

# LeanAttention: Hardware-Aware Scalable Attention Mechanism for the Decode-Phase of Transformers

Rya Sanovar

Srikant Bharadwaj

Renee St. Amant

Victor Rühle

Saravan Rajmohan

Microsoft

**Abstract**—Transformer-based models have emerged as one of the most widely used architectures for natural language processing, natural language generation, and image generation. The size of state-of-the-art models has increased steadily reaching billions of parameters. These huge models are memory hungry and incur significant inference latencies even on cutting edge AI-accelerators, such as GPUs. Specifically, the time and memory complexity of the attention operation is quadratic in terms of the total context length, i.e., prompt and output tokens. Thus, several optimizations such as key-value tensor caching and FlashAttention have been proposed to deliver the low latency demands of applications relying on such large models. However, these techniques do not cater to the computationally distinct nature of different phases during inference.

To that end, we propose LeanAttention, a scalable and generalized attention partitioning mechanism for the decode-phase of transformer-based models. LeanAttention enables scaling the attention mechanism for the challenging case of long context lengths by re-designing the attention execution flow for the decode-phase. We identify and prove the associative property of the softmax re-scaling operator, which allows it to function as a reduction operator. This property enables us to extend a "stream-K"-style reduction of tiled calculation to self-attention, which efficiently parallelizes attention computation over large context lengths and achieves near-perfect hardware occupancy irrespective of context size. As a result, we achieve an average of 1.73x speedup in attention execution compared to FlashDecoding, with up to 2.18x speedup for 256k context length.

## I. INTRODUCTION

Transformer-based [36] language models [9], [14], [25], [35], [43] have revolutionized the field of natural language processing (NLP) and found applications across diverse domains [19], [33]. These powerful models, fueled by massive amounts of data and sophisticated architectures, have become indispensable tools for tasks such as machine translation [22], question answering [7], text generation [7], and sentiment analysis.

The core of the transformer architecture is the powerful component of self-attention. However, execution of the self-attention mechanism is slow and suffers from a large memory footprint, especially when dealing with longer contexts. A standard implementation of self-attention has quadratic time and memory complexity with respect to total sequence length, which leads to scalability challenges as model sizes [12] and supported context lengths increase [8], [26], [42]. Despite these scalability challenges, we see a trend of state-of-the-art models supporting greater and greater context lengths, with some production models supporting contexts hundreds

of thousands of tokens long. Support for long context lengths can improve a model's utility by allowing for an increasingly rich context, which is particularly beneficial in a range of applications (e.g. RAG involving numerous or long documents) allowing improved relevance, coherence, and user experience.

To mitigate LLM scalability challenges, mechanisms like FlashAttention [16] and FlashAttention-2/3 [15], [32] have been developed. FlashAttention brings IO-awareness to attention computation by reducing slow reads and writes to and from the GPU's global memory [20]. Instead, it computes attention in the faster shared memory using a tiling strategy. It allows for parallelization over batch size and number of heads. FlashAttention-2 builds on FlashAttention to further optimize attention computation by enabling parallelization over input sequence length (or query length).

While these optimizations provide significant improvements, these mechanisms only provide performance benefits for a subset of problem sizes (i.e. query length, context length, batch size, and number of heads). Their utilization of underlying hardware resources is mostly optimized for problem sizes encountered in the prefill-phase of transformer-based models, and often results in critically low hardware utilization for problem sizes found in the decode-phase. By overlooking the distinct behavior of attention during the decode phase versus the prefill phase, these mechanisms miss out on potential performance gains that could be achieved by efficiently exploiting the parallelization capabilities of the underlying hardware.

In decoder-only transformer models, the inference process for a single request involves multiple forward passes through the model where output tokens are generated sequentially [23]. This inference procedure inherently comprises of two distinct computational phases due to the practice of reusing (caching) the key-value tensors of the previously computed tokens [31]. The first phase is the *prefill-phase* (sometimes known as *prompt-computation phase*) where attention is computed of the entire input prompt against itself to generate the first output token. This phase is computationally intensive and demands high FLOPS/s (floating point operations per second) [23]. Following the prefill-phase, the *decode-phase* (sometimes known as *token-generation phase*) begins in an auto-regressive manner [36]. Each subsequent token is produced based on the attention computed of the preceding token against itself and the entire cached context (*kv-cache*) of previous tokens in the

sequence. With the push towards longer contexts, this cached context length can get extremely long, exceeding more than hundreds of thousands of tokens in length [17], [24], [26], [42]. Despite state-of-the-art batching techniques [40] and attention partitioning mechanisms [4], [15], [16], [39], the lack of a smart parallelized execution of attention along this long context length makes the decode-phase slow, bound by memory bandwidth [37] and capacity [23]. Importantly, as we discuss in section III, even when the query length is significantly longer than the total number of output tokens produced during inference, the majority of the overall processing time is consumed by the decode-phase.

It's evident now that efficient parallelization of attention computation over the context length dimension is highly necessary. Although mechanisms like FlashDecoding [4] and FlashInfer [39] enable this parallelization through the fixed-split partitioning strategy, they provide the hardware with imbalanced loads and as a consequence still suffer from hardware resource underutilization and unnecessary reduction overheads - both being highly contingent on problem size. Further, attention optimizations [10] are increasingly relying on batching requests of unequal context lengths to improve overall throughput but suffer from hardware under-utilization because of the partitioning strategies adopted by FlashAttention and FlashDecoding/FlashInfer.

In this work, we aim to address the limitations of previous work as it relates to the decode-phase of inference, which we find exhibits unique computational characteristics in comparison to the prefill-phase.

We introduce LeanAttention, a generalized exact-attention mechanism which enables parallelization across all problem size dimensions, ensures perfect quantization efficiency, i.e. 100% GPU occupancy, for all problem sizes with constant reduction overheads, delivers a runtime speedup in attention computation for long context lengths, and is scalable to multi-GPU scenarios with tensor parallelism.

Overall, our contributions are as follows:

- Identify the limitations of state-of-the-art attention execution optimizations on GPUs during the decode-phase of transformer-based models. (subsection III-B)
- Identify the *associative nature* of the softmax re-scaling operator that enables it to function as a reductive operator, and leverage this crucial property in LeanAttention to split an independent attention workload into unequal sizes along the context length dimension when needed. (subsection IV-A)
- Leverage a stream-K style [30] partitioning in LeanAttention that *always* provides *equalized* compute loads to every compute unit in the hardware system (as shown in Figure 1), thus giving near 100% hardware occupancy and delivering speedup irrespective of problem size and hardware architecture. (subsection IV-C)
- Expatriate LeanAttention's versatility and generalizability, where FlashAttention-2 and FlashDecoding (as well as FlashInfer) can be recovered as special cases of it. (subsection IV-C)

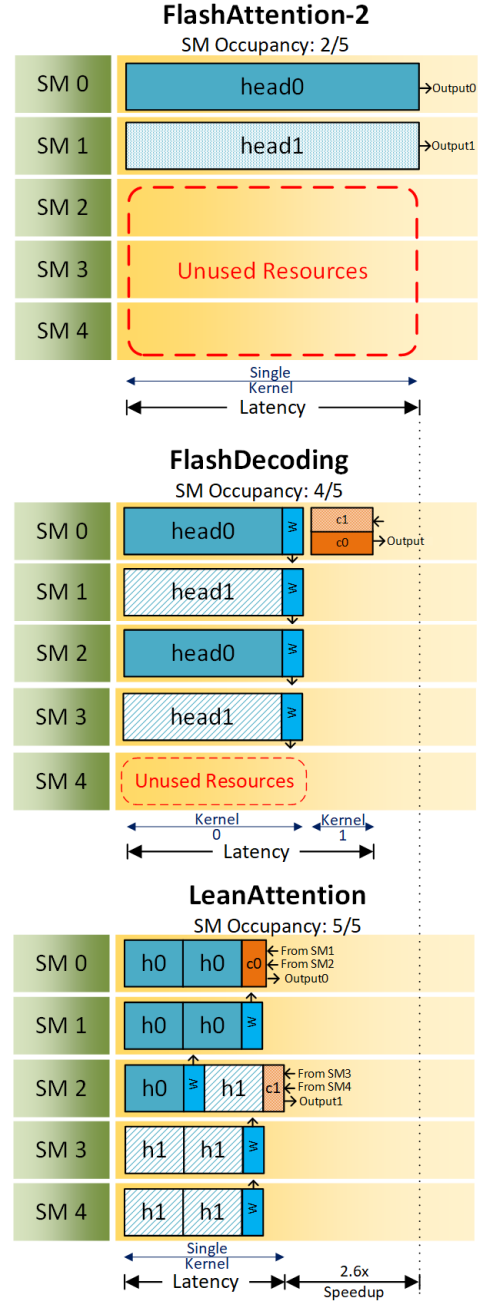


Fig. 1. Execution schedule of FlashAttention-2 [15], FlashDecoding [4] and FlashInfer [39] (fixed-split), and LeanAttention across a hypothetical five SM GPU executing attention of 2 heads. LeanAttention splits the context into optimal LeanTiles (shown here with 5 tiles per head).

As detailed in section VI, LeanAttention results in an average of 1.73x latency speedup over FlashDecoding for the decode phase of transformer-based models and up to 2.18x speedup over FlashDecoding for 256k context length, while maintaining a near 100% GPU occupancy irrespective of problem size.

## II. BACKGROUND

In this section, we provide required background on Standard Attention [36] and FlashAttention-2 [15].

### A. Standard Attention

For a given input tensor with dimensions of batch size  $B$ , query length  $N_q$ , key/value sequence length (also known as context length)  $N_k$ , and hidden dimension  $D$ , multi-head attention typically splits attention computation into  $h$  number of heads along the hidden dimension, with each head responsible for computing attention independently for a head dimension of size  $d = D/h$ .

Unlike standard transformer execution, the query length and context lengths may not always be equal, where key-value tensors are cached [31]. For instance, the prefill-phase of generative decoder-only transformers such as GPT-4 [9] or Phi-2 [25] has sequence lengths  $N_q = N_k = N$ , but in their decode-phase the context length increments by 1 after every autoregressive step of decode generation, while the query length, for a given batch instance and head, is the singular token that was generated in the previous  $n$ -th time step, i.e.,  $N_q = 1$  and  $N_k = N + n$ .

The query matrix  $Q \in R^{N_q \times d}$  and key  $K$  and value  $V$  matrices  $\in R^{N_k \times d}$  are inputs to the following equation which is computed independently by the different batch instances and heads. The output matrix  $O \in R^{N_q \times d}$  is obtained in essentially three steps as shown in Equation 1. Table I summarises the three operations involved in self-attention along with their corresponding dimensions in both prefill and decode-phase at time step  $n = 0$ .

$$S = QK^T, P = \text{softmax}\left(\frac{S}{\sqrt{d}}\right), O = PV \quad (1)$$

Standard attention implementation involves computing the large intermediate matrices, namely the attention score matrix  $S \in R^{N_q \times N_k}$  and the softmax matrix  $P \in R^{N_q \times N_k}$  and storing them in global memory. These intermediate matrices need to be stored in the global memory because the computation of the softmax matrix  $P$  requires *a priori* knowledge of all tokens in a given row, specifically the row-wise maximum and exponential sum of tokens in a row need to be computed beforehand to calculate the softmax-ed value of each element in the row.

The computational complexity of standard attention is on the order of  $O(N_q N_k d)$ , with the two matrix multiplications (MatMul's) contributing to the majority of it. Due to slow global memory access speeds, storing and retrieving these intermediate matrices [20] is costly in terms of latency and incurs a large memory footprint, both in the order of  $O(N_q N_k)$ .

Operation	Type	Operation Dimension	
		Prefill	Decode
$query \times key$	MatMul	$N \times d \times N$	$1 \times d \times N$
$softmax$	EleWise	$N \times N$	$1 \times N$
$attn\_score \times value$	MatMul	$N \times N \times d$	$1 \times N \times d$

TABLE I  
OPERATIONS IN SELF-ATTENTION. MATRIX MULTIPLICATIONS ARE DESCRIBED IN THE  $M \times N \times K$  FORMAT.

### B. Flash Attention-2

To mitigate the memory footprint and access overhead [20] associated with storing the  $S$  and  $P$  matrices, FlashAttention employs kernel fusion of the three operations as shown in Equation 1:  $query \times key$  MatMul, softmax and  $attn\_score \times value$  MatMul, requiring no intermediate global memory reads and writes. To this end, it employs the tiling strategy.

By utilizing the online softmax algorithm [29], FlashAttention only requires a single pass over an entire row of tokens to compute their softmax, bypassing the issue of *a priori* knowledge in standard attention. This helps leverage the tiling strategy which partitions the attention output matrix  $O$  into independent output tiles (attention computation an output tile is independent of the computation of other output tiles). A grid of *cooperative thread arrays* (CTAs)<sup>1</sup> is launched, each computing a given output tile of the output matrix  $O$ . The input matrices  $Q$ ,  $K$  and  $V$  are partitioned into smaller tiles too. While computing the output tile corresponding to a given query tile, the key/value tiles are brought into shared memory in a *sequential* manner and the attention output tile is continuously updated and corrected by the right scaling factors. This on-chip updation avoids the need of storing the intermediate  $S$  and  $P$  matrices in global memory. In addition to parallelizing computation over batches and heads like FlashAttention, FlashAttention-2 further parallelizes over the query length dimension, as the attention computation of output tiles along this length is independent. This results in a 2x speedup over FlashAttention.

Thus, the tiling strategy ensures that the extra global memory space required by FlashAttention-2 is  $O(N_q)$  (needed to store the logexpsum  $L$  for backward pass), an impressive improvement in memory footprint over the  $O(N_q \times N_k)$  in traditional attention. The additional parallelism over query length helps it reach 50-70% of peak theoretical FLOPS/s and increases hardware occupancy in the prefill phase. FlashAttention-2 was augmented to FlashAttention-3 [32], specifically fine-tuned for execution on Hopper H100 GPU's [28] to exploit its low-precision and asynchronous hardware capabilities.

Like FlashAttention-2, other related techniques such as Ring Attention [26] and Striped Attention [11], are optimized for prefill-phase problem sizes and thus suffer from longer latencies during the decode phase.

## III. CHALLENGES IN THE DECODE PHASE

Prior to outlining our methodology for LeanAttention, to set the stage for our approach, we delve into some of the challenges encountered in the decode phase of LLM inference, as well as the limitations of FlashAttention-2 optimizations in the decode phase.

### A. Time Spent in Decode Phase

As we've discussed, modern generative LLM inference comprises of two computationally distinct phases: the prefill phase followed by the decode phase. In the prefill phase, self-attention is computed for the entire input prompt. The query

<sup>1</sup>Blocks of GPU threads are coscheduled in CTAs, which virtualize the hardware's streaming multiprocessor cores (SMs)

length,  $N_q$ , in this phase is the same as the context length,  $N_k$ , i.e., ( $N_q = N_k = N$ ). Whereas, the decode phase begins generating each subsequent output token in autoregressive iterations. For each iteration of the decode phase, its query length is a single token,  $N_q = 1$ , and its context length,  $N_k$ , could be very long, in the order of more than thousands of tokens depending on the auto-regressive step and input query length.

Figure 2 depicts the processing time breakdown of the prefill and decode phase, with the decode phase’s further breakdown into time spent in the  $Q/K/V$  activation layer, the decode attention layer and the feed-forward linear layers.

While the large matrix multiplications found in the linear layers of the prefill phase of inference are heavily optimized (all the model layers during prefill phase only taking up 10% of the timeshare even for a high prompt:output ratio), the decode phase presents a different challenge. During the decode phase, where the query length is only 1 token long, linear layers perform matrix multiplications on very narrow matrices which do not provide enough work to occupy the GPU. MatMul partitioning strategies like Stream-K [30] can be leveraged to efficiently partition these narrow matrices and accelerate their computation, preventing the linear layers from becoming a bottleneck during decode phase. However, the attention layer, with the existing attention partitioning techniques [4], [15], [16], [39] experiences longer latencies along with significant underutilization of hardware resources during this phase. This makes leveraging efficient parallelism along context length ( $N_k$ ) during attention a crucial aspect in increasing SM occupancy and reducing decode phase processing time.

As the number of output tokens generated rises, the context length becomes longer and thus the proportion of time spent in the decode phase relative to the prefill phase becomes larger. Figure 2 depicts this imbalance in processing time spent in the prefill phase and the decode phase attention. Even with a prompt input to output token ratio of 8:1, more than 50% of the processing time is consumed by the decode phase, taking up to nearly 80% of the timeshare for longer prompt sizes. Additionally, other layers of the decode phase such as QKV and MLP (FFN layers) can be optimized using state-of-the-art MatMul partitioning techniques such as Stream-k. These operations are typically quantized to lower data formats such as INT8 to further enhance their efficiency. As a result, the attention operation can constitute up to 40-50% of the total duration of decode phase inference as shown.

### B. Limitations of FlashAttention-2 for Decode

In both the prefill and decode phase, FlashAttention-2 traverses the context length dimension ( $N_k$ ) sequentially, i.e., it updates the attention output for a given query tile by bringing in the key/value tiles into shared memory in a sequential manner. While FlashAttention-2 does parallelize over query length ( $N_q$ ) to increase SM occupancy, this additional mode of parallelism has limited parallelization capacity in the decode phase where the query is a single token ( $N_q = 1$ ). Not parallelizing attention computation along context length makes FlashAttention-2 [15]

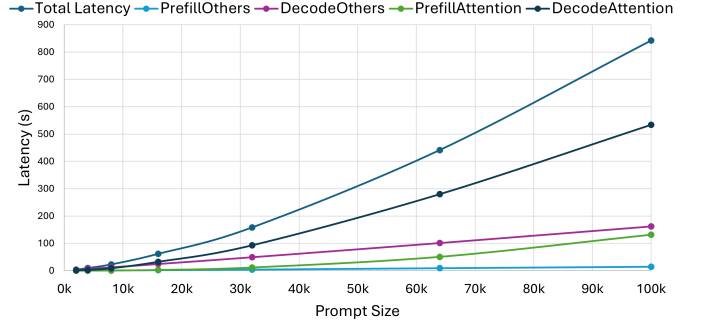


Fig. 2. Timeshare of decode attention compared to other stages for different prompt sizes with 8:1 token ratio for Phi-3 Medium model with single batch size.

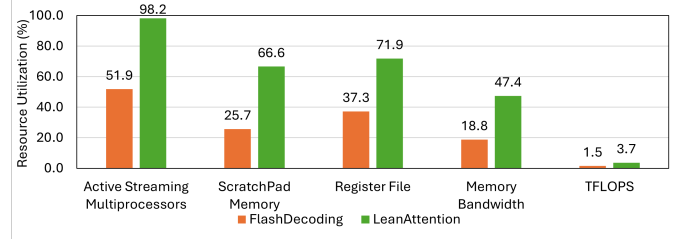


Fig. 3. Utilization of various resources on a single Nvidia-A100-80GB GPU in LeanAttention compared to FlashDecoding kernel at Heads=56 and BS=1 measured using Nsight Compute. FlashDecoding has a quantization efficiency issue with the 108 SMs on the GPU. LeanAttention occupies all SMs available in the system.

suffer from extremely low SM occupancy during decode as depicted in Figure 3. This means that at any given point in time, the number of CTAs in flight on the GPU is directly proportional to the number of query tiles, and, therefore, to the query length - regardless of the context length.

More explicitly, for a single batch instance, the maximum number of heads for state of the art LLMs barely occupy the compute resources of modern hardware architecture systems during the decode phase where query length  $N_q = 1$ . For example, for a model with 128 heads, its decode phase would suffer from severe under-utilization of an 8 GPU A100 system that has 864 compute cores at its disposal. Unlike the prefill phase, decode phase can offer parallelization only across batch size and number of heads for FlashAttention-2.

Processor occupancy in FlashAttention-2 could be improved by increasing the batch size or number of heads, the other two modes of parallelization it addresses. Intuitively, having larger batch sizes in the decode phase could provide enough work to every compute resource to fully occupy the GPU, but this introduces other challenges and limitations. Due to increasingly large model sizes, the need to independently cache KV context for every batch instance would likely exceed the memory capacity of the hardware system. Moreover, scheduling overheads [41] for efficiently batching queries along with the challenges of batching low SLA queries would increase inference latency and challenge utilization.

Without having to resort to larger batch sizes as the sole solution to resolving the GPU occupancy issue (which is limited by available memory capacity), the large context

length in the decode phase would benefit from partitioning its workload across different SMs efficiently. This motivates the need for smarter attention decomposition techniques which can efficiently distribute the workload across the cores without resorting to larger batch sizes.

### C. Limitations of Related Work

FlashDecoding, which is FlashAttention-2 with fixed-split partitioning, has recently been proposed [4], [5], [18], where attention computation is also partitioned along context length  $N_k$ . FlashInfer [39] implements an identical fixed-split partitioning of attention for single-request's in decode phase in pure CUDA. For the case of batched-requests in decode, FlashInfer implements an optimized version of PagedAttention for efficient KV cache storing and fetching.

Fixed-split is a general matrix multiplication decomposition scheme that we briefly describe here. Given a MatMul computation problem with matrices  $A (M \times K)$  and  $B (N \times K)$  to obtain a matrix  $C (M \times N)$  where  $C = AB^T$ , to optimize concurrent computation, the fixed-split mechanism [6] partitions the K-mode of the  $A$  and  $B$  matrices into  $s$  batches based on a fixed splitting factor  $s$  provided dynamically at run time. This launches  $s$  times the CTAs (Cooperative Thread Arrays, equivalent to a threadblock) as launched without fixed-split, which are computing partial products of the output tiles of the  $C$  matrix concurrently. Fixed-split utilizes the associativity of addition in the inner product of a MatMul to later reduce or “fix-up” the partially computed  $C$  matrices to produce the final  $C$  matrix. The concurrency from fixed-split reduces latency and simultaneously increases hardware occupancy at the cost of an additional reduction at the end.

FlashDecoding++ [18] achieves speedup over FlashDecoding by approximating the softmax operation to remove the sequential dependencies it creates in attention. Notably, this approach compromises on accuracy and its implementation is limited to certain model architectures. In contrast, LeanAttention computes exact attention with no loss in accuracy and can be used in any transformer-based model. FlashDecoding++, as well as other works that focus on softmax approximations to achieve speedup (like ConSmax [27] and Softmax [34]), can be seamlessly integrated into LeanAttention.

Despite these improvements, fixed-split used in these mechanisms [4], [18], [39] is a non-optimal load balancing strategy. While this method of partitioning would provide speedup and occupy the GPU well for some attention workloads, it is an inefficient load balancing strategy for the entire problem space and often results in partially full waves of attention computation that suffer from quantization inefficiencies, i.e. low GPU occupancy due to imbalanced loads, and loses out on performance gains it could get from the idle resources otherwise (depicted in Figure 1). While increasing the number of splits could help occupy the GPU better, it would result in reduction overheads that scale with the split factor and would allocate minimal work to each SM making it an inefficient use of register space.

This fixed-split partitioning along the context length does occupy a larger number of compute resources on the GPU compared to vanilla FlashAttention-2, but it often provides imbalanced loads to the compute units, and GPU occupancy then varies greatly with problem size, split factor and number of compute units in the hardware system as shown in Figure 3, making it unlikely for FlashDecoding and its variants to reach perfect quantization efficiency for all problem sizes and hardware systems. Moreover, the problem of quantization inefficiencies with these mechanisms would be particularly exacerbated in the common cases of processing a batch of requests of heterogenous context lengths [10]. In contrast, LeanAttention, with its stream-K-style decomposition discussed in section IV, will *always* provide well-balanced loads to each compute unit in the hardware system and reach near 100% GPU occupancy for all problem sizes and hardware architectures, making it perfectly adept to handle batched requests of unequally sized contexts.

### D. Multi-GPU Execution with Tensor Parallelism

FlashAttention-2 not only severely underutilizes GPU cores in the decode phase, but is also not adaptable to multi-GPU scenarios due to its lack of support for tensor parallelism. This makes FlashAttention-2 less scalable to multi-GPU systems which has become an imperative due to capacity-boundedness of contemporary large language models [12] and the support they require for increasingly long context lengths [8]. This asserts the need for an attention mechanism that also scales well to multi-GPU scenarios.

These challenges motivate the need for a generalized attention mechanism that works for a vast set of problem sizes (in both prefill and decode phase) and is closely aligned with the memory and compute hierarchies of modern hardware systems. We formulate this scalable and generalized exact attention mechanism as **LeanAttention**, which computes exact attention faster in a single fused kernel launch, has optimal quantization efficiency for all kinds of problem sizes, whilst also being scalable to multi-GPU scenarios through its support for tensor parallelism.

## IV. LEANATTENTION

LeanAttention is an optimized scalable attention execution mechanism. It provides extensive parallelism across all modes of the attention tensor, with well-balanced computation workloads to each CTA ensuring close to 100% SM occupancy while delivering a runtime speedup in attention execution.

First, we identify the smallest optimal granularity of decomposition in attention computation, termed as *LeanTile* (subsection IV-B), which can be linearly mapped on the hardware resources in a flexible style akin to stream-k decomposition of matrix multiplications subsection IV-C). Multiple such LeanTiles belonging to either single or multiple attention outputs will constitute a workload assigned to a CTA. By the nature of Stream-K's equalized load balancing strategy, each CTA will compute equal number of LeanTiles, ensuring no idle SMs during the entire duration of attention computation.



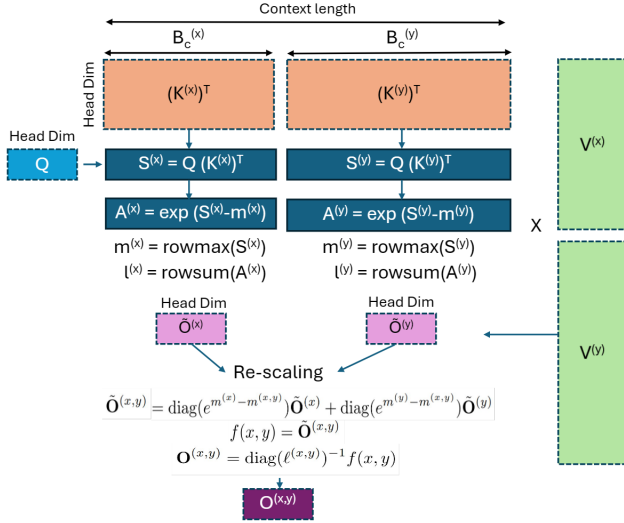


Fig. 4. Illustrative diagram showing LeanAttention’s partitioning strategy with two differently sized work volumes of a head assigned to different CTAs. The un-scaled outputs are independently computed and re-scaled later in a reduction operation. Note that this can be generalized to any arbitrary-sized work volume split.

Second, we identify that the associative property of softmax re-scaling enables us to treat it as a reduction operation along the context length dimension and allows us to split the workload (i.e. KV tensors) of a single head into unequally sized blocks (described in subsection IV-A). In the most common case where processor width (number of SMs) is not a multiple of the total number of heads, we must split the workload of each head into unequal sizes on the SMs *in order* for the *total* workload per SM to be equal (unlike FlashDecoding, FlashInfer and variants), and by virtue of this associative nature we can correctly reduce the partial attention outputs corresponding to these unequally sized blocks.

In the following subsections, we first outline the identification of softmax re-scaling as a reduction operation, followed by a conceptualization of a *LeanTile* as a unit granularity in a CTA block and the stream-K style mapping within these CTAs, followed by an explanation of the overall execution flow of LeanAttention. Figure 5 shows the high-level architectural execution flow of a single CTA computing LeanAttention.

#### A. Softmax Re-scaling as Reduction

In LeanAttention, we propose computation of partial attention outputs of a given query tile concurrently on different hardware units, while ensuring that we have a well-balanced work distribution across all hardware units through a Stream-K style decomposition of attention (discussed later in subsection IV-C). This decomposition results in splits of work for a given SM that are not always equal in size, i.e., the key/value tensors of a given query tile are not dispatched in same-sized blocks to different SMs (unlike FlashDecoding [4], FlashInfer [39]). For example, in Figure 1, for computing LeanAttention for a query tile  $h_0$ , SM0 and SM1 receive same-sized KV blocks (each KV block consists of 2 LeanTiles), but SM2 receives

half the amount of work for  $h_0$  (the KV block for it consists of 1 LeanTile) than SM0 or SM1 received.

To reduce these partial attention outputs that result from differently sized blocks, we use a softmax re-scaling operation. This requires us to identify softmax re-scaling’s associativity property that allows it to correctly reduce blocks of unequal sizes, i.e., application of softmax re-scaling as a reduction operator will give the same exact attention output with no loss in accuracy, *regardless of the way the work might be split, whether in same-sized blocks or arbitrary differently sized blocks*.

Without loss of generality, we describe this process of reduction to obtain one row vector of the attention score matrix  $\mathbf{S}$ , of the form  $[\mathbf{S}^{(x)} \quad \mathbf{S}^{(y)}]$  consisting of some unequal length vectors  $\mathbf{S}^{(x)}, \mathbf{S}^{(y)}$  where  $\mathbf{S}^{(x)} \in \mathbb{R}^{1 \times B_c^{(x)}}$  and  $\mathbf{S}^{(y)} \in \mathbb{R}^{1 \times B_c^{(y)}}$ , where 1 is the query length and  $B_c^{(x)}$  and  $B_c^{(y)}$  are the unequal context lengths. The vectors  $\mathbf{S}^{(x)}$  and  $\mathbf{S}^{(y)}$  were computed from  $\mathbf{Q} \times (\mathbf{K}^{(x)})^T$  and  $\mathbf{Q} \times (\mathbf{K}^{(y)})^T$  as shown in Figure 4. Note that, to generalize this procedure for blocks of any size, the context length of  $\mathbf{K}^{(x)}$  and  $\mathbf{K}^{(y)}$  are  $B_c^{(x)}$  and  $B_c^{(y)}$  and are not necessarily equal.

The attention computation is split into two parts. The first part involves calculation of an “un-scaled” version of  $\mathbf{O}^{(i)}$  (where  $i$  is either  $x$  or  $y$ ) along with statistics  $m^{(i)}$  and  $\ell^{(i)}$ :

$$\begin{aligned} \mathbf{S}^{(i)} &= \mathbf{Q}(\mathbf{K}^{(i)})^T \in \mathbb{R}^{1 \times B_c^{(i)}} \\ m^{(i)} &= \text{rowmax}(\mathbf{S}^{(i)}) \in \mathbb{R}^{1 \times 1} \\ \ell^{(i)} &= \text{rowsum}(e^{\mathbf{S}^{(i)} - m^{(i)}}) \in \mathbb{R}^{1 \times 1} \\ \mathbf{A}^{(i)} &= \exp(\mathbf{S}^{(i)} - m^{(i)}) \in \mathbb{R}^{1 \times B_c^{(i)}} \\ \tilde{\mathbf{O}}^{(i)} &= \mathbf{A}^{(i)} \mathbf{V}^{(i)} \in \mathbb{R}^{1 \times d} \end{aligned}$$

**Softmax Re-scaling Operation.** The second part involves re-scaling the “un-scaled” outputs  $\tilde{\mathbf{O}}^{(i)}$  using the previously computed statistics  $m^{(i)}$  and  $\ell^{(i)}$ .

We define the softmax re-scaling operation  $f(x, y)$  for two intermediate outputs  $\mathbf{O}^{(x)}$  and  $\mathbf{O}^{(y)}$  as follows:

$$\begin{aligned} m^{(x,y)} &= \max(m^{(x)}, m^{(y)}) \\ \ell^{(x,y)} &= e^{m^{(x)} - m^{(x,y)}} \ell^{(x)} + e^{m^{(y)} - m^{(x,y)}} \ell^{(y)} \\ f(x, y) &= \text{diag}(e^{m^{(x)} - m^{(x,y)}}) \tilde{\mathbf{O}}^{(x)} + \text{diag}(e^{m^{(y)} - m^{(x,y)}}) \tilde{\mathbf{O}}^{(y)} \\ f(x, y) &= \tilde{\mathbf{O}}^{(x,y)} \\ \mathbf{O}^{(x,y)} &= \text{diag}(\ell^{(x,y)})^{-1} f(x, y) \end{aligned}$$

**Proof of Associativity** The associative nature of softmax re-scaling  $f(x, y)$  allows us to reduce intermediate outputs produced from key/value vectors of different lengths in LeanAttention. We shall briefly prove that  $f(f(x, y), z) = f(x, f(y, z)) = f(x, y, z)$ , where:  $f(x, y) = \tilde{\mathbf{O}}^{(x,y)}$ ,  $f(y, z) = \tilde{\mathbf{O}}^{(y,z)}$  and  $f(x, y, z) = \tilde{\mathbf{O}}^{(x,y,z)}$ .

Proving that  $f(f(x, y), z) = f(x, y, z)$ :

$$f(x, y) = \tilde{\mathbf{O}}^{(x,y)}$$

$$\begin{aligned}
f(f(x, y), z) &= \text{diag}(e^{m^{(x,y)} - m^{((x,y),z)}}) \tilde{\mathbf{O}}^{(x,y)} \\
&+ \text{diag}(e^{m^{(z)} - m^{((x,y),z)}}) \tilde{\mathbf{O}}^{(z)} \\
&= \text{diag}(e^{m^{(x,y)} - m^{(x,y,z)}}) \tilde{\mathbf{O}}^{(x,y)} \\
&+ \text{diag}(e^{m^{(z)} - m^{(x,y,z)}}) \tilde{\mathbf{O}}^{(z)} \\
&= \text{diag}(e^{m^{(x,y)} - m^{(x,y,z)}}) \\
&\times (\text{diag}(e^{m^{(x)} - m^{(x,y)}}) \tilde{\mathbf{O}}^{(x)} \\
&+ \text{diag}(e^{m^{(y)} - m^{(x,y)}}) \tilde{\mathbf{O}}^{(y)} \\
&+ \text{diag}(e^{m^{(z)} - m^{(x,y,z)}}) \tilde{\mathbf{O}}^{(z)} \\
&= \text{diag}(e^{m^{(x)} - m^{(x,y,z)}}) \tilde{\mathbf{O}}^{(x)} \\
&+ \text{diag}(e^{m^{(y)} - m^{(x,y,z)}}) \tilde{\mathbf{O}}^{(y)} \\
&+ \text{diag}(e^{m^{(z)} - m^{(x,y,z)}}) \tilde{\mathbf{O}}^{(z)} \\
&= \tilde{\mathbf{O}}^{(x,y,z)} = f(x, y, z)
\end{aligned}$$

Therefore,  $f(f(x, y), z) = f(x, y, z)$  and similarly  $\ell^{((x,y),z)} = \ell^{(x,y,z)}$ . For brevity, we omit the proof of  $f(x, f(y, z)) = f(x, y, z)$ , but it can be deduced in a similar manner.

This associativity of softmax re-scaling is leveraged in LeanAttention to concurrently calculate the “partial” outputs produced from unequally sized KV blocks and then “reduce” them to obtain exact attention.

### B. LeanTile

To enable us to efficiently distribute the work of computing the attention output tiles, we define the smallest granularity of a KV block as a *LeanTile*. A single LeanTile iteration computes “local attention” across a subset of tokens along the  $N_k$  dimension as shown in the grey box of Figure 5. Thus, a LeanTile takes in a query, key, and value tensor and computes local attention to generate an un-scaled attention output.

Algorithm 1 depicts the subroutine for computing the partial attention outputs for a sequence of LeanTile’s. This LeanTile() subroutine is called when computing each partial output tile in a CTA launched in LeanAttention, as will be discussed later (Algorithm 2).

To efficiently split attention into smaller tiles, it is necessary to identify the smallest tile size capable of achieving the highest compute efficiency. LeanTile size depends on the computational power and memory access costs and, thus, are fixed for a particular hardware architecture. After an extensive empirical sweep through various sizes for a *LeanTile*, we found a tile size granularity of 256 and 128 tokens along the  $N_k$  dimension to be the most optimal for a head size of 64 and 128 respectively for FP16→32 problems while experimenting on an A100 GPU [13], [21]. This optimal size can similarly be identified for other head dimensions and hardware architectures.

### C. Decomposition and Mapping of LeanTiles

Finally, LeanAttention uses a stream-K [30] style decomposition and mapping of these LeanTiles to deliver

---

#### Algorithm 1 LeanTile() for a sequence of lean tile iterations.

---

```

1: function LeanTile(tile_idx, iter_begin, iter_end)
2:   _shared_  $O_{acc}[T_m, d]$ 
3:   _shared_  $Q_f[T_m, d]$ 
4:   _shared_  $K_f[T_n, d]$ 
5:   _shared_  $V_f[T_n, d]$ 
6:   _shared_  $m[T_m, 1]$ 
7:   _shared_  $l[T_m, 1]$ 
8:   Initialize  $O_{acc}$  to  $(0)_{T_m \times d} \in R^{T_m \times d}$  in SMEM.
9:   Initialize  $m$  to  $(-\infty)_{T_m \times 1}$  and  $l$  to  $(0)_{T_m \times 1} \in R^{T_m \times 1}$  in SMEM.
10:   $mm = T_m \times (\text{tile\_idx} / 1)$ 
11:   $nn = d \times (\text{tile\_idx} \% 1)$ 
12:  Perform lean tile iterations for this output tile.
13:  for  $iter = \text{iter\_begin}$  to  $\text{iter\_end}$  do
14:     $kk = iter \times T_n$ 
15:    load fragments from GMEM to SMEM
16:     $Q_f = \text{LoadFragment}(Q, mm, nn)$ 
17:     $K_f = \text{LoadFragment}(K, nn, kk)$ 
18:     $V_f = \text{LoadFragment}(V, nn, kk)$ 
19:    Compute on chip:
20:     $S_f = Q_f K_f$  where  $S_f \in R^{T_m \times T_n}$ 
21:     $m^{new} = \max(m, \text{rowmax}(S_f))$ 
22:     $P_f = \exp(S_f - m^{new})$  where  $P_f \in R^{T_m \times T_n}$ 
23:     $l^{new} = e^{m - m^{new}} l + \text{rowsum}(P_f)$ 
24:     $O_{acc} = P_f V_f + \text{diag}(e^{m - m^{new}}) O_{acc}$ 
25:     $l = l^{new}, m = m^{new}$ 
26:  end for
27:  return  $O_{acc}, l, m$ 
28: end function

```

---

efficient execution of attention.

**Stream-K Decomposition.** Stream-K is a parallel decomposition technique for dense matrix-matrix multiplication on GPUs. Stream-k partitioning addresses the inefficiencies in fixed-split by dividing the total workload (MAC operations) equally to all the CTAs using a pre-determined optimal tile size. It does this by rolling out the inner mode iterations of all output tiles and appending them along the inner mode to form a linear mapping. With the given grid size, it divides this total work into buckets demarcated appropriately such that each CTA has equal amount of MAC operations to perform. The grid size is fixed for a given tile size as LeanAttention provides equal work to all SMs. For example, for a tile of 256 tokens, two CTAs can be co-executed in a single wave with the available shared memory of A100 GPU. This would result in a total grid size of  $108(\text{NumSMs}) \times 2 = 216$ . Number of tiles to be computed by each CTA can, thus be calculated as follows:

$$\text{TilesPerCTA} = \frac{\text{BatchSize} \times \text{NumHeads} \times \text{ContextLength}}{\text{TileSize} \times \text{NumSMs} \times \text{MaxCTAsPerSM}} \quad (2)$$

LeanAttention extends Stream-K style of linear mapping of iterations by rolling out LeanTile iterations in a similar fashion, assigning equal number of LeanTiles to every CTA as shown in Figure 1. Each CTAs range of LeanTile iterations

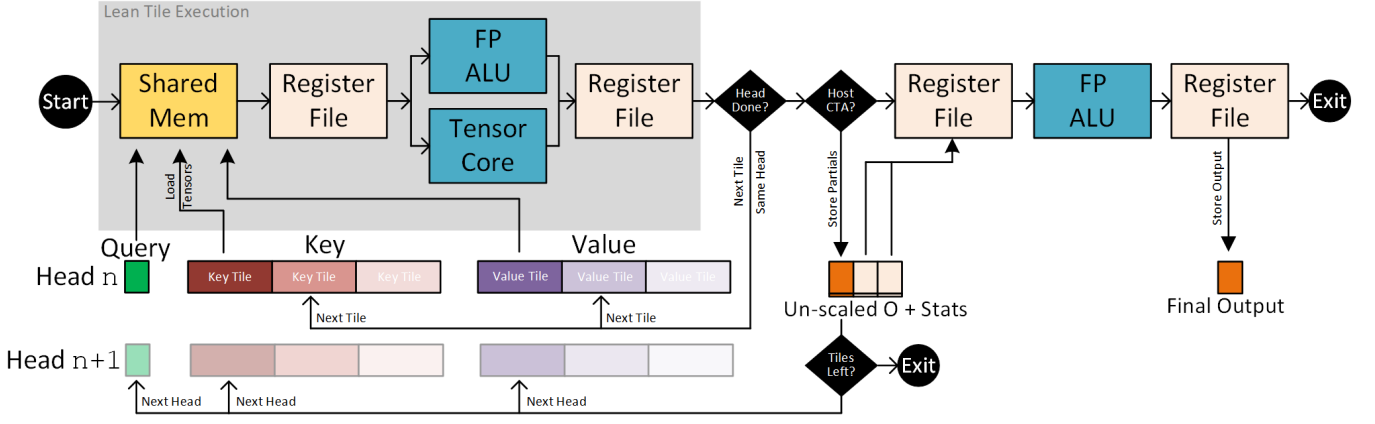


Fig. 5. Control and dataflow of a single CTA in LeanAttention utilizing various hardware resources. The tensors are loaded to shared memory in a tiled manner. At the end of a head, a reduction is performed if it is a host CTA or the partial un-scaled results are written to memory before moving to the next head.

is mapped contiguously into the batch size  $\rightarrow$  heads  $\rightarrow$  context length linearization, crossing the head and query boundary as it may. Should a given CTA’s starting and/or ending LeanTile not coincide with the head’s boundary, it must consolidate its partial output with those of the other CTA(s) also covering that head’s output tile. In our implementation of LeanAttention, each attention output tile is consolidated by the CTA that performed that output’s first LeanTile (called as a host block). Before it can do so, however, it must accumulate the un-scaled output tensors from other CTA(s) in temporary global storage, as shown in Figure 1. The negligible synchronization overhead of original stream-K implementation also extends to LeanAttention, thus leading to near 100% occupancy of SMs (not tensor core utilization) during the execution of a single CTA. Note that the temporary global storage overhead is minimal in the case of decode-phase where the output tensors are of dimensions  $1 \times \text{head\_dim}$ , where  $\text{head\_dim}$  is typically in the range of 64 to 256.

Further, since we distribute the overall attention problem into optimal LeanTiles, we achieve a near 100% quantization efficiency irrespective of problem size (context length). This cohesive implementation of parallel computation and reduction happens in a single kernel launch in LeanAttention, avoiding the reduction kernel launch overheads that FlashDecoding suffers from. A difference in Stream-K decomposition in LeanAttention is in the reduction or “fix-up” phase. While Stream-K for MatMuls has addition as its reductive operator, LeanAttention has softmax re-scaling as its reductive operator.

Naturally, some CTA’s will be computing LeanTile iterations of more than one independent output tile. In such cases, stream-K’s equalized partitioning makes lean attention more adept for problem sizes which would not occupy the hardware well if executed using its counterparts, FlashAttention-2 and FlashDecoding. To enable such a smooth transition between tiles, the input tensor view is also different in LeanAttention compared to FlashAttention-2. This requires a constant stride moving between different heads as we transition from the LeanTile of one head to another requiring query, key, and value tensors be of the shape  $(\text{batch\_size}, \text{heads}, \text{query}/\text{ctx\_length}, \text{head\_dim})$

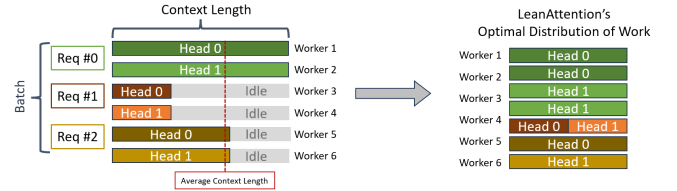


Fig. 6. Illustrative diagram of LeanAttention’s optimal distribution of work in the ragged batching case. Each SM receives equal amount of LeanTiles.

compared to FlashAttention-2’s requirement of  $(\text{batch\_size}, \text{query}/\text{ctx\_length}, \text{heads}, \text{head\_dim})$ .

With this design of execution, we must point out that LeanAttention behaves as a versatile attention partitioning mechanism which generalizes to FlashAttention-2 in the case where the number of output tiles is equal to grid size, and generalizes to FlashDecoding when grid size is an even multiple of number of output tiles. Finally, for all other cases (most common) LeanAttention efficiently distributes the work across the compute resources available in the system. Thus, LeanAttention will either always perform better or the same as FlashAttention-2 and FlashDecoding.

**Lean Ragged Batching.** For the special case of dealing with unequal context lengths within a batch of requests, the number of LeanTiles per request becomes unique, resulting in a total workload that is smaller compared to the non-ragged case. To account for this difference, ragged key/value tensor inputs to LeanAttention are first prepared with unpadded dimensions of  $(\text{NumHeads}, \text{TotalContextLength}, \text{HeadDim})$ , where  $\text{TotalContextLength}$  is the sum of all distinct context lengths within the batch. As seen, the batch dimension is eliminated, and both batch indices and the true context lengths of each request are tracked through pointers to a cumulative sequence lengths array for each input tensor. These pointers have size  $(\text{BatchSize} + 1)$  each, introducing minimal memory overhead.

With the total workload of LeanTiles correctly determined, Lean ragged batching functions identically to the non-ragged case as shown in . The workload is distributed evenly across the grid as shown in Figure 6, ensuring each CTA receives the same number of LeanTiles to process. The range of LeanTile iterations assigned to each CTA is mapped contiguously in a



Heads  $\rightarrow$  TotalContextLength linearization and partial outputs for each head are consolidated in the same manner as in the non-ragged case.

#### D. Execution Flow

Algorithm 2 depicts a Stream-K style execution of LeanAttention. For a fixed grid size  $G$ , CTAs are launched and given equal amount of LeanTile iterations to work with (Line 7). Each CTA computes the LeanTile() subroutine for every distinct output tile that comes under its boundaries (Line 16). Figure 5 shows the execution flow of a single CTA computing partial and final attention outputs for the assigned heads.

The unique reduction phase of LeanAttention, characterized by its softmax re-scaling operator, is performed by the host CTA block (Lines 24-40). A host CTA (Line 17) is the CTA responsible for computing the first ever LeanTile for a given output tile, and it behaves as the reducing CTA during parallel reduction of partial un-scaled outputs.

All non-host CTAs will share their partials through a store to global memory and signal their arrival (Lines 20-23). On the other hand, a host block, which is a non-finishing block (Lines 24-25), needs to wait for other contributing peer CTA blocks to signal their completion (Line 28) and then proceed to carry out the reduction (Lines 29-35).

A CTA that is computing the attention for a head exclusively completes all the LeanTile iterations for its output tile in a single CTA and so can directly store its results from the register file to global memory (Line 38-39) without any need for further reduction.

#### V. EVALUATION METHODOLOGY

**Implementation.** We implement LeanAttention using the CUTE abstractions [1]–[3] provided by Nvidia’s CUTLASS library [6]. For comparative measurements we utilize FlashAttention-2’s implementation of FlashDecoding as it is available on their Github repository [5]<sup>2</sup> and FlashInfer’s implementation from their Github repository [39]<sup>3</sup>. For the end-to-end inference results we use OPT models as available in the HuggingFace Transformers repository [38] and modify them to allow execution via LeanAttention wherever necessary. Note that optimizations such as FlashAttention-3 [32] are orthogonal to this work and targeted specifically for H100. The core computation of LeanAttention can adopt FlashAttention-3’s optimizations for further benefits on the Nvidia-H100 GPU [28].

**System.** We benchmark the attention mechanisms on Nvidia-A100-80GB-GPU [13] system with up to 8 GPUs. We measure runtime using a single GPU as well as 8xGPUs for larger models and context lengths. A single A100 GPU consists of 108 streaming multiprocessors (SMs) with an 80GB HBM global memory. To demonstrate LeanAttention’s versatility across hardware architectures we benchmark it similarly on

---

#### Algorithm 2 Lean Attention

---

```

1: _shared_  $O[T_m, d]$ 
2: _shared_  $m[T_m, 1]$ 
3: _shared_  $l[T_m, 1]$ 
4: Number of output tiles:  $C_m = \lceil N_q/T_m \rceil$ 
5: Number of iterations for each output tile:  $C_n = \lceil N_k/T_n \rceil$ 

6: Total number of iterations:  $I = C_m C_n$ 
7: Number of iterations per CTA:  $I_G = I/G$ 
8: fork CTA $g$  in  $G$  do
9:   cta_start =  $g \cdot I_G$  and cta_end = cta_start +  $I_G$ 
10:  for iter = cta_start to cta_end do
11:    Index of current output tile: tile_idx = iter /  $C_n$ 
12:    tile_iter = tile_idx  $\times C_n$ 
13:    tile_iter_end = tile_iter +  $C_n$ 
14:    local_iter = iter - tile_iter
15:    local_iter_end = min(tile_iter_end, cta_end) - tile_iter
16:    O, m, l = LeanTile(tile_idx, local_iter, local_iter_end)
17:    host-block if: iter == tile_iter
18:    finishing-block if: cta_end  $\geq$  tile_iter_end
19:    if !(host-block) then
20:      StorePartials(Op[g], O)
21:      StorePartials(mp[g], m)
22:      StorePartials(lp[g], l)
23:      Signal(flags[g])
24:    else
25:      if !(finishing-block) then
26:        last_cta = tile_iter_end /  $C_n$ 
27:        for cta = (g + 1) to last_cta do
28:          Wait(flags[cta])
29:           $m_{cta} = \text{LoadPartials}(mp[cta])$ 
30:           $l_{cta} = \text{LoadPartials}(lp[cta])$ 
31:           $O_{cta} = \text{LoadPartials}(Op[cta])$ 
32:           $m^{new} = \max(m_{cta}, m)$ 
33:           $l^{new} = e^{m_{cta} - m^{new}} l_{cta} + e^{m - m^{new}} l$ 
34:           $O^{new} = e^{m_{cta} - m^{new}} O_{cta} + e^{m - m^{new}} O$ 
35:          Update  $m = m_i^{new}, l = l_i^{new}$ 
36:        end for
37:      end if
38:      Write  $O = \text{diag}(l)^{-1} O$  to GMEM.
39:      Write  $L = m + \log(l)$  to GMEM.
40:    end if
41:    iter = tile_iter_end
42:  end for
43: join

```

---

a single Nvidia-H100-SXM-80GB-GPU [28] which has 132 streaming multiprocessors and an 80GB HBM global memory.

**Batching.** The evaluations assume the same context length for all queries in a batch working in tandem with batching techniques such as Orca [41]. However, LeanAttention supports varied context length execution including heterogeneous batching such as prefill queries with decode.

**Multi-GPU Tensor Parallelism.** We utilize Tensor Parallelism for the multi-GPU measurements to reflect the large model

<sup>2</sup>Version 2.5.6

<sup>3</sup>Version 0.1.6 - Note that we increase the number of heads in the single batch decoding API to simulate batch size as the batch API performance is too low

# 1x Nvidia-A100-80GB

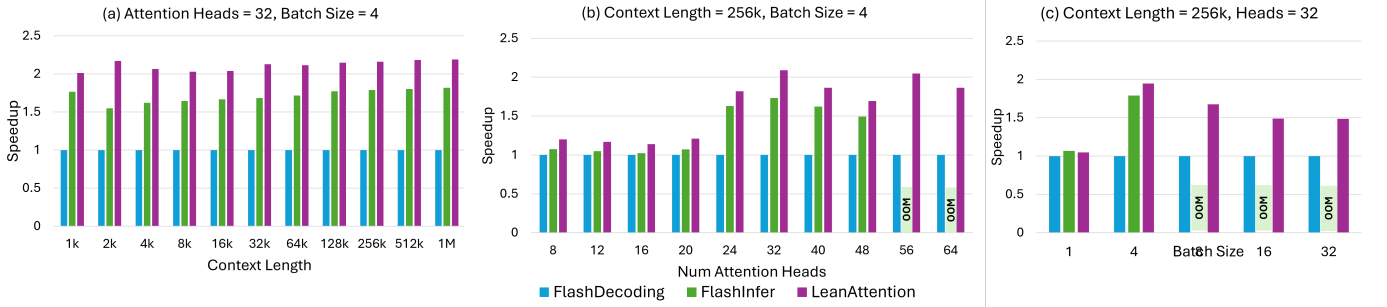


Fig. 7. Speedup of LA compared to state-of-the-art Attention execution mechanisms at different context lengths, batch sizes and attention heads with head dimension = 64 on a single Nvidia-A100-80GB GPU.

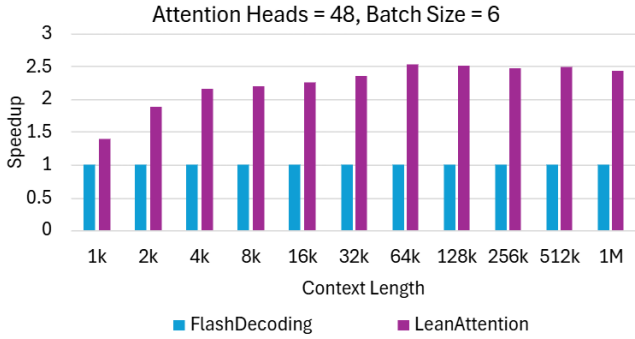


Fig. 8. Speedup of LA compared to state-of-the-art Attention execution mechanisms at varying context lengths, at a fixed batch size and attention heads with head dimension = 64 on a single Nvidia-H100-SXM-80GB GPU.

executions. Since FlashDecoding does not support Tensor Parallelism, we scale the implementation to the total number of SMs available in the system.

**Attention Mechanism.** In addition to FlashDecoding [4], we also benchmark FlashInfer for comparison against LeanAttention. FlashInfer has two kernel implementations for decoding: one for single-request decoding and another for batched-request decoding. The single-request decode kernel implements Fixed-Split, while the batched decode kernel implements PagedAttention. In our experiments, we benchmarked the batched decode kernel using a page size of 16. We’ve observed no impact of page size on FlashInfer’s latency, a finding consistent with their blog post [39]. Moreover, since FI reserves extra GPU memory to store buffers for managing its paged KV cache, we were unable to test it on certain large problem sizes, which we indicate as Out-of-Memory (“OOM”) errors in our evaluation figures. For the rest of the paper we refer to FlashDecoding as FD, FlashInfer as FI and LeanAttention as LA.

## VI. EVALUATION RESULTS

In this section, we evaluate the impact of LeanAttention (LA) at the attention operation-level as well as end-to-end inference performance.

### A. Benchmarking Attention - Decode-Phase

We benchmark the runtime of just the attention operation using the different mechanisms at varying context lengths,

number of attention heads, head dimensions (64:default and 128), and inference batch sizes on a single Nvidia A100-80GB GPU and a single Nvidia H100-SXM-80GB GPU.

**Increasing Context Length.** Figure 7(a) shows the speedup of different attention mechanisms for a model with 32 attention heads with a batch size of 4 on a single A100 GPU. LA delivers close to 2x speedup compared to FD even at smaller context lengths, reaching up to 2.18x speedup as the context lengths grows to 256k tokens. When context lengths exceed 16k, we observe more than 1.46x speedup over FI. Repeating a similar exercise on an H100 GPU, we observe the speedups of FD versus LA at a fixed batch size of 6 and 48 attention heads as shown in Figure 8(a). LA delivers more than 2x speedup even for contexts over 4k tokens, reaching upto a maximum of 2.52x speedup over FD at a 64k context length and 4.48x speedup over FI which more or less plateaus at the context lengths increase.

**Increasing Attention Heads.** Figure 7(b) and Figure 8(b) shows the speedup delivered by LA compared to FD and FI for models with an increasing number of heads. LA delivers comparable speedups to FD and FI at smaller model sizes. With fewer attention heads, FD’s fixed-split mechanism can distribute the workload as evenly as LA. However, as the number of heads increases, FD resorts to fewer splits per head, resulting in partially filled waves of attention on the SMs. In contrast, LA maintains even workload distribution at both small and large model sizes. The speedups over FD and FI vary depending on model size as depicted, with LA achieving an impressive 2.53x speedup over FD on an H100 GPU with 6 batched contexts of length 64k. Additionally, LA delivers more than 2x speedup over FD on an A100 GPU when the number of attention heads exceeds 24 for 4 batched contexts, each with a length of 256k. This shows that LA is able to scale well for both a small and large number of heads.

**Effect of Batching.** Figure 7(c) and Figure 8(c) shows the performance improvement of LA at varying batch sizes. As expected, we observe that LA gives comparable speedup to FD and FI. This is because both FD and FI are able to employ a higher number of splits at smaller batch sizes to occupy all the SMs in the GPU. However, as batch sizes increase, FD selects fewer splits per head. For instance, FD opts not to split at batch sizes above 4 in Figure 7(c) and Figure 8(c) because

### 8x Nvidia-A100-80GB

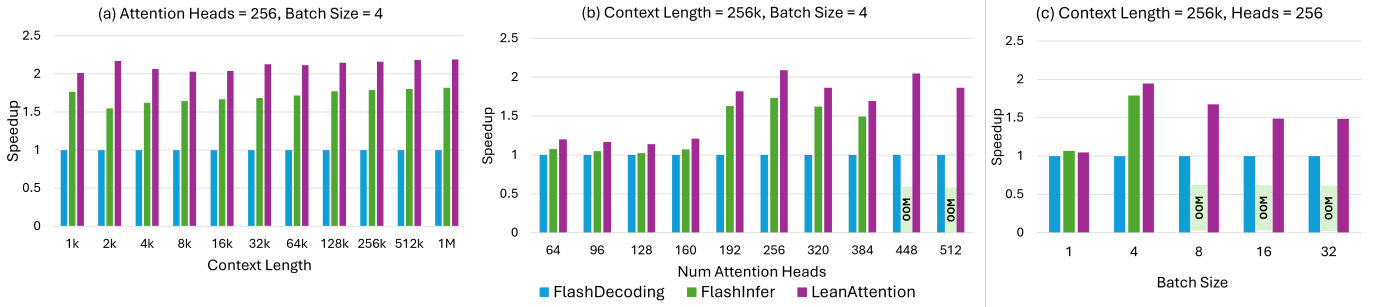


Fig. 9. Speedup of LA compared to state-of-the-art Attention execution mechanisms at different context lengths, batch sizes, attention heads with head dimension = 64 on an 8x Nvidia-A100-80GB system.

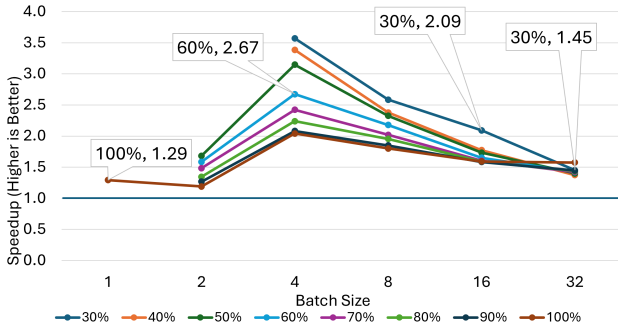


Fig. 10. Speedup offered by LA over FD at different batch sizes with heterogeneous context lengths. Batch context ratio(%) shows the ratio of average context length over maximum context length

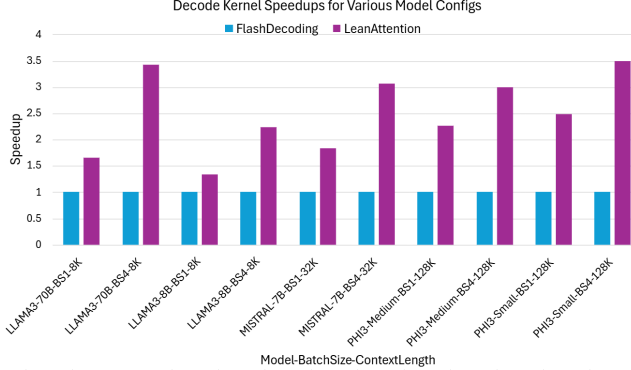


Fig. 11. Speedup offered by LA for decode attention across models, batch sizes, context lengths integrated with ONNXRT

the total number of heads in the batch exceeds the number of SMs available in the system. As a result, it behaves like vanilla FlashAttention-2, missing out on potential performance gains by leaving some SMs idle in its final, partially full wave. Consequently, LA achieves more than 1.5x speedup compared to FD through its stream-K-ed decomposition.

Overall, we benchmarked the system on more than a 1,000 samples with varying batch sizes, context lengths, and attention heads. On an A100 GPU, we observed an average speedup of 1.73x over FD (Max: 2.18x for 56 heads, batch size 2, 256k context) and an average of 3.42x over FI (Max: 5.66x for 12 heads, batch size 8, 512k context). On the H100 GPU, we

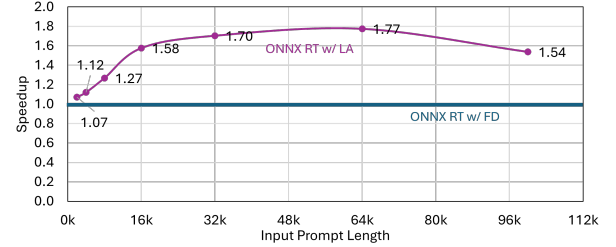


Fig. 12. End-to-End Speedup of LA compared to FD in ONNXRT running Phi-3 Medium model at different context lengths, batch size = 1, prompt size : output tokens = 8 : 1

recorded an average speedup of 1.52x over FD (Max: 2.53x for 48 heads, batch size 6, 64k context) and an average speedup of 3.63x over FI (Max: 4.59x for 56 heads, batch size 4, 128k context).

**Ragged Batching in Decode.** For the purpose of our evaluations, we quantify the heterogeneity of a ragged batch as the ratio of average context length to the maximum context length present in the batch. Figure 10 shows the speedup of LA over FD. We observe that as the heterogeneity of batch increases, LA delivers a higher speedup because of better distribution of work across SMs.

**Multi-GPU Execution** Repeating a similar benchmarking process on an 8xA100 GPU system, we vary the context lengths from 1k to 1M, with 256 attention heads at a batch size of 4 as shown in Figure 9(a). LeanAttention reaches a speedup of more than 2x even at smaller contexts. This is because parallelizing only over the batch and heads (total heads =  $256 \times 4 = 1024$ ) does not provide sufficient work for each SM (total SMs =  $8 \times 108 = 864$ ) as  $1024 - (9 \times 108) = 52$  SMs remain idle in the last wave. Furthermore, FD behaves identically to vanilla FlashAttention-2, opting for a split factor of 1. In contrast, LeanAttention computes attention in fully quantized waves for all problem sizes.

To observe this effect in greater detail, we evaluate across a varying number of attention heads in Figure 9(b) with a context length of 256k and batch size of 4. We observe comparable speedups to FD and FI at a smaller number of heads (64, 160). This is because at these dimensions, fixed-split is able to produce enough splits to occupy most of the SMs. We can also clearly see that FD resorts to vanilla execution when we increase the number of heads from 160 to 512. LA, on the

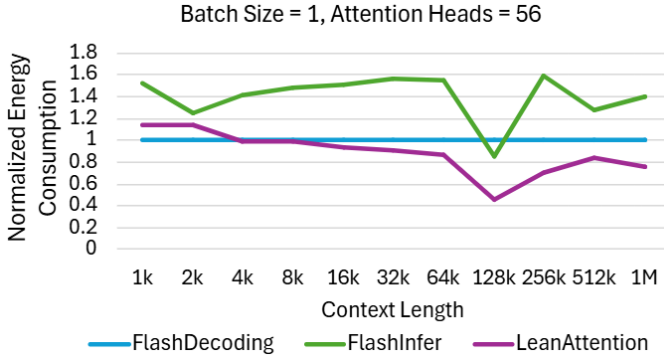


Fig. 13. Ratio of Energy consumed by attention kernel to energy consumed by FlashDecoding kernel of different attention mechanisms for batch size = 1, number of heads = 56, head dimension = 64 on a single Nvidia-A100-80GB GPU as measured using NVML APIs.

other hand scales well as we increase the number of heads, showcasing its hardware-aware scalable execution algorithm. With 256 heads, as we increase the batch size from 1 to 32, we can see that LA starts to outperform FA2 variants as batch size increases.

**Effect of Head Dimension.** Figure 11 shows the speedup offered by LA for models with a LLAMA-2, Mistral and Phi3-like config with a head dimension of 128. We utilize a 128-token wide LeanTile for decomposition of each problem instead of 256. We observe a similar trend in performance, where LA delivers a speedup of 3.5x compared to FD at 128k context length. Even at smaller context lengths of 8k tokens, we observed an improved performance of 1.34x over FD.

To summarize, LA not only outperforms FD at lower batch sizes, long context lengths, but also delivers better performance at higher batch sizes, and with higher number of attention heads. This is mainly due to the lean decomposition of the problem on the hardware compute resource. For cases where there are enough parallelizable dimensions, LA automatically generalizes to FA2-like execution. We can thus treat FA2’s execution algorithm as a special case of LA, which occurs depending on the optimal grid size that LeanAttention chooses depending on the hardware resources and LeanTile dimensions.

Fixed-split partitioning results in imbalanced workloads across the SMs, leaving many of them idle during the final stages of computation. This inefficiency makes fixed-split attention mechanisms energy-inefficient. As shown in Figure 13, the disparity in energy consumption between FlashDecoding, FlashInfer, and LeanAttention increases as context lengths grow over 128k. LeanAttention, with its well-balanced load partitioning strategy, ensures more consistent utilization of SMs, making it significantly more energy-efficient.

### B. End-to-End Inference Performance

We measure the end-to-end inference runtime using Phi-3 Medium model (with 40 attention heads) as shown in Figure 12 at different prompt sizes (total context = prompt tokens + tokens generated so far) with a prompt to output token ratio of 8:1. This includes the prefill-stage latency as well as the total runtime of decode-phase. LeanAttention offers a 1.12x

speedup with Phi-3 Medium as compared to FlashDecoding for first 1k output tokens. However, the LA offers a higher speedup as the output tokens increase beyond 16k delivering an average of 1.73x speedup compared to FA2. As we note, the inference-level runtime improvement delivered by LA will change heavily on the number of heads, total context length, batch size, etc.

## VII. CONCLUSION

The attention mechanism in transformer-based language models is a slow and memory hungry process. State-of-the-art optimization mechanisms, such as FlashAttention-2, FlashDecoding, and FlashInfer, have cleverly addressed this challenge; however, they fail to adapt to the computationally distinct phases of inference. We observe that these techniques fail at parallelization along the context length dimension during the decode phase of inference, resulting in low occupancy of the underlying hardware and slower inference. As state-of-the-art models continue to push the limits on supporting increasingly long context lengths and heterogeneous context-length batching [10], the importance of optimization techniques that effectively parallelize across this dimension is becoming increasingly critical.

To address this challenge, we propose *LeanAttention*, a scalable and generalized exact attention execution mechanism that ensures lower runtimes and almost 100% hardware occupancy during attention computation irrespective of problem size. LeanAttention leverages the associative property of softmax re-scaling, treating it as a reductive operator in a “stream-K”-style partitioning of the attention mechanism. This enables an efficient parallelized execution of the attention mechanism, particularly during the decode phase of inference, which traditionally suffers from longer runtimes and critically low hardware utilization.

Our measurements indicate that LeanAttention delivers an average speedup of 1.73x over FlashDecoding, with up to 2.18x speedup for a 256k context size. Notably, in a multi-GPU execution scenario with numerous attention heads, the speedup realized by LeanAttention continues to increase as context length grows. LeanAttention intelligently utilizes the underlying hardware, enabling efficient scaling of next-generation LLMs that leverage large context lengths.

## REFERENCES

- [1] “Cute layouts.” [https://github.com/NVIDIA/cutlass/blob/main/media/docs/cute/01\\_layout.md](https://github.com/NVIDIA/cutlass/blob/main/media/docs/cute/01_layout.md), [Accessed 19-04-2024].
- [2] “Cute tensors.” [https://github.com/NVIDIA/cutlass/blob/main/media/docs/cute/03\\_tensor.md](https://github.com/NVIDIA/cutlass/blob/main/media/docs/cute/03_tensor.md), [Accessed 19-04-2024].
- [3] “Cute’s support for matrix multiply-accumulate instructions.” [https://github.com/NVIDIA/cutlass/blob/main/media/docs/cute/0t\\_mma\\_atom.md](https://github.com/NVIDIA/cutlass/blob/main/media/docs/cute/0t_mma_atom.md), [Accessed 19-04-2024].
- [4] “Flashdecoding: Stanford CRFM — crfm.stanford.edu,” <https://crfm.stanford.edu/2023/10/12/flashdecoding.html>, [Accessed 22-04-2024].
- [5] “GitHub - Dao-AILab/flash-attention: Fast and memory-efficient exact attention — github.com,” <https://github.com/Dao-AILab/flash-attention>, [Accessed 19-04-2024].
- [6] “GitHub - NVIDIA/cutlass: CUDA Templates for Linear Algebra Sub-routines — github.com,” <https://github.com/NVIDIA/cutlass>, [Accessed 01-04-2024].
- [7] “Introducing ChatGPT — openai.com,” <https://openai.com/blog/chatgpt>, [Accessed 01-04-2024].
- [8] “Introducing the next generation of claude,” <https://www.anthropic.com/news/claude-3-family>, [Accessed 19-04-2024].
- [9] J. Achiam, S. Adler, S. Agarwal, L. Ahmad, I. Akkaya, F. L. Aleman, D. Almeida, J. Altschmidt, S. Altman, S. Anadkat et al., “Gpt-4 technical report,” *arXiv preprint arXiv:2303.08774*, 2023.
- [10] A. Agrawal, N. Kedia, A. Panwar, J. Mohan, N. Kwatra, B. S. Gulavani, A. Tumanov, and R. Ramjee, “Taming throughput-latency tradeoff in llm inference with sarathi-serve,” 2024. [Online]. Available: <https://arxiv.org/abs/2403.02310>
- [11] W. Brandon, A. Nrusimha, K. Qian, Z. Ankner, T. Jin, Z. Song, and J. Ragan-Kelley, “Striped attention: Faster ring attention for causal transformers,” *arXiv preprint arXiv:2311.09431*, 2023.
- [12] T. Brown, B. Mann, N. Ryder, M. Subbiah, J. D. Kaplan, P. Dhariwal, A. Neelakantan, P. Shyam, G. Sastry, A. Askell et al., “Language models are few-shot learners,” *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.
- [13] J. Choquette, E. Lee, R. Krashinsky, V. Balan, and B. Khailany, “3.2 the a100 datacenter gpu and ampere architecture,” in *2021 IEEE International Solid-State Circuits Conference (ISSCC)*, vol. 64. IEEE, 2021, pp. 48–50.
- [14] A. Chowdhery, S. Narang, J. Devlin, M. Bosma, G. Mishra, A. Roberts, P. Barham, H. W. Chung, C. Sutton, S. Gehrmann et al., “Palm: Scaling language modeling with pathways,” *Journal of Machine Learning Research*, vol. 24, no. 240, pp. 1–113, 2023.
- [15] T. Dao, “Flashattention-2: Faster attention with better parallelism and work partitioning,” *arXiv preprint arXiv:2307.08691*, 2023.
- [16] T. Dao, D. Fu, S. Ermon, A. Rudra, and C. Ré, “Flashattention: Fast and memory-efficient exact attention with io-awareness,” *Advances in Neural Information Processing Systems*, vol. 35, pp. 16 344–16 359, 2022.
- [17] Y. Fu, R. Panda, X. Niu, X. Yue, H. Hajishirzi, Y. Kim, and H. Peng, “Data engineering for scaling language models to 128k context,” *arXiv preprint arXiv:2402.10171*, 2024.
- [18] K. Hong, G. Dai, J. Xu, Q. Mao, X. Li, J. Liu, Y. Dong, Y. Wang et al., “Flashdecoding++: Faster large language model inference with asynchronization, flat gemm optimization, and heuristics,” *Proceedings of Machine Learning and Systems*, vol. 6, pp. 148–161, 2024.
- [19] G. Inc., “An important next step on our AI journey — blog.google,” <https://blog.google/technology/ai/bard-google-ai-search-updates/>, [Accessed 31-03-2024].
- [20] A. Ivanov, N. Dryden, T. Ben-Nun, S. Li, and T. Hoefer, “Data movement is all you need: A case study on optimizing transformers,” *Proceedings of Machine Learning and Systems*, vol. 3, pp. 711–732, 2021.
- [21] Z. Jia and P. Van Sandt, “Dissecting the ampere gpu architecture via microbenchmarking,” in *GPU Technology Conference*, 2021.
- [22] J. D. M.-W. C. Kenton and L. K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in *Proceedings of naacL-HLT*, vol. 1, 2019, p. 2.
- [23] S. Kim, C. Hooper, T. Wattanawong, M. Kang, R. Yan, H. Genc, G. Dinh, Q. Huang, K. Keutzer, M. W. Mahoney et al., “Full stack optimization of transformer inference: a survey,” *arXiv preprint arXiv:2302.14017*, 2023.
- [24] T. Li, G. Zhang, Q. D. Do, X. Yue, and W. Chen, “Long-context llms struggle with long in-context learning,” *arXiv preprint arXiv:2404.02060*, 2024.
- [25] Y. Li, S. Bubeck, R. Eldan, A. Del Giorno, S. Gunasekar, and Y. T. Lee, “Textbooks are all you need ii: phi-1.5 technical report,” *arXiv preprint arXiv:2309.05463*, 2023.
- [26] H. Liu, M. Zaharia, and P. Abbeel, “Ring attention with blockwise transformers for near-infinite context,” *arXiv preprint arXiv:2310.01889*, 2023.
- [27] S. Liu, G. Tao, Y. Zou, D. Chow, Z. Fan, K. Lei, B. Pan, D. Sylvester, G. Kielian, and M. Saligane, Eds., *ConSmax: Hardware-Friendly Alternative Softmax with Learnable Parameters*, 2024. [Online]. Available: <https://arxiv.org/abs/2402.10930>
- [28] W. Luo, R. Fan, Z. Li, D. Du, Q. Wang, and X. Chu, “Benchmarking and dissecting the nvidia hopper gpu architecture,” 2024. [Online]. Available: <https://arxiv.org/abs/2402.13499>
- [29] M. Milakov and N. Gimelshein, “Online normalizer calculation for softmax,” *arXiv preprint arXiv:1805.02867*, 2018.
- [30] M. Osama, D. Merrill, C. Cecka, M. Garland, and J. D. Owens, “Stream-k: Work-centric parallel decomposition for dense matrix-matrix multiplication on the gpu,” in *Proceedings of the 28th ACM SIGPLAN Annual Symposium on Principles and Practice of Parallel Programming*, 2023, pp. 429–431.
- [31] R. Pope, S. Douglas, A. Chowdhery, J. Devlin, J. Bradbury, J. Heek, K. Xiao, S. Agrawal, and J. Dean, “Efficiently scaling transformer inference,” *Proceedings of Machine Learning and Systems*, vol. 5, 2023.
- [32] J. Shah, G. Bikshandi, Y. Zhang, V. Thakkar, P. Ramani, and T. Dao, “Flashattention-3: Fast and accurate attention with asynchrony and low-precision,” 2024. [Online]. Available: <https://arxiv.org/abs/2407.08608>
- [33] J. Spataro and M. Inc., “Introducing Microsoft 365 Copilot – your copilot for work - The Official Microsoft Blog — blogs.microsoft.com,” <https://blogs.microsoft.com/blog/2023/03/16/introducing-microsoft-365-copilot-your-copilot-for-work/>, [Accessed 31-03-2024].
- [34] J. R. Stevens, R. Venkatesan, S. Dai, B. Khailany, and A. Raghunathan, “Softmax: Hardware/software co-design of an efficient softmax for transformers,” in *2021 58th ACM/IEEE Design Automation Conference (DAC)*. IEEE, 2021, pp. 469–474.
- [35] H. Touvron, L. Martin, K. Stone, P. Albert, A. Almahairi, Y. Babaei, N. Bashlykov, S. Batra, P. Bhargava, S. Bhosale et al., “Llama 2: Open foundation and fine-tuned chat models,” *arXiv preprint arXiv:2307.09288*, 2023.
- [36] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, “Attention is all you need,” *Advances in neural information processing systems*, vol. 30, 2017.
- [37] S. Williams, A. Waterman, and D. Patterson, “Roofline: an insightful visual performance model for multicore architectures,” *Communications of the ACM*, vol. 52, no. 4, pp. 65–76, 2009.
- [38] T. Wolf, L. Debut, V. Sanh, J. Chaumond, C. Delangue, A. Moi, P. Cistac, T. Rault, R. Louf, M. Funtowicz et al., “Huggingface’s transformers: State-of-the-art natural language processing,” *arXiv preprint arXiv:1910.03771*, 2019.
- [39] Z. Ye, L. Chen, R. Lai, Y. Zhao, S. Zheng, J. Shao, B. Hou, H. Jin, Y. Zuo, L. Yin, T. Chen, and L. Ceze, “Accelerating self-attentions for llm serving with flashinfer,” February 2024. [Online]. Available: <https://flashinfer.ai/2024/02/02/introduce-flashinfer.html>
- [40] G.-I. Yu, J. S. Jeong, G.-W. Kim, S. Kim, and B.-G. Chun, “Orca: A distributed serving system for {Transformer-Based} generative models,” in *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*, 2022, pp. 521–538.
- [41] G.-I. Yu, J. S. Jeong, G.-W. Kim, S. Kim, and B.-G. Chun, “Orca: A distributed serving system for Transformer-Based generative models,” in *16th USENIX Symposium on Operating Systems Design and Implementation (OSDI 22)*. Carlsbad, CA: USENIX Association, Jul. 2022, pp. 521–538. [Online]. Available: <https://www.usenix.org/conference/osdi22/presentation/yu>
- [42] P. Zhang, Z. Liu, S. Xiao, N. Shao, Q. Ye, and Z. Dou, “Soaring from 4k to 400k: Extending llm’s context with activation beacon,” *arXiv preprint arXiv:2401.03462*, 2024.
- [43] S. Zhang, S. Roller, N. Goyal, M. Artetxe, M. Chen, S. Chen, C. Dewan, M. Diab, X. Li, X. V. Lin et al., “Opt: Open pre-trained transformer language models,” *arXiv preprint arXiv:2205.01068*, 2022.