

to 2.

**SAN** (Xiao et al., 2019) shares attention maps across multiple adjacent layers. But different from Lazyformer, SAN does not use a uniform subblock size (i.e., the number of transformer layers in a subblock). The subblock size is dynamically determined based on the similarity of layers in terms of the JS divergence (Menéndez et al., 1997).

**AttnCache-f** is a variant of AttnCache. AttnCache-f applies memoization at the transformer layer level instead of the whole model level (as AttnCache does). In particular, at each layer, AttnCache-f searches the attention map database for similar attention maps, hence applying a fine-grained memoization. Moreover, AttnCache-f does not consider sequence length when training the feature projector, meaning that the  $y$  in Equation 2 does not take into account the computation of  $\|s_1 - s_2\|_1$ .

To quantify the trade-off between speed and performance degradation, we adopt the Speedup Degradation Ratio  $\gamma$  (He et al., 2024) as an evaluation metric.

$$\gamma = \frac{\text{Avg}_{\text{full}} - \text{Avg}_{\text{method}}}{\text{Speedup}_{\text{method}} - \text{Speedup}_{\text{full}}} \quad (6)$$

where  $\text{Avg}_{\text{full}}$  and  $\text{Avg}_{\text{method}}$  are the average performance of LLMs and each method across the seven tasks respec-



Table 3. MMLU Math performance on Llama-3-3B under different n-shot settings.

Type	Context Length		Accuracy		Attn Speedup		E2E Speedup	
	Max Len.	Avg. Len.	Baseline	AttnCache	CPU	GPU	CPU	GPU
0-shot	373	93.85	28.70	29.57	2.09x	3.01x	1.22x	1.66x
1-shot	483	178.86	40.87	38.26	2.04x	2.99x	1.23x	1.62x
2-shot	613	284.79	41.74	42.61	2.11x	2.97x	1.19x	1.64x
3-shot	784	393.12	40.00	39.13	2.13x	3.01x	1.27x	1.62x
4-shot	954	490.28	43.48	41.74	2.07x	3.02x	1.25x	1.63x
5-shot	1019	560.28	47.83	46.96	2.12x	3.00x	1.24x	1.65x

tively, and  $\text{Speedup}_{full}$  and  $\text{Speedup}_{method}$  represent the corresponding speedup respectively. A smaller  $\gamma$  indicates that the method is more efficient.

## 5.6 Main Results

Table 2 summarizes the experimental results on the STS datasets. Across various models (Llama2-7B, Llama3-8B, and Mistral-7B), SAN and LazyFormer both lead to notable performance declines, despite achieving higher speedups. For instance, LazyFormer results in an average 25.39% performance decline (from 71.88% to 46.49%) for Llama-3-8B, with a speedup of  $1.42\times$ , corresponding to a  $\gamma$  of 0.60. We also notice that the inter-sentence methods (i.e. AttnCache-f and AttnCache) exhibit higher performance but lower speedup compared to intra-sentence methods because they only reuse attention maps with high similarity. For example, for Llama-2-7B, AttnCache-f and AttnCache achieve average performance of 50.39% and 67.75% with corresponding speedups of  $1.14\times$  and  $1.19\times$ , while SAN and LazyFormer yield 30.34% and 34.48% performance with speedups of  $1.39\times$  and  $1.45\times$  separately. Moreover, AttnCache maintains near full model performance on various datasets and strikes a better balance between speed and performance, with  $\gamma$  values of 0.04, 0.11 and 0.09 for three LLMs, making it a superior method for the acceleration of self attention.

As shown in Table 4, AttnCache-f performs embedding and vector search at each layer, even when no reusable attention maps are found, increasing latency. In contrast, AttnCache performs this computation only once at the beginning of inference to determine whether to reuse attention maps, eliminating embedding and vector search overhead during subsequent layers. Consequently, AttnCache achieves a higher  $\gamma$  than AttnCache-f. As illustrated in Figure 5, AttnCache speeds up Llama-3-3B on MMLU STEM, achieving up to  $2\times/3\times$  attention and  $1.2\times/1.6\times$  end-to-end speedups on CPU/GPU, with minimal accuracy loss. Similar results are observed under varying context lengths by adjusting the n-shot input settings, as shown in Table 3.

Table 4. Time (ms) breakdown for AttnCache-f and AttnCache in a layer of Llama-3-8B.

Time (ms)	Full Model	AttnCache-f	AttnCache
Embedding	N/A	32	N/A
Vector Searching	N/A	2	N/A
APM Fetching	N/A	18	N/A
Q Computation	73	N/A	N/A
K Computation	41	N/A	N/A
Rotary Pos Encoding	124	N/A	N/A
V Computation	41	41	41
AM Computation	88	N/A	N/A
Other (e.g. AM • V)	115	115	115
<b>Attention</b>	482	208	156
<b>FFN</b>	830	830	830
<b>Total</b>	1312	1038	986

## 5.7 Evaluation on Dense and MoE Language Models.

To fit the LLMs into a single GPU for MoE models, we use bitsandbytes (bitsandbytes, 2025) quantization to reduce the GPU memory footprint of LLMs. Specifically, for Llama-3-8B, Deepseek-MoE-16B, and Qwen1.5-MoE-A2.7B, we apply NF4 (Normal Float 4) quantization (Dettmers et al., 2023). The experimental results are presented in Table 5. AttnCache achieves up to a  $2.43\times$  attention speedup and a  $1.48\times$  end-to-end speedup with only minor accuracy degradation, demonstrating the robustness of the method. These results confirm that AttnCache is both effective and generalizable across Dense and MoE model architectures.

Table 5. MMLU Math performance on Dense and MoE models.

Models	Llama-3-8B	Deepseek-MoE-16B	Qwen1.5-MoE-A2.7B
# Parameters	8.03B	16.40B	14.30B
Model Type	Dense	MoE	MoE
Quant Type	NF4	NF4	NF4
Memory Footprint	7G	10G	9G
Baseline Accuracy	41.74	30.43	35.65
AttnCache Accuracy	40.00	28.70	34.78
Attn Speedup	1.67x	2.43x	2.28x
E2E Speedup	1.48x	1.12x	1.15x

Table 6. Model size and architecture type.

Models	BERT base	RoBERTa	DeBERTa
# Parameters	110M	125M	139M
Model Type	Encoder	Encoder	Encoder

Table 7. Accuracy and speedup under different thresholds on SST-2.

%	Baseline	Conservative	Moderate	Aggressive	Avg. Diff.
BERT base	91.3	91.1	90.2	85.7	-2.3
RoBERTa	94.8	93.2	92.6	90.4	-2.7
DeBERTa	95.0	95.5	95.2	90.5	-1.3
E2E Speedup	N/A	1.10x	1.18x	1.34x	1.21x

Table 8. Integration with model Quantization and Pruning. “w/Quant”, “w/AttnDrop” and “w/BlockDrop” denotes integration with the quantized model, attention pruning and layer pruning respectively.

Llama-3.2-3B					
Method	STS13	STS14	STS15	STS16	Avg.
Full Model	76.56	60.05	74.76	79.30	72.67
AttnCache	74.74	59.95	74.19	77.38	71.57
Quanto	75.27	57.55	74.41	76.96	71.05
w/Quanto	74.25	54.75	74.49	76.92	70.10
AttnDrop	75.33	59.04	69.92	78.37	70.67
w/AttnDrop	73.21	56.01	69.48	75.49	68.55
BlockDrop	67.98	50.44	72.42	75.52	66.59
w/BlockDrop	67.18	50.49	70.44	73.67	65.45

## 5.8 Evaluation on Transformer Encoder Models.

We evaluate the effectiveness of AttnCache on Transformer Encoder models (BERT base, RoBERTa, and DeBERTa) on the SST-2 dataset. We collect the hidden states and attention maps from the SST-2 training set to train the feature projector, store the attention maps in the database, and test on the validation set. We set the thresholds for Conservative, Moderate, and Aggressive to 0.995, 0.99, and 0.95, respectively. As shown in Table 7, AttnCache yields an average  $1.21\times$  inference speedup with a modest performance degradation of 1.3% to 2.7%. Notably, while some accuracy drop is observed (e.g., in BERT base and RoBERTa), DeBERTa unexpectedly shows slight improvements under conservative reuse, suggesting that cached attention maps can, in some cases, enhance representation quality. These findings confirm that with carefully chosen similarity thresholds, AttnCache can balance efficiency and accuracy even in full-attention encoder settings.

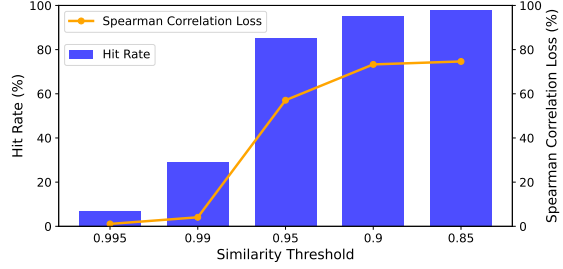


Figure 6. Impact of Threshold on Spearman Correlation.

## 6 ANALYSIS

### 6.1 Impacts of Model Quantization and Pruning

Model quantization represents weights and activations with lower-precision data type, and can improve efficiency in memory usage and inference speed. We integrate AttnCache with quantization and apply Quanto (Optimum, 2024) to all weights, and use 4-bit quantization. We also combine AttnCache with recent LLM pruning methods, AttnDrop and BlockDrop (He et al., 2024), which remove redundant attentions and layers by measuring the similarity between input and output of each layer. Table 8 shows the results. The integration of model quantization and pruning with AttnCache maintains performance: the difference between AttnCache and Quanto/BlockDrop is only 1%, and the difference between AttnCache and AttnDrop is only 2%, on average.

### 6.2 Impact of Similarity Thresholds

Assume that there are  $N$  input sentences for an LLM to generate sentence embeddings, we count how many times AttnCache is successfully applied (indicating similar attention maps are found), denoted as  $M$ . We use the ratio  $M/N$  as the hit rate. We randomly select 100 sentences from STS15, and change the similarity threshold  $\theta$  from 0.995 to 0.85. We measure the hit rate and loss in the Spearman correlation score. As shown in Figure 6. When we reduce  $\theta$ , the hit rate increases, which means that more attention maps are found and AttnCache leads to higher acceleration. However, this might lead to replacement with less similarity, decreasing the performance. By setting  $\theta$  to 0.99, our results show that AttnCache provides 30% hit rate with only 2% reduction in the Spearman correlation score.

## 7 CONCLUSIONS

In this paper, we propose AttnCache to accelerate self attention inference during the prefill stage of LLM inference. Our work is based on the observation that semantically different input sentences can exhibit highly similar attention maps across layers or heads during inference computation. By

---

pre-storing similar attention maps in a database, when generating a new sentence embedding, the most similar attention map can be retrieved from the attention map database and reused to reduce self-attention computation. AttnCache provides an average  $1.2\times$  end-to-end and  $2\times$  attention speedup on CPU, and  $1.6\times$  end-to-end and  $3\times$  attention speedup on GPU, with negligible accuracy loss.

## REFERENCES

- Anthropic. Claude ai, 2025. URL <https://claude.ai/new>.
- BehnamGhader, P., Adlakha, V., Mosbach, M., Bahdanau, D., Chapados, N., and Reddy, S. Llm2vec: Large language models are secretly powerful text encoders. *arXiv preprint arXiv:2404.05961*, 2024.
- Beltagy, I., Peters, M. E., and Cohan, A. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*, 2020.
- Bhojanapalli, S., Chakrabarti, A., Veit, A., Lukasik, M., Jain, H., Liu, F., Chang, Y.-W., and Kumar, S. Leveraging redundancy in attention with reuse transformers. *arXiv preprint arXiv:2110.06821*, 2021.
- bitsandbytes. bitsandbytes. <https://github.com/bitsandbytes-foundation/bitsandbytes>, 2025. Accessed: 2025-05-19.
- ByteDance. Trae ai, 2025. URL <https://www.trae.ai/>.
- Conneau, A. and Kiela, D. Senteval: An evaluation toolkit for universal sentence representations. *arXiv preprint arXiv:1803.05449*, 2018.
- Cursor. Cursor official website, 2025. URL <https://cursor.com/>.
- DeepSeek. Deepseekmoe: Scaling vision-language models with mixture-of-experts. <https://deepseekcoder.github.io/>, 2024. Accessed 2024.
- DeepSeek. Deepseek, 2025. URL <https://www.deepseek.com/>.
- Dettmers, T., Pagnoni, A., Holtzman, A., and Zettlemoyer, L. Qlora: Efficient finetuning of quantized llms. *Advances in neural information processing systems*, 36:10088–10115, 2023.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*, 2018.
- Du, K., Wang, B., Zhang, C., Cheng, Y., Lan, Q., Sang, H., Cheng, Y., Yao, J., Liu, X., Qiao, Y., et al. Prefillonly: An inference engine for prefill-only workloads in large language model applications. In *Proceedings of the ACM SIGOPS 31st Symposium on Operating Systems Principles*, pp. 399–414, 2025.
- Firooz, H., Sanjabi, M., Englhardt, A., Gupta, A., Levine, B., Olgiati, D., Polatkan, G., Melnychuk, I., Ramgopal, K., Talanine, K., et al. 360brew: A decoder-only foundation model for personalized ranking and recommendation. *arXiv preprint arXiv:2501.16450*, 2025.
- Ge, S., Zhang, Y., Liu, L., Zhang, M., Han, J., and Gao, J. Model tells you what to discard: Adaptive kv cache compression for llms. *arXiv preprint arXiv:2310.01801*, 2023.
- Gholamian, S., Romani, G., Rudnikowicz, B., and Skylaki, S. Llm-based robust product classification in commerce and compliance. *arXiv preprint arXiv:2408.05874*, 2024.
- Gim, I., Chen, G., Lee, S.-s., Sarda, N., Khandelwal, A., and Zhong, L. Prompt cache: Modular attention reuse for low-latency inference. *Proceedings of Machine Learning and Systems*, 6:325–338, 2024.
- Girshick, R. Fast r-cnn. In *Proceedings of the IEEE international conference on computer vision*, pp. 1440–1448, 2015.
- GitHub. Github copilot – write code faster, 2025. URL <https://copilot.github.com/>.
- Ham, T. J., Jung, S. J., Kim, S., Oh, Y. H., Park, Y., Song, Y., Park, J.-H., Lee, S., Park, K., Lee, J. W., et al. A<sup>3</sup>: Accelerating attention mechanisms in neural networks with approximation. In *2020 IEEE International Symposium on High Performance Computer Architecture (HPCA)*, pp. 328–341. IEEE, 2020.
- He, P., Liu, X., Gao, J., and Chen, W. DeBERTa: Decoding-enhanced bert with disentangled attention. *arXiv preprint arXiv:2006.03654*, 2020.
- He, S., Sun, G., Shen, Z., and Li, A. What matters in transformers? not all attention is needed. *arXiv preprint arXiv:2406.15786*, 2024.
- He, X., Lin, Z., Gong, Y., Jin, A., Zhang, H., Lin, C., Jiao, J., Yiu, S. M., Duan, N., Chen, W., et al. Annollm: Making large language models to be better crowdsourced annotators. *arXiv preprint arXiv:2303.16854*, 2023.
- Hendrycks, D., Burns, C., Basart, S., Zou, A., Mazeika, M., Song, D., and Steinhardt, J. Measuring massive multitask language understanding. *arXiv preprint arXiv:2009.03300*, 2020.

- 
- Hunter, R., Dudziak, Ł., Abdelfattah, M. S., Mehrotra, A., Bhattacharya, S., and Wen, H. Fast inference through the reuse of attention maps in diffusion models. *arXiv preprint arXiv:2401.01008*, 2023.
- Jiang, H., Li, Y., Zhang, C., Wu, Q., Luo, X., Ahn, S., Han, Z., Abdi, A. H., Li, D., Lin, C.-Y., et al. Minference 1.0: Accelerating pre-filling for long-context llms via dynamic sparse attention. *arXiv preprint arXiv:2407.02490*, 2024.
- Jiang, T., Huang, S., Luan, Z., Wang, D., and Zhuang, F. Scaling sentence embeddings with large language models. *arXiv preprint arXiv:2307.16645*, 2023.
- Johnson, J., Douze, M., and Jégou, H. Billion-scale similarity search with gpus. *IEEE Transactions on Big Data*, 7(3):535–547, 2019.
- Koch, G., Zemel, R., Salakhutdinov, R., et al. Siamese neural networks for one-shot image recognition. In *ICML deep learning workshop*, volume 2, pp. 1–30. Lille, 2015.
- Köpkülü, O., Babaee, M., Hörmann, S., and Rigoll, G. Convolutional neural networks with layer reuse. In *2019 IEEE International Conference on Image Processing (ICIP)*, pp. 345–349. IEEE, 2019.
- Lample, G., Conneau, A., et al. Mistral 7b. <https://mistral.ai/news/introducing-mistral-7b/>, 2023.
- Lan, X., Cheng, Y., Sheng, L., Gao, C., and Li, Y. Depression detection on social media with large language models. *arXiv preprint arXiv:2403.10750*, 2024.
- Lee, C., Roy, R., Xu, M., Raiman, J., Shoeybi, M., Catanzaro, B., and Ping, W. Nv-embed: Improved techniques for training llms as generalist embedding models. *arXiv preprint arXiv:2405.17428*, 2024a.
- Lee, J., Dai, Z., Ren, X., Chen, B., Cer, D., Cole, J. R., Hui, K., Boratko, M., Kapadia, R., Ding, W., et al. Gecko: Versatile text embeddings distilled from large language models. *arXiv preprint arXiv:2403.20327*, 2024b.
- Lei, Y., Wu, D., Zhou, T., Shen, T., Cao, Y., Tao, C., and Yates, A. Meta-task prompting elicits embedding from large language models. *arXiv preprint arXiv:2402.18458*, 2024.
- Lei, Y., Shen, T., Cao, Y., and Yates, A. Enhancing lexicon-based text embeddings with large language models. *arXiv preprint arXiv:2501.09749*, 2025.
- Li, C., Qin, M., Xiao, S., Chen, J., Luo, K., Shao, Y., Lian, D., and Liu, Z. Making text embedders few-shot learners. *arXiv preprint arXiv:2409.15700*, 2024.
- Li, Z. and Zhou, T. Your mixture-of-experts llm is secretly an embedding model for free. *arXiv preprint arXiv:2410.10814*, 2024.
- Liao, B. and Vargas, D. V. Beyond kv caching: Shared attention for efficient llms. *arXiv preprint arXiv:2407.12866*, 2024.
- Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- Malkov, Y. A. and Yashunin, D. A. Efficient and robust approximate nearest neighbor search using hierarchical navigable small world graphs. *IEEE transactions on pattern analysis and machine intelligence*, 42(4):824–836, 2018.
- Men, X., Xu, M., Zhang, Q., Wang, B., Lin, H., Lu, Y., Han, X., and Chen, W. Shortgpt: Layers in large language models are more redundant than you expect. *arXiv preprint arXiv:2403.03853*, 2024.
- Menéndez, M. L., Pardo, J., Pardo, L., and Pardo, M. The jensen-shannon divergence. *Journal of the Franklin Institute*, 334(2):307–318, 1997.
- Muennighoff, N., Tazi, N., Magne, L., and Reimers, N. Mteb: Massive text embedding benchmark. *arXiv preprint arXiv:2210.07316*, 2022.
- Muennighoff, N., Su, H., Wang, L., Yang, N., Wei, F., Yu, T., Singh, A., and Kiela, D. Generative representational instruction tuning. *arXiv preprint arXiv:2402.09906*, 2024.
- Ni, J., Abrego, G. H., Constant, N., Ma, J., Hall, K. B., Cer, D., and Yang, Y. Sentence-t5: Scalable sentence encoders from pre-trained text-to-text models. *arXiv preprint arXiv:2108.08877*, 2021.
- Ning, L. and Shen, X. Deep reuse: Streamline cnn inference on the fly via coarse-grained computation reuse. In *Proceedings of the ACM International Conference on Supercomputing*, pp. 438–448, 2019.
- Ning, L., Guan, H., and Shen, X. Adaptive deep reuse: Accelerating cnn training on the fly. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*, pp. 1538–1549. IEEE, 2019.
- OpenAI. Chatgpt, 2023. URL <https://chatgpt.com/>.
- Optimum. Optimum-quanto, 2024. URL <https://huggingface.co/docs/transformers/main/quantization/quanto>.



- 
- Oren, M., Hassid, M., Yarden, N., Adi, Y., and Schwartz, R. Transformers are multi-state rnns. *arXiv preprint arXiv:2401.06104*, 2024.
- Phan, L., Gatti, A., Han, Z., Li, N., Hu, J., Zhang, H., Zhang, C. B. C., Shaaban, M., Ling, J., Shi, S., et al. Humanity’s last exam. *arXiv preprint arXiv:2501.14249*, 2025.
- Qin, Z., Jagerman, R., Hui, K., Zhuang, H., Wu, J., Yan, L., Shen, J., Liu, T., Liu, J., Metzler, D., et al. Large language models are effective text rankers with pairwise ranking prompting. *arXiv preprint arXiv:2306.17563*, 2023.
- Silfa, F., Dot, G., Arnau, J.-M., and González, A. Neuron-level fuzzy memoization in rnns. In *Proceedings of the 52nd Annual IEEE/ACM International Symposium on Microarchitecture*, pp. 782–793, 2019.
- Song, J., Oh, K., Kim, T., Kim, H., Kim, Y., and Kim, J.-J. Sleb: Streamlining llms through redundancy verification and elimination of transformer blocks. *arXiv preprint arXiv:2402.09025*, 2024.
- Talmor, A., Herzig, J., Lourie, N., and Berant, J. CommonsenseQA: A question answering challenge targeting commonsense knowledge. In Burstein, J., Doran, C., and Solorio, T. (eds.), *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pp. 4149–4158, Minneapolis, Minnesota, June 2019. Association for Computational Linguistics. doi: 10.18653/v1/N19-1421. URL <https://aclanthology.org/N19-1421/>.
- Tang, H., Lin, Y., Lin, J., Han, Q., Hong, S., Yao, Y., and Wang, G. Razorattention: Efficient kv cache compression through retrieval heads. *arXiv preprint arXiv:2407.15891*, 2024.
- Team, Q. Qwen1.5: Enhancing multilingual and multimodal capabilities of language models. <https://github.com/QwenLM/Qwen>, 2024. Accessed 2024.
- Touvron, H., Martin, L., Stone, K., Albert, P., Almahairi, A., Babaei, Y., Bashlykov, N., Batra, S., Bhargava, P., Bhosale, S., et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- Touvron, H., Anastasopoulos, A., et al. Llama 3: Open and efficient foundation language models, 2024.
- Vajjala, S. and Shimangaud, S. Text classification in the llm era—where do we stand? *arXiv preprint arXiv:2502.11830*, 2025.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, Ł., and Polosukhin, I. Attention is all you need. In *NeurIPS*, 2017.
- Wang, A., Singh, A., Michael, J., Hill, F., Levy, O., and Bowman, S. R. Glue: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pp. 353–355, 2018.
- Wang, H., Zhang, Z., and Han, S. Spatten: Efficient sparse attention architecture with cascade token and head pruning. In *2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*, pp. 97–110. IEEE, 2021.
- Wang, L., Yang, N., Huang, X., Yang, L., Majumder, R., and Wei, F. Improving text embeddings with large language models. *arXiv preprint arXiv:2401.00368*, 2023a.
- Wang, Y., Chu, Z., Ouyang, X., Wang, S., Hao, H., Shen, Y., Gu, J., Xue, S., Zhang, J. Y., Cui, Q., et al. Enhancing recommender systems with large language model reasoning graphs. *arXiv preprint arXiv:2308.10835*, 2023b.
- Wu, L., Zheng, Z., Qiu, Z., Wang, H., Gu, H., Shen, T., Qin, C., Zhu, C., Zhu, H., Liu, Q., et al. A survey on large language models for recommendation. *World Wide Web*, 27(5):60, 2024a.
- Wu, R., Zhang, F., Guan, J., Zheng, Z., Du, X., and Shen, X. Drew: Efficient winograd cnn inference with deep reuse. In *Proceedings of the ACM Web Conference 2022*, pp. 1807–1816, 2022.
- Wu, W., Wang, Y., Xiao, G., Peng, H., and Fu, Y. Retrieval head mechanistically explains long-context factuality. *arXiv preprint arXiv:2404.15574*, 2024b.
- Xiao, G., Tian, Y., Chen, B., Han, S., and Lewis, M. Efficient streaming language models with attention sinks. *arXiv preprint arXiv:2309.17453*, 2023.
- Xiao, T., Li, Y., Zhu, J., Yu, Z., and Liu, T. Sharing attention weights for fast transformer. *arXiv preprint arXiv:1906.11024*, 2019.
- Ye, L., Tao, Z., Huang, Y., and Li, Y. Chunkattention: Efficient self-attention with prefix-aware kv cache and two-phase partition. *arXiv preprint arXiv:2402.15220*, 2024.
- Ying, C., Ke, G., He, D., and Liu, T.-Y. Lazyformer: Self attention with lazy update. *arXiv preprint arXiv:2102.12702*, 2021.



- 
- Zhang, B., Chang, K., and Li, C. Simple techniques for enhancing sentence embeddings in generative language models. In *International Conference on Intelligent Computing*, pp. 52–64. Springer, 2024a.
- Zhang, R., Li, Y., Ma, Y., Zhou, M., and Zou, L. Llmaaa: Making large language models as active annotators. *arXiv preprint arXiv:2310.19596*, 2023a.
- Zhang, Y., Li, Y., Wang, X., Shen, Q., Plank, B., Bischl, B., Rezaei, M., and Kawaguchi, K. Finercut: Finer-grained interpretable layer pruning for large language models. *arXiv preprint arXiv:2405.18218*, 2024b.
- Zhang, Z., Sheng, Y., Zhou, T., Chen, T., Zheng, L., Cai, R., Song, Z., Tian, Y., Ré, C., Barrett, C., et al. H2o: Heavy-hitter oracle for efficient generative inference of large language models. *Advances in Neural Information Processing Systems*, 36:34661–34710, 2023b.
- Zhuang, S., Ma, X., Koopman, B., Lin, J., and Zuccon, G. Promptreps: Prompting large language models to generate dense and sparse representations for zero-shot document retrieval. *arXiv preprint arXiv:2404.18424*, 2024.