FLASHATTENTION: Fast and Memory-Efficient Exact Attention with IO-Awareness

Tri Dao[†], Daniel Y. Fu[†], Stefano Ermon[†], Atri Rudra[‡], and Christopher Ré[†]

July 1, 2023

Abstract

Transformers are slow and memory-hungry on long sequences, since the time and memory complexity of self-attention are quadratic in sequence length. Approximate attention methods have attempted to address this problem by trading off model quality to reduce the compute complexity, but often do not achieve wall-clock speedup. We argue that a missing principle is making attention algorithms IOaware—accounting for reads and writes between levels of GPU memory. We propose FlashAttention, an IO-aware exact attention algorithm that uses tiling to reduce the number of memory reads/writes between GPU high bandwidth memory (HBM) and GPU on-chip SRAM. We analyze the IO complexity of FlashAttention, showing that it requires fewer HBM accesses than standard attention, and is optimal for a range of SRAM sizes. We also extend FlashAttention to block-sparse attention, yielding an approximate attention algorithm that is faster than any existing approximate attention method. FLASHATTENTION trains Transformers faster than existing baselines: 15% end-to-end wall-clock speedup on BERT-large (seq. length 512) compared to the MLPerf 1.1 training speed record, 3x speedup on GPT-2 (seq. length 1K), and 2.4× speedup on long-range arena (seq. length 1K-4K). FLASHATTENTION and block-sparse FlashAttention enable longer context in Transformers, yielding higher quality models (0.7 better perplexity on GPT-2 and 6.4 points of lift on long-document classification) and entirely new capabilities: the first Transformers to achieve better-than-chance performance on the Path-X challenge (seq. length 16K, 61.4% accuracy) and Path-256 (seq. length 64K, 63.1% accuracy).

1 Introduction

Transformer models [86] have emerged as the most widely used architecture in applications such as natural language processing and image classification. Transformers have grown larger [5] and deeper [87], but equipping them with longer context remains difficult [83], since the self-attention module at their heart has time and memory complexity quadratic in sequence length. An important question is whether making attention faster and more memory-efficient can help Transformer models address their runtime and memory challenges for long sequences.

Many approximate attention methods have aimed to reduce the compute and memory requirements of attention. These methods range from sparse-approximation [53, 77] to low-rank approximation [13, 52, 88], and their combinations [3, 9, 96]. Although these methods reduce the compute requirements to linear or near-linear in sequence length, many of them do not display wall-clock speedup against standard attention and have not gained wide adoption. One main reason is that they focus on FLOP reduction (which may not correlate with wall-clock speed) and tend to ignore overheads from memory access (IO).

In this paper, we argue that a missing principle is making attention algorithms *IO-aware* [1]—that is, carefully accounting for reads and writes to different levels of fast and slow memory (e.g., between fast GPU on-chip SRAM and relatively slow GPU high bandwidth memory, or HBM [47], Figure 1 left). On modern

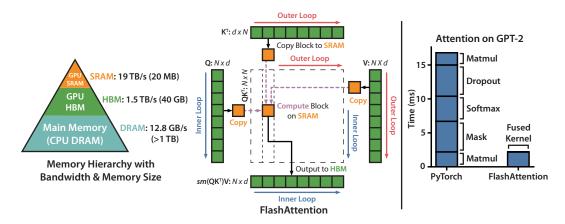


Figure 1: **Left:** FLASHATTENTION uses tiling to prevent materialization of the large $N \times N$ attention matrix (dotted box) on (relatively) slow GPU HBM. In the outer loop (red arrows), FLASHATTENTION loops through blocks of the **K** and **V** matrices and loads them to fast on-chip SRAM. In each block, FLASHATTENTION loops over blocks of **Q** matrix (blue arrows), loading them to SRAM, and writing the output of the attention computation back to HBM. **Right:** Speedup over the PyTorch implementation of attention on GPT-2. FLASHATTENTION does not read and write the large $N \times N$ attention matrix to HBM, resulting in an 7.6× speedup on the attention computation.

GPUs, compute speed has out-paced memory speed [64, 65, 66], and most operations in Transformers are bottlenecked by memory accesses [45]. IO-aware algorithms have been critical for similar memory-bound operations, when reading and writing data can account for a large portion of the runtime—such as database joins [74], image processing [73], numerical linear algebra [4], and more [42, 89]. However, common Python interfaces to deep learning such as PyTorch and Tensorflow do not allow fine-grained control of memory access.

We propose FLASHATTENTION, a new attention algorithm that computes exact attention with far fewer memory accesses. Our main goal is to avoid reading and writing the attention matrix to and from HBM. This requires (i) computing the softmax reduction without access to the whole input (ii) not storing the large intermediate attention matrix for the backward pass. We apply two well-established techniques to address these challenges. (i) We restructure the attention computation to split the input into blocks and make several passes over input blocks, thus incrementally performing the softmax reduction (also known as **tiling**). (ii) We store the softmax normalization factor from the forward pass to quickly **recompute** attention on-chip in the backward pass, which is faster than the standard approach of reading the intermediate attention matrix from HBM. We implement Flashattention in CUDA to achieve fine-grained control over memory access and fuse all the attention operations into one GPU kernel. Even with the increased FLOPs due to recomputation, our algorithm both **runs faster** (up to 7.6x on GPT-2 [70], Figure 1 right) and **uses less memory**—linear in sequence length—than standard attention, thanks to the massively reduced amount of HBM access.

We analyze the IO complexity [1] of FLASHATTENTION, proving that it requires $O(N^2d^2M^{-1})$ HBM accesses where d is the head dimension and M is the size of SRAM, as compared to $\Omega(Nd+N^2)$ of standard attention. For typical values of d and M, FLASHATTENTION requires many times fewer HBM accesses compared to standard attention (up to $9\times$ fewer, as shown in Fig. 2). Moreover, we provide a lower bound, showing that no exact attention algorithm can asymptotically improve on the number of HBM accesses over all SRAM sizes.

We also show that FlashAttention can serve as a useful primitive for realizing the potential of approximate attention algorithms by overcoming their issues with memory access overhead. As a proof of concept, we implement block-sparse FlashAttention, a sparse attention algorithm that is 2-4× faster than even FlashAttention, scaling up to sequence length of 64k. We prove that block-sparse FlashAttention has better IO complexity than FlashAttention by a factor proportional to the sparsity ratio. We discuss further extensions to other operations (attention on multi-GPU, kernel regression, block-sparse matrix

multiply) in Section 5. We open-source FLASHATTENTION to make it easier to build on this primitive¹.

We empirically validate that FlashAttention speeds up model training and improves model quality by modeling longer context. We also benchmark the runtime and memory footprint of FlashAttention and block-sparse FlashAttention compared to prior attention implementations.

- Faster Model Training. FLASHATTENTION trains Transformer models faster in wall-clock time. We train BERT-large (seq. length 512) 15% faster than the training speed record in MLPerf 1.1 [60], GPT2 (seq. length 1K) 3× faster than baseline implementations from HuggingFace [91] and Megatron-LM [80], and long-range arena (seq. length 1K-4K) 2.4× faster than baselines.
- Higher Quality Models. FlashAttention scales Transformers to longer sequences, which improves their quality and enables new capabilities. We observe a 0.7 improvement in perplexity on GPT-2 and 6.4 points of lift from modeling longer sequences on long-document classification [14]. FlashAttention enables the first Transformer that can achieve better-than-chance performance on the Path-X [83] challenge, solely from using a longer sequence length (16K). Block-sparse FlashAttention enables a Transformer to scale to even longer sequences (64K), resulting in the first model that can achieve better-than-chance performance on Path-256.
- Benchmarking Attention. FlashAttention is up to 3× faster than the standard attention implementation across common sequence lengths from 128 to 2K and scales up to 64K. Up to sequence length of 512, FlashAttention is both faster and more memory-efficient than any existing attention method, whereas for sequence length beyond 1K, some approximate attention methods (e.g., Linformer) start to become faster. On the other hand, block-sparse FlashAttention is faster than all existing approximate attention methods that we know of.

2 Background

We provide some background on the performance characteristics of common deep learning operations on modern hardware (GPUs). We also describe the standard implementation of attention.

2.1 Hardware Performance

We focus here on GPUs. Performance on other hardware accelerators are similar [48, 50].

GPU Memory Hierarchy. The GPU memory hierarchy (Fig. 1 left) comprises multiple forms of memory of different sizes and speeds, with smaller memory being faster. As an example, the A100 GPU has 40-80GB of high bandwidth memory (HBM) with bandwidth 1.5-2.0TB/s and 192KB of on-chip SRAM per each of 108 streaming multiprocessors with bandwidth estimated around 19TB/s [46, 47]. The on-chip SRAM is an order of magnitude faster than HBM but many orders of magnitude smaller in size. As compute has gotten faster relative to memory speed [64, 65, 66], operations are increasingly bottlenecked by memory (HBM) accesses. Thus exploiting fast SRAM becomes more important.

Execution Model. GPUs have a massive number of threads to execute an operation (called a kernel). Each kernel loads inputs from HBM to registers and SRAM, computes, then writes outputs to HBM.

Performance characteristics. Depending on the balance of computation and memory accesses, operations can be classified as either compute-bound or memory-bound. This is commonly measured by the *arithmetic intensity* [89], which is the number of arithmetic operations per byte of memory access.

- 1. Compute-bound: the time taken by the operation is determined by how many arithmetic operations there are, while time accessing HBM is much smaller. Typical examples are matrix multiply with large inner dimension, and convolution with large number of channels.
- 2. Memory-bound: the time taken by the operation is determined by the number of memory accesses, while time spent in computation is much smaller. Examples include most other operations: elementwise (e.g., activation, dropout), and reduction (e.g., sum, softmax, batch norm, layer norm).

Kernel fusion. The most common approach to accelerate memory-bound operations is kernel fusion: if there are multiple operations applied to the same input, the input can be loaded once from HBM, instead of multiple times for each operation. Compilers can automatically fuse many elementwise operations [55, 68, 78].

 $^{{}^1\}mathrm{FlashAttention}\ code\ is\ available\ at\ \mathtt{https://github.com/HazyResearch/flash-attention}$

However, in the context of model training, the intermediate values still need to be written to HBM to save for the backward pass, reducing the effectiveness of naive kernel fusion.

2.2 Standard Attention Implementation

Given input sequences $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ where N is the sequence length and d is the head dimension, we want to compute the attention output $\mathbf{O} \in \mathbb{R}^{N \times d}$:

$$\mathbf{S} = \mathbf{Q}\mathbf{K}^{\mathsf{T}} \in \mathbb{R}^{N \times N}, \quad \mathbf{P} = \operatorname{softmax}(\mathbf{S}) \in \mathbb{R}^{N \times N}, \quad \mathbf{O} = \mathbf{P}\mathbf{V} \in \mathbb{R}^{N \times d},$$

where softmax is applied row-wise.

Standard attention implementations materialize the matrices **S** and **P** to HBM, which takes $O(N^2)$ memory. Often $N \gg d$ (e.g., for GPT2, N = 1024 and d = 64). We describe the standard attention implementation in Algorithm 0. As some or most of the operations are memory-bound (e.g., softmax), the large number of memory accesses translates to slow wall-clock time.

This problem is exacerbated by other elementwise operations applied to the attention matrix, such as masking applied to \mathbf{S} or dropout applied to \mathbf{P} . As a result, there have been many attempts to fuse several elementwise operations, such as fusing masking with softmax [80].

In Section 3.2, we will show that the standard attention implementation performs HBM accesses quadratic in the sequence length N. We also compare the number of FLOPs and number of HBM accesses of standard attention and of our method (FLASHATTENTION).

Algorithm 0 Standard Attention Implementation

Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM.

- 1: Load \mathbf{Q}, \mathbf{K} by blocks from HBM, compute $\mathbf{S} = \mathbf{Q}\mathbf{K}^{\mathsf{T}}$, write \mathbf{S} to HBM.
- 2: Read **S** from HBM, compute P = softmax(S), write **P** to HBM.
- 3: Load **P** and **V** by blocks from HBM, compute $\mathbf{O} = \mathbf{PV}$, write \mathbf{O} to HBM.
- 4: Return **O**.

3 FLASHATTENTION: Algorithm, Analysis, and Extensions

We show how to compute exact attention with fewer HBM reads/writes and without storing large intermediate matrices for the backward pass. This yields an attention algorithm that is both memory efficient and faster in wall-clock time. We analyze its IO complexity, showing that our method requires much fewer HBM accesses compared to standard attention. We further show that FlashAttention can serve as a useful primitive by extending it to handle block-sparse attention.

We focus here on the forward pass for ease of exposition; Appendix B contains details for the backward.

3.1 An Efficient Attention Algorithm With Tiling and Recomputation

Given the inputs $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ in HBM, we aim to compute the attention output $\mathbf{O} \in \mathbb{R}^{N \times d}$ and write it to HBM. Our goal is to reduce the amount of HBM accesses (to sub-quadratic in N).

We apply two established techniques (tiling, recomputation) to overcome the technical challenge of computing exact attention in sub-quadratic HBM accesses. We describe this in Algorithm 1. The main idea is that we split the inputs $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ into blocks, load them from slow HBM to fast SRAM, then compute the attention output with respect to those blocks. By scaling the output of each block by the right normalization factor before adding them up, we get the correct result at the end.

Tiling. We compute attention by blocks. Softmax couples columns of **K**, so we decompose the large softmax with scaling [53, 62, 69]. For numerical stability, the softmax of vector $x \in \mathbb{R}^B$ is computed:

$$m(x) := \max_i \ x_i, \quad f(x) := \left[e^{x_1 - m(x)} \ \dots \ e^{x_B - m(x)} \right], \quad \ell(x) := \sum_i f(x)_i, \quad \operatorname{softmax}(x) := \frac{f(x)}{\ell(x)}.$$

For vectors $x^{(1)}, x^{(2)} \in \mathbb{R}^B$, we can decompose the softmax of the concatenated $x = [x^{(1)}, x^{(2)}] \in \mathbb{R}^{2B}$ as:

$$\begin{split} m(x) &= m(\left[x^{(1)} \ x^{(2)}\right]) = \max(m(x^{(1)}), m(x^{(2)})), \quad f(x) = \left[e^{m(x^{(1)}) - m(x)} f(x^{(1)}) \quad e^{m(x^{(2)}) - m(x)} f(x^{(2)})\right], \\ \ell(x) &= \ell(\left[x^{(1)} \ x^{(2)}\right]) = e^{m(x^{(1)}) - m(x)} \ell(x^{(1)}) + e^{m(x^{(2)}) - m(x)} \ell(x^{(2)}), \quad \text{softmax}(x) = \frac{f(x)}{\ell(x)}. \end{split}$$

Therefore if we keep track of some extra statistics $(m(x), \ell(x))$, we can compute softmax one block at a time.² We thus split the inputs $\mathbf{Q}, \mathbf{K}, \mathbf{V}$ into blocks (Algorithm 1 line 3), compute the softmax values along with extra statistics (Algorithm 1 line 10), and combine the results (Algorithm 1 line 12).

Recomputation. One of our goals is to not store $O(N^2)$ intermediate values for the backward pass. The backward pass typically requires the matrices $S, P \in \mathbb{R}^{N \times N}$ to compute the gradients with respect to Q, K, V. However, by storing the output O and the softmax normalization statistics (m, ℓ) , we can recompute the attention matrix O and O easily in the backward pass from blocks of O and O in SRAM. This can be seen as a form of selective gradient checkpointing [10, 36]. While gradient checkpointing has been suggested to reduce the maximum amount of memory required [69], all implementations (that we know off) have to trade speed for memory. In contrast, even with more FLOPs, our recomputation speeds up the backward pass due to reduced HBM accesses (Fig. 2). The full backward pass description is in Appendix B.

Implementation details: Kernel fusion. Tiling enables us to implement our algorithm in one CUDA kernel, loading input from HBM, performing all the computation steps (matrix multiply, softmax, optionally masking and dropout, matrix multiply), then write the result back to HBM (masking and dropout in Appendix B). This avoids repeatedly reading and writing of inputs and outputs from and to HBM.

Algorithm 1 FlashAttention

```
Require: Matrices \mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d} in HBM, on-chip SRAM of size M.
  1: Set block sizes B_c = \left\lceil \frac{M}{4d} \right\rceil, B_r = \min\left(\left\lceil \frac{M}{4d} \right\rceil, d\right).

2: Initialize \mathbf{Q} = (0)_{N \times d} \in \mathbb{R}^{N \times d}, \ell = (0)_N \in \mathbb{R}^N, m = (-\infty)_N \in \mathbb{R}^N in HBM.

3: Divide \mathbf{Q} into T_r = \left\lceil \frac{N}{B_r} \right\rceil blocks \mathbf{Q}_1, \dots, \mathbf{Q}_{T_r} of size B_r \times d each, and divide \mathbf{K}, \mathbf{V} in to T_c = \left\lceil \frac{N}{B_c} \right\rceil blocks
          \mathbf{K}_1, \ldots, \mathbf{K}_{T_c} and \mathbf{V}_1, \ldots, \mathbf{V}_{T_c}, of size B_c \times d each.
   4: Divide \mathbf{0} into T_r blocks \mathbf{0}_i, \dots, \mathbf{0}_{T_r} of size B_r \times d each, divide \ell into T_r blocks \ell_i, \dots, \ell_{T_r} of size B_r each,
          divide m into T_r blocks m_1, \ldots, m_{T_r} of size B_r each.
         for 1 \le j \le T_c do
                Load \mathbf{K}_j, \mathbf{V}_j from HBM to on-chip SRAM.
                for 1 \le i \le T_r do
   7:
                      Load \mathbf{Q}_i, \mathbf{O}_i, \ell_i, m_i from HBM to on-chip SRAM.
   8:
                      On chip, compute \mathbf{S}_{ij} = \mathbf{Q}_i \mathbf{K}_i^T \in \mathbb{R}^{B_r \times B_c}.
  9:
                      On chip, compute \tilde{m}_{ij} = \text{rowmax}(\mathbf{S}_{ij}) \in \mathbb{R}^{B_r}, \tilde{\mathbf{P}}_{ij} = \exp(\mathbf{S}_{ij} - \tilde{m}_{ij}) \in \mathbb{R}^{B_r \times B_c} (pointwise), \tilde{\ell}_{ij} = \exp(\mathbf{S}_{ij} - \tilde{m}_{ij})
 10:
                      \operatorname{rowsum}(\tilde{\mathbf{P}}_{i\,i}) \in \mathbb{R}^{B_r}.
                      On chip, compute m_i^{\text{new}} = \max(m_i, \tilde{m}_{ij}) \in \mathbb{R}^{B_r}, \ell_i^{\text{new}} = e^{m_i - m_i^{\text{new}}} \ell_i + e^{\tilde{m}_{ij} - m_i^{\text{new}}} \tilde{\ell}_{ij} \in \mathbb{R}^{B_r}. Write \mathbf{O}_i \leftarrow \operatorname{diag}(\ell_i^{\text{new}})^{-1} (\operatorname{diag}(\ell_i) e^{m_i - m_i^{\text{new}}} \mathbf{O}_i + e^{\tilde{m}_{ij} - m_i^{\text{new}}} \tilde{\mathbf{P}}_{ij} \mathbf{V}_j) to HBM. Write \ell_i \leftarrow \ell_i^{\text{new}}, m_i \leftarrow m_i^{\text{new}} to HBM.
 11:
 12:
 13:
                end for
 14:
 15: end for
16: Return O
```

We show FlashAttention's correctness, runtime, and memory requirement (proof in Appendix C).

Theorem 1. Algorithm 1 returns $\mathbf{O} = \operatorname{softmax}(\mathbf{Q}\mathbf{K}^{\mathsf{T}})\mathbf{V}$ with $O(N^2d)$ FLOPs and requires O(N) additional memory beyond inputs and output.

3.2 Analysis: IO Complexity of FlashAttention

We analyze the IO complexity of FlashAttention, showing significant reduction in HBM accesses compared to standard attention. We also provide a lower bound, proving that no exact attention algorithm can

²This style of aggregation is called algebraic aggregation [35].

			<u>@</u> _	Effect	t of B	lock Size
Attention	Standard	FLASHATTENTION	96-			
GFLOPs	66.6	75.2	sesses			Runtime
HBM R/W (GB)	35.3	4.4	ÿ ₂ -	A- HBM	$\overline{}$	Runtine
Runtime (ms)	35.1	11.7	BM /	ccesses		
			岩 6	4 128	256	

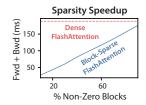


Figure 2: Left: Forward + backward runtime of standard attention and FlashAttention for seq. length 1024, head dim. 64, 16 heads, batch size 64, key-padding mask and no dropout on A100 GPU. HBM access is one of the primary factors affecting runtime. Middle: Forward runtime of FlashAttention (seq. length 1024, head dim. 64, 16 heads, batch size 64) on A100 GPU. Fewer HBM accesses result in faster runtime, up to a point. Right: The runtime (for seq. length 4K) of block-sparse FlashAttention is faster than FlashAttention by a factor proportional to the sparsity.

Block Size

asymptotically improve on HBM accesses over all SRAM sizes. Proofs are in Appendix C.

Theorem 2. Let N be the sequence length, d be the head dimension, and M be size of SRAM with $d \le M \le Nd$. Standard attention (Algorithm 0) requires $\Theta(Nd + N^2)$ HBM accesses, while FlashAttention (Algorithm 1) requires $\Theta(N^2d^2M^{-1})$ HBM accesses.

For typical values of d (64-128) and M (around 100KB), d^2 is many times smaller than M, and thus FLASHATTENTION requires many times fewer HBM accesses than standard implementation. This leads to both faster execution and lower memory footprint, which we validate in Section 4.3.

The main idea of the proof is that given the SRAM size of M, we can load blocks of K, V of size $\Theta(M)$ each (Algorithm 1 line 6). For each block of K and V, we iterate over all blocks of Q (Algorithm 1 line 8) to compute the intermediate values, resulting in $\Theta(NdM^{-1})$ passes over Q. Each pass loads $\Theta(Nd)$ elements, which amounts to $\Theta(N^2d^2M^{-1})$ HBM accesses. We similarly prove that the backward pass of standard attention requires $\Theta(Nd+N^2)$ HBM accesses while the backward pass of FLASHATTENTION requires $\Theta(N^2d^2M^{-1})$ HBM accesses (Appendix B).

We prove a lower-bound: one cannot asymptotically improve on the number of HBM accesses for all values of M (the SRAM size) when computing exact attention.

Proposition 3. Let N be the sequence length, d be the head dimension, and M be size of SRAM with $d \leq M \leq Nd$. There does not exist an algorithm to compute exact attention with $o(N^2d^2M^{-1})$ HBM accesses for all M in the range [d, Nd].

The proof relies on the fact that for $M = \Theta(Nd)$ any algorithm must perform $\Omega(N^2d^2M^{-1}) = \Omega(Nd)$ HBM accesses. This type of lower bound over a subrange of M is common in the streaming algorithms literature [92]. We leave proving parameterized complexity [29] lower bounds in terms of M as exciting future work.

We validate that the number of HBM accesses is the main determining factor of attention run-time. In Fig. 2 (left), we see that even though FlashAttention has higher FLOP count compared to standard attention (due to recomputation in the backward pass), it has much fewer HBM accesses, resulting in much faster runtime. In Fig. 2 (middle), we vary the block size B_c of FLASHATTENTION, which results in different amounts of HBM accesses, and measure the runtime of the forward pass. As block size increases, the number of HBM accesses decreases (as we make fewer passes over the input), and runtime decreases. For large enough block size (beyond 256), the runtime is then bottlenecked by other factors (e.g., arithmetic operations). Moreover, larger block size will not fit into the small SRAM size.

Extension: Block-Sparse FlashAttention

We extend FlashAttention to approximate attention: we propose block-sparse FlashAttention, whose IO complexity is smaller than FlashAttention by a factor proportional to the sparsity.

Given inputs $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ and a mask matrix $\tilde{\mathbf{M}} \in \{0, 1\}^{N \times N}$, we want to compute:

$$\mathbf{S} = \mathbf{Q}\mathbf{K}^{\top} \in \mathbb{R}^{N \times N}, \quad \mathbf{P} = \operatorname{softmax}(\mathbf{S} \odot \mathbb{1}_{\tilde{\mathbf{M}}}) \in \mathbb{R}^{N \times N}, \quad \mathbf{O} = \mathbf{P}\mathbf{V} \in \mathbb{R}^{N \times d},$$

where $(\mathbf{S} \odot \mathbb{I}_{\tilde{\mathbf{M}}})_{kl} = \mathbf{S}_{kl}$ if $\tilde{\mathbf{M}}_{kl} = 1$ and $-\infty$ if $\mathbf{M}_{kl} = 0$. We require $\tilde{\mathbf{M}}$ to have block form: for some block sizes B_r, B_c , for all k, l, $\tilde{\mathbf{M}}_{k,l} = \mathbf{M}_{lj}$ with $i = \lfloor k/B_r \rfloor$, $j = \lfloor l/B_c \rfloor$ for some $\mathbf{M} \in \{0,1\}^{N/B_r \times N/B_c}$.

Given a predefined block sparsity mask $\mathbf{M} \in \{0,1\}^{N/B_r \times N/B_c}$ we can easily adapt Algorithm 1 to only compute the nonzero blocks of the attention matrix. The algorithm is identical to Algorithm 1, except we skip zero blocks. We reproduce the algorithm description in Algorithm 5 in Appendix B.

We also analyze the IO complexity of block-sparse FlashAttention.

Proposition 4. Let N be the sequence length, d be the head dimension, and M be size of SRAM with $d \leq M \leq Nd$. Block-sparse FlashAttention (Algorithm 5) requires $\Theta(Nd + N^2d^2M^{-1}s)$ HBM accesses where s is the fraction of nonzero blocks in the block-sparsity mask.

We see that applying block-sparsity yields a direct improvement by the sparsity to the larger term in the IO complexity. For large sequence lengths N, s is often set to $N^{-1/2}$ [12] or $N^{-1} \log N$ [3, 18, 96], resulting in $\Theta(N\sqrt{N})$ or $\Theta(N \log N)$ IO complexity. For downstream experiments, we use the fixed butterfly sparsity pattern [18], which has been shown to be able to approximate arbitrary sparsity [17].

In Fig. 2 (right), we validate that as the sparsity increases, the runtime of block-sparse FlashAttention improves proportionally. On the LRA benchmark, block-sparse FlashAttention achieves 2.8× speedup, while performing on par with standard attention (Section 4).

4 Experiments

We evaluate the impact of using FlashAttention to train Transformer models. We validate two claims about training time and model accuracy, and report attention runtime and memory benchmarks.

- Training Speed. FLASHATTENTION outperforms the MLPerf 1.1 [60] speed record for BERT by 15%, and speeds up GPT-2 up to 3× over HuggingFace [91] and 1.8× over Megatron [80] over standard Transformers. FLASHATTENTION speeds up the long-range arena (LRA) benchmark 2.4×.
- Quality. FlashAttention scales Transformers to longer sequences, yielding higher quality. FlashAttention trains GPT-2 with context length 4K faster than Megatron trains GPT-2 with context length 1K, while achieving 0.7 better perplexity. Modeling longer sequences yields 6.4 points of lift on two long-document classification tasks. Finally, FlashAttention yields the first Transformer that can achieve better-than-random performance on the challenging Path-X task (sequence length 16K), and block-sparse FlashAttention yields the first sequence model that we know of that can achieve better-than-random performance on Path-256 (sequence length 64K).
- Benchmarking Attention. We measure the runtime and memory performance of FlashAttention and block-sparse FlashAttention based on sequence length. We confirm that the memory footprint of FlashAttention scales linearly with seq. length and is up to 3× faster than standard attention for common seq. lengths (up to 2K). We confirm that runtime of block-sparse FlashAttention scales linearly in seq. length and is faster than all existing approximate attention baselines.

Additional experiment details are in Appendix E.

4.1 Faster Models with FlashAttention

BERT. FLASHATTENTION yields the fastest single-node BERT training speed that we know of. We train a BERT-large [24] model with FLASHATTENTION on Wikipedia. Table 1 compares our training time to the implementation from Nvidia that set the training speed record for MLPerf 1.1 [60, 63]. Our implementation is 15% faster.

Table 1: Training time of BERT-large, starting from the same initialization provided by the MLPerf benchmark, to reach the target accuracy of 72.0% on masked language modeling. Averaged over 10 runs on 8×A100 GPUs.

BERT Implementation	Training time (minutes)
Nvidia MLPerf 1.1 [63]	20.0 ± 1.5
FLASHATTENTION (ours)	17.4 ± 1.4

GPT-2. FLASHATTENTION yields faster training times for GPT-2 [70] on the large OpenWebtext dataset [34] than the widely used HuggingFace [91] and Megatron-LM [80] implementations. Table 2 shows up to 3× end-to-end speedup compared to Huggingface and 1.7× speedup compared to Megatron-LM. FLASHATTENTION achieves the same perplexity as the other two implementations, as we do not change the model definition. Appendix E includes plots of the validation perplexity throughout training, confirming that FLASHATTENTION is as numerically stable as the baselines and produces the same training / validation curves.

Table 2: GPT-2 small and medium using FlashAttention achieve up to 3× speed up compared to Huggingface implementation and up to 1.7× compared to Megatron-LM. Training time reported on 8×A100s GPUs.

Model implementations	OpenWebText (ppl)	Training time (speedup)
GPT-2 small - Huggingface [91]	18.2	$9.5 \text{ days } (1.0 \times)$
GPT-2 small - Megatron-LM [80]	18.2	$4.7 \text{ days } (2.0 \times)$
GPT-2 small - FlashAttention	18.2	$\textbf{2.7 days} \textbf{(3.5} \times \textbf{)}$
GPT-2 medium - Huggingface [91]	14.2	21.0 days (1.0×)
GPT-2 medium - Megatron-LM [80]	14.2	11.5 days $(1.8\times)$
GPT-2 medium - FLASHATTENTION	14.2	$6.9 ext{ days } (3.0 \times)$

Long-range Arena. We compare vanilla Transformer (with either standard implementation or FLASHAT-TENTION) on the long-range arena (LRA [83]) benchmark. We measure accuracy, throughput, and training time of all models. Each task has a different sequence length varying between 1024 and 4096. We follow the implementation and experimental setting in Tay et al. [83] and Xiong et al. [94]. Table 3 shows that FLASHAT-TENTION achieves up 2.4× speed-up compared to standard attention. Block-sparse FLASHATTENTION is faster than all of the approximate attention methods that we have tested.

Table 3: The performance of standard attention, FlashAttention, block-sparse FlashAttention, and approximate attention baselines on the Long-Range-Arena benchmarks.

Models	ListOps	Text	Retrieval	Image	Pathfinder	Avg	Speedup
Transformer	36.0	63.6	81.6	42.3	72.7	59.3	-
FLASHATTENTION	37.6	63.9	81.4	43.5	72.7	59.8	$2.4 \times$
Block-sparse FlashAttention	37.0	63.0	81.3	43.6	73.3	59.6	2.8 imes
Linformer [88]	35.6	55.9	77.7	37.8	67.6	54.9	2.5×
Linear Attention [52]	38.8	63.2	80.7	42.6	72.5	59.6	$2.3 \times$
Performer [13]	36.8	63.6	82.2	42.1	69.9	58.9	$1.8 \times$
Local Attention [83]	36.1	60.2	76.7	40.6	66.6	56.0	$1.7 \times$
Reformer [53]	36.5	63.8	78.5	39.6	69.4	57.6	$1.3 \times$
Smyrf [20]	36.1	64.1	79.0	39.6	70.5	57.9	$1.7 \times$

4.2 Better Models with Longer Sequences

Language Modeling with Long Context. The runtime and memory-efficiency of FlashAttention allow us to increase the context length of GPT-2 by 4× while still running faster than the optimized implementation from Megatron-LM. Table 4 shows that that GPT-2 with FlashAttention and context length 4K is still 30% faster than GPT-2 from Megatron with context length 1K, while achieving 0.7 better perplexity.

Table 4: GPT-2 small with FlashAttention, with $4\times$ larger context length compared to Megatron-LM, is still 30% faster while achieving 0.7 better perplexity. Training time on $8\times$ A100 GPUs is reported.

Model implementations	Context length	OpenWebText (ppl)	Training time (speedup)
GPT-2 small - Megatron-LM	1k	18.2	$4.7 \text{ days } (1.0 \times)$
GPT-2 small - FlashAttention	1k	18.2	$2.7~\mathrm{days}~(1.7\times)$
GPT-2 small - FlashAttention	2k	17.7	$3.0 \text{ days } (1.6 \times)$
GPT-2 small - FlashAttention	4k	17.2	$3.6 \text{ days } (1.3\times)$

³LRA accuracy results are known to be highly dependent on the tuning procedure [94]. Our reproduced baselines perform better than as reported in the original comparison [83].

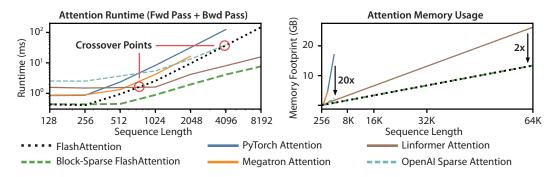


Figure 3: Left: runtime of forward pass + backward pass. Right: attention memory usage.

Long Document Classification. Training Transformers with longer sequences with FLASHATTENTION improves performance on the MIMIC-III [49] and ECtHR [6, 7] datasets. MIMIC-III contains intensive care unit patient discharge summaries, each annotated with multiple labels. ECtHR contains legal cases from the European Court of Human Rights, each of which is mapped to articles of the Convention of Human Rights that were allegedly violaged. Both of these datasets contain very long text documents; the average number of tokens in MIMIC is 2,395 tokens, and the longest document contains 14,562 tokens, while the average and longest numbers in ECtHR are 2,197 and 49,392, respectively. We evaluate lift from increasing the sequence length of a pretrained RoBERTa model [58] (we repeat the positional embeddings, as in Beltagy et al. [3]).

Table 5 shows that sequence length 16K outperforms length 512 by 4.3 points on MIMIC, and that length 8K outperforms length 512 by 8.5 points on ECtHR. The discrepancies may be due to subtle distribution shifts: MIMIC-III contains specialized medical text and thus may be more susceptible to a distribution shift in the document length, whereas ECtHR contains general language.

Table 5: Long Document performance (micro F_1) at different sequence lengths using FLASHATTENTION.

						16384
MIMIC-III [49]						
ECtHR [6]	72.2	74.3	77.1	78.6	80.7	79.2

Table 6: We report the first Transformer model that can achieve non-random performance on Path-X and Path-256.

Model	Path-X	Path-256
Transformer	Х	Х
Linformer [88]	Х	×
Linear Attention [52]	Х	×
Performer [13]	Х	×
Local Attention [83]	Х	×
Reformer [53]	Х	×
SMYRF [20]	Х	X
FLASHATTENTION	61.4	Х
Block-sparse FlashAttention	56.0	63.1

Path-X and Path-256. The Path-X and Path-256 benchmarks are challenging tasks from the long-range arena benchmark designed to test long context. The task is to classify whether two points in a black and white 128×128 (or 256×256) image have a path connecting them, and the images are fed to the transformer one pixel at a time. In prior work, all transformer models have either run out of memory, or only achieved random performance [83]. There has been a search for alternative architectures that can model such long context [39]. We present here the first result of Transformer models being able to solve Path-X and Path-256 (Table 6). We pretrain a transformer on Path-64, and then transfer to Path-X by spatially interpolating the positional embeddings. FlashAttention achieves 61.4 accuracy on Path-X. Additionally, block-sparse FlashAttention enables the Transformers to scale to sequence length 64K, achieving 63.1 accuracy 4 on Path-256.

4.3 Benchmarking Attention

We vary sequence length and measure runtime and memory usage of FlashAttention and block-sparse FlashAttention against various attention baselines on one A100 GPU with 40 GB HBM, with dropout and

⁴Path-256 requires longer sequences but has relatively shorter paths than Path-X, so it is easier to obtain a higher accuracy.

a padding mask. We compare against reference implementations for exact attention, approximate attention, and sparse attention. We report a subset of baselines in the main body; Appendix E contains more baselines and full details.

Runtime. Figure 3 (left) reports the runtime in milliseconds of the forward + backward pass of FlashAttention and block-sparse FlashAttention compared to the baselines in exact, approximate, and sparse attention (exact numbers in Appendix E). Runtime grows quadratically with sequence length, but FlashAttention runs significantly faster than **exact attention** baselines, up to 3× faster than the PyTorch implementation. The runtimes of many approximate/sparse attention mechanisms grow linearly with sequence length, but FlashAttention still runs faster than approximate and sparse attention for short sequences due to fewer memory accesses. The **approximate attention** runtimes begin to cross over with FlashAttention at sequences between 512 and 1024. On the other hand, block-sparse FlashAttention is faster than all implementations of exact, sparse, and approximate attention that we know of, across all sequence lengths.

Memory Footprint. Figure 3 (right) shows the memory footprint of FlashAttention and block-sparse FlashAttention compared to various exact, approximate, and sparse attention baselines. FlashAttention and block-sparse FlashAttention have the same memory footprint, which grows linearly with sequence length. FlashAttention is up to 20× more memory efficient than exact attention baselines, and is more memory-efficient than the approximate attention baselines. All other algorithms except for Linformer run out of memory on an A100 GPU before 64K, and FlashAttention is still 2× more efficient than Linformer.

5 Limitations and Future Directions

We discuss limitations of our approach and future directions. Related work is given in Appendix A.

Compiling to CUDA. Our current approach to building IO-aware implementations of attention requires writing a new CUDA kernel for each new attention implementation. This requires writing the attention algorithm in a considerably lower-level language than PyTorch, and requires significant engineering effort. Implementations may also not be transferrable across GPU architectures. These limitations suggest the need for a method that supports writing attention algorithms in a high-level language (e.g., PyTorch), and compiling to IO-aware implementations in CUDA—similar to efforts such as Halide in image processing [73].

IO-Aware Deep Learning. We believe that the IO-aware approach can extend beyond attention. Attention is the most memory-intensive computation in Transformers, but every layer in a deep network touches GPU HBM. We hope our work inspires IO-aware implementations of additional modules. We discuss these potential extensions in Appendix D.

Multi-GPU IO-Aware Methods. Our IO-aware implementation of attention is optimal within constants for computing attention on a single GPU. However, the attention computation may be parallelizable across multiple GPUs [75]. Using multiple GPUs adds an additional layer to IO analysis—accounting for data transfer between GPUs. We hope our work inspires future work in this direction.

Acknowledgments

Our implementation uses Apex's FMHA code (https://github.com/NVIDIA/apex/tree/master/apex/contrib/csrc/fmha) as a starting point. We thank Young-Jun Ko for the in-depth explanation of his FMHA implementation and for his thoughtful answers to our questions about CUDA. We thank Sabri Eyuboglu, Megan Leszczynski, Laurel Orr, Yuhuai Wu, Beidi Chen, and Xun Huang for their constructive feedback and suggestions on early drafts of the paper. We thank Markus Rabe and Charles Staats for helpful discussion of their attention algorithm.

We gratefully acknowledge the support of NIH under No. U54EB020405 (Mobilize), NSF under Nos. CCF1763315 (Beyond Sparsity), CCF1563078 (Volume to Velocity), and 1937301 (RTML); ARL under No. W911NF-21-2-0251 (Interactive Human-AI Teaming); ONR under No. N000141712266 (Unifying Weak Supervision); ONR N00014-20-1-2480: Understanding and Applying Non-Euclidean Geometry in Machine Learning; N000142012275 (NEPTUNE); NXP, Xilinx, LETI-CEA, Intel, IBM, Microsoft, NEC, Toshiba,

TSMC, ARM, Hitachi, BASF, Accenture, Ericsson, Qualcomm, Analog Devices, Google Cloud, Salesforce, Total, the HAI-GCP & HAI-Azure Cloud Credits for Research program, the Stanford Data Science Initiative (SDSI), Department of Defense (DoD) through the National Defense Science and Engineering Graduate Fellowship (NDSEG) Program, and members of the Stanford DAWN project: Facebook, Google, and VMWare. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright notation thereon. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views, policies, or endorsements, either expressed or implied, of NIH, ONR, or the U.S. Government. Atri Rudra's research is supported by NSF grant CCF-1763481.

References

- [1] Alok Aggarwal and S Vitter, Jeffrey. The input/output complexity of sorting and related problems. Communications of the ACM, 31(9):1116–1127, 1988.
- [2] Irwan Bello. LambdaNetworks: Modeling long-range interactions without attention. arXiv preprint arXiv:2102.08602, 2021.
- [3] Iz Beltagy, Matthew E Peters, and Arman Cohan. Longformer: The long-document transformer. arXiv preprint arXiv:2004.05150, 2020.
- [4] L Susan Blackford, Antoine Petitet, Roldan Pozo, Karin Remington, R Clint Whaley, James Demmel, Jack Dongarra, Iain Duff, Sven Hammarling, Greg Henry, et al. An updated set of basic linear algebra subprograms (blas). ACM Transactions on Mathematical Software, 28(2):135–151, 2002.
- [5] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. Advances in neural information processing systems, 33:1877–1901, 2020.
- [6] Ilias Chalkidis, Ion Androutsopoulos, and Nikolaos Aletras. Neural legal judgment prediction in English. In Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics, pages 4317–4323, Florence, Italy, 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1424. URL https://www.aclweb.org/anthology/P19-1424.
- [7] Ilias Chalkidis, Manos Fergadiotis, Dimitrios Tsarapatsanis, Nikolaos Aletras, Ion Androutsopoulos, and Prodromos Malakasiotis. Paragraph-level rationale extraction through regularization: A case study on european court of human rights cases. In Proceedings of the Annual Conference of the North American Chapter of the Association for Computational Linguistics, Mexico City, Mexico, 2021. Association for Computational Linguistics.
- [8] Benjamin Charlier, Jean Feydy, Joan Alexis Glaunès, François-David Collin, and Ghislain Durif. Kernel operations on the gpu, with autodiff, without memory overflows. *Journal of Machine Learning Research*, 22(74):1-6, 2021. URL http://jmlr.org/papers/v22/20-275.html.
- [9] Beidi Chen, Tri Dao, Eric Winsor, Zhao Song, Atri Rudra, and Christopher Ré. Scatterbrain: Unifying sparse and low-rank attention. In Advances in Neural Information Processing Systems (NeurIPS), 2021.
- [10] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear memory cost. arXiv preprint arXiv:1604.06174, 2016.
- [11] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Meghan Cowan, Leyuan Wang, Yuwei Hu, Luis Ceze, et al. {TVM}: An automated {End-to-End} optimizing compiler for deep learning. In 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI 18), pages 578-594, 2018.
- [12] Rewon Child, Scott Gray, Alec Radford, and Ilya Sutskever. Generating long sequences with sparse transformers. arXiv preprint arXiv:1904.10509, 2019.

- [13] Krzysztof Marcin Choromanski, Valerii Likhosherstov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Quincy Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. In *International Conference on Learning Representations (ICLR)*, 2020.
- [14] Xiang Dai, Ilias Chalkidis, Sune Darkner, and Desmond Elliott. Revisiting transformer-based models for long document classification. arXiv preprint arXiv:2204.06683, 2022.
- [15] Zihang Dai, Zhilin Yang, Yiming Yang, Jaime G Carbonell, Quoc Le, and Ruslan Salakhutdinov. Transformer-XL: Attentive language models beyond a fixed-length context. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 2978–2988, 2019.
- [16] Tri Dao, Albert Gu, Matthew Eichhorn, Atri Rudra, and Christopher Ré. Learning fast algorithms for linear transforms using butterfly factorizations. In *International Conference on Machine Learning (ICML)*, 2019.
- [17] Tri Dao, Nimit Sohoni, Albert Gu, Matthew Eichhorn, Amit Blonder, Megan Leszczynski, Atri Rudra, and Christopher Ré. Kaleidoscope: An efficient, learnable representation for all structured linear maps. In *International Conference on Learning Representations (ICLR)*, 2020.
- [18] Tri Dao, Beidi Chen, Kaizhao Liang, Jiaming Yang, Zhao Song, Atri Rudra, and Christopher Ré. Pixelated butterfly: Simple and efficient sparse training for neural network models. In *International Conference on Learning Representations (ICLR)*, 2022.
- [19] Tri Dao, Beidi Chen, Nimit Sohoni, Arjun Desai, Michael Poli, Jessica Grogan, Alexander Liu, Aniruddh Rao, Atri Rudra, and Christopher Ré. Monarch: Expressive structured matrices for efficient and accurate training. In *International Conference on Machine Learning (ICML)*, 2022.
- [20] Giannis Daras, Nikita Kitaev, Augustus Odena, and Alexandros G Dimakis. Smyrf-efficient attention using asymmetric clustering. Advances in Neural Information Processing Systems, 33:6476–6489, 2020.
- [21] Christopher De Sa, Albert Gu, Rohan Puttagunta, Christopher Ré, and Atri Rudra. A two-pronged progress in structured dense matrix vector multiplication. In *Proceedings of the Twenty-Ninth Annual ACM-SIAM Symposium on Discrete Algorithms*, pages 1060–1079. SIAM, 2018.
- [22] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE conference on computer vision and pattern recognition, pages 248–255. Ieee, 2009.
- [23] Peter J Denning. The working set model for program behavior. Communications of the ACM, 11(5): 323–333, 1968.
- [24] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. 2019.
- [25] Xin Dong, Shangyu Chen, and Sinno Jialin Pan. Learning to prune deep neural networks via layer-wise optimal brain surgeon. arXiv preprint arXiv:1705.07565, 2017.
- [26] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *International Conference on Learning Representations*, 2020.
- [27] Y Eidelman and I Gohberg. On a new class of structured matrices. *Integral Equations and Operator Theory*, 34(3):293–324, 1999.
- [28] Jean Feydy, Joan Glaunès, Benjamin Charlier, and Michael Bronstein. Fast geometric learning with symbolic matrices. Advances in Neural Information Processing Systems, 33, 2020.
- [29] Jörg Flum and Martin Grohe. Parameterized Complexity Theory. Springer, 2006.

- [30] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. In *International Conference on Learning Representations*, 2018.
- [31] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M Roy, and Michael Carbin. Stabilizing the lottery ticket hypothesis. arXiv preprint arXiv:1903.01611, 2019.
- [32] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel Roy, and Michael Carbin. Linear mode connectivity and the lottery ticket hypothesis. In *International Conference on Machine Learning*, pages 3259–3269 PMLR, 2020.
- [33] Karan Goel, Albert Gu, Chris Donahue, and Christopher Ré. It's raw! audio generation with state-space models. In *International Conference on Machine Learning (ICML)*, 2022.
- [34] Aaron Gokaslan, Vanya Cohen, Pavlick Ellie, and Stefanie Tellex. Openwebtext corpus, 2019.
- [35] Jim Gray, Surajit Chaudhuri, Adam Bosworth, Andrew Layman, Don Reichart, Murali Venkatrao, Frank Pellow, and Hamid Pirahesh. Data cube: A relational aggregation operator generalizing group-by, cross-tab, and sub-totals. *Data mining and knowledge discovery*, 1(1):29–53, 1997.
- [36] Andreas Griewank and Andrea Walther. Evaluating derivatives: principles and techniques of algorithmic differentiation. SIAM, 2008.
- [37] Albert Gu, Tri Dao, Stefano Ermon, Atri Rudra, and Christopher Ré. Hippo: Recurrent memory with optimal polynomial projections. In Advances in neural information processing systems (NeurIPS), 2020.
- [38] Albert Gu, Isys Johnson, Karan Goel, Khaled Saab, Tri Dao, Atri Rudra, and Christopher Ré. Combining recurrent, convolutional, and continuous-time models with linear state space layers. *Advances in Neural Information Processing Systems*, 34, 2021.
- [39] Albert Gu, Karan Goel, and Christopher Ré. Efficiently modeling long sequences with structured state spaces. In *The International Conference on Learning Representations (ICLR)*, 2022.
- [40] Song Han, Jeff Pool, John Tran, and William J Dally. Learning both weights and connections for efficient neural networks. arXiv preprint arXiv:1506.02626, 2015.
- [41] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding. In *International Conference on Learning Representations*, 2016.
- [42] John Hennessy and David Patterson. Memory hierarchy design. Computer Architecture: A Quantitative Approach, pages 390–525, 2003.
- [43] Sara Hooker. The hardware lottery. arXiv preprint arXiv:2009.06489, 2020.
- [44] Weizhe Hua, Zihang Dai, Hanxiao Liu, and Quoc V Le. Transformer quality in linear time. arXiv preprint arXiv:2202.10447, 2022.
- [45] Andrei Ivanov, Nikoli Dryden, Tal Ben-Nun, Shigang Li, and Torsten Hoefler. Data movement is all you need: A case study on optimizing transformers. *Proceedings of Machine Learning and Systems*, 3: 711–732, 2021.
- [46] Zhe Jia and Peter Van Sandt. Dissecting the Ampere GPU architecture via microbenchmarking. GPU Technology Conference, 2021.
- [47] Zhe Jia, Marco Maggioni, Benjamin Staiger, and Daniele P Scarpazza. Dissecting the nvidia Volta GPU architecture via microbenchmarking. arXiv preprint arXiv:1804.06826, 2018.
- [48] Zhe Jia, Blake Tillman, Marco Maggioni, and Daniele Paolo Scarpazza. Dissecting the graphcore IPU architecture via microbenchmarking. arXiv preprint arXiv:1912.03413, 2019.

- [49] Alistair EW Johnson, Tom J Pollard, Lu Shen, Li-wei H Lehman, Mengling Feng, Mohammad Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G Mark. Mimic-iii, a freely accessible critical care database. *Scientific data*, 3(1):1–9, 2016.
- [50] Norman P Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, et al. In-datacenter performance analysis of a tensor processing unit. In *Proceedings of the 44th annual international symposium on computer architecture*, pages 1–12, 2017.
- [51] Thomas Kailath, Sun-Yuan Kung, and Martin Morf. Displacement ranks of matrices and linear equations. Journal of Mathematical Analysis and Applications, 68(2):395–407, 1979.
- [52] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are RNNs: Fast autoregressive transformers with linear attention. In *International Conference on Machine Learning*, pages 5156–5165. PMLR, 2020.
- [53] Nikita Kitaev, Łukasz Kaiser, and Anselm Levskaya. Reformer: The efficient transformer. In *The International Conference on Machine Learning (ICML)*, 2020.
- [54] Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. Albert: A lite BEDRT for self-supervised learning of language representations. In *The International Conference on Learning Representations (ICLR)*, 2020.
- [55] Mingzhen Li, Yi Liu, Xiaoyan Liu, Qingxiao Sun, Xin You, Hailong Yang, Zhongzhi Luan, Lin Gan, Guangwen Yang, and Depei Qian. The deep learning compiler: A comprehensive survey. IEEE Transactions on Parallel and Distributed Systems, 32(3):708-727, 2020.
- [56] Valerii Likhosherstov, Krzysztof Choromanski, Jared Davis, Xingyou Song, and Adrian Weller. Sub-linear memory: How to make performers slim. arXiv preprint arXiv:2012.11346, 2020.
- [57] Ji Lin, Yongming Rao, Jiwen Lu, and Jie Zhou. Runtime neural pruning. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, Advances in Neural Information Processing Systems, volume 30. Curran Associates, Inc., 2017.
- [58] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692, 2019.
- [59] Xuezhe Ma, Xiang Kong, Sinong Wang, Chunting Zhou, Jonathan May, Hao Ma, and Luke Zettlemoyer. Luna: Linear unified nested attention. Advances in Neural Information Processing Systems, 34, 2021.
- [60] Peter Mattson, Christine Cheng, Gregory Diamos, Cody Coleman, Paulius Micikevicius, David Patterson, Hanlin Tang, Gu-Yeon Wei, Peter Bailis, Victor Bittorf, et al. Mlperf training benchmark. Proceedings of Machine Learning and Systems, 2:336–349, 2020.
- [61] Frank McSherry, Michael Isard, and Derek G Murray. Scalability! but at what {COST}? In 15th Workshop on Hot Topics in Operating Systems (HotOS XV), 2015.
- [62] Maxim Milakov and Natalia Gimelshein. Online normalizer calculation for softmax. arXiv preprint arXiv:1805.02867, 2018.
- [63] MLCommons. Mlperf 1.1 training results, 2021. URL https://mlcommons.org/en/training-normal-11/.
- [64] NVIDIA. Nvidia Tesla V100 GPU architecture, 2017.
- [65] NVIDIA. Nvidia A100 tensor core GPU architecture, 2020.
- [66] NVIDIA. Nvidia H100 tensor core GPU architecture, 2022.

- [67] D Stott Parker. Random butterfly transformations with applications in computational linear algebra. 1995.
- [68] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. Advances in neural information processing systems, 32, 2019.
- [69] Markus N Rabe and Charles Staats. Self-attention does not need $O(n^2)$ memory. $arXiv\ preprint$ $arXiv:2112.05682,\ 2021.$
- [70] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [71] Jack Rae and Ali Razavi. Do transformers need deep long-range memory? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, Online, July 2020. Association for Computational Linguistics. URL https://www.aclweb.org/anthology/2020.acl-main.672.
- [72] Jack W Rae, Anna Potapenko, Siddhant M Jayakumar, and Timothy P Lillicrap. Compressive transformers for long-range sequence modelling. In *The International Conference on Learning Representations* (ICLR), 2020.
- [73] Jonathan Ragan-Kelley, Connelly Barnes, Andrew Adams, Sylvain Paris, Frédo Durand, and Saman Amarasinghe. Halide: a language and compiler for optimizing parallelism, locality, and recomputation in image processing pipelines. Acm Sigplan Notices, 48(6):519–530, 2013.
- [74] Raghu Ramakrishnan, Johannes Gehrke, and Johannes Gehrke. Database management systems, volume 3. McGraw-Hill New York, 2003.
- [75] Benjamin Recht and Christopher Ré. Parallel stochastic gradient algorithms for large-scale matrix completion. *Mathematical Programming Computation*, 5(2):201–226, 2013.
- [76] Hongyu Ren, Hanjun Dai, Zihang Dai, Mengjiao Yang, Jure Leskovec, Dale Schuurmans, and Bo Dai. Combiner: Full attention transformer with sparse computation cost. Advances in Neural Information Processing Systems, 34, 2021.
- [77] Aurko Roy, Mohammad Saffar, Ashish Vaswani, and David Grangier. Efficient content-based sparse attention with routing transformers. *Transactions of the Association for Computational Linguistics*, 9: 53–68, 2021.
- [78] Amit Sabne. XLA: Compiling machine learning for peak performance. 2020.
- [79] Victor Sanh, Thomas Wolf, and Alexander M Rush. Movement pruning: Adaptive sparsity by fine-tuning. arXiv preprint arXiv:2005.07683, 2020.
- [80] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan Catanzaro. Megatron-LM: Training multi-billion parameter language models using model parallelism. arXiv preprint arXiv:1909.08053, 2019.
- [81] Vikas Sindhwani, Tara Sainath, and Sanjiv Kumar. Structured transforms for small-footprint deep learning. In Advances in Neural Information Processing Systems, pages 3088–3096, 2015.
- [82] Sainbayar Sukhbaatar, Edouard Grave, Piotr Bojanowski, and Armand Joulin. Adaptive attention span in transformers. In *Proceedings of the Annual Meeting of the Association for Computational Linguistics*, 2019.
- [83] Yi Tay, Mostafa Dehghani, Samira Abnar, Yikang Shen, Dara Bahri, Philip Pham, Jinfeng Rao, Liu Yang, Sebastian Ruder, and Donald Metzler. Long range arena: A benchmark for efficient transformers. In *International Conference on Learning Representations*, 2020.

- [84] Yi Tay, Mostafa Dehghani, Dara Bahri, and Donald Metzler. Efficient transformers: A survey. arXiv preprint arXiv:2009.06732, 2020.
- [85] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021.
- [86] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. Advances in neural information processing systems, 30, 2017.
- [87] Hongyu Wang, Shuming Ma, Li Dong, Shaohan Huang, Dongdong Zhang, and Furu Wei. Deepnet: Scaling transformers to 1,000 layers. arXiv preprint arXiv:2203.00555, 2022.
- [88] Sinong Wang, Belinda Z Li, Madian Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. arXiv preprint arXiv:2006.04768, 2020.
- [89] Samuel Williams, Andrew Waterman, and David Patterson. Roofline: an insightful visual performance model for multicore architectures. *Communications of the ACM*, 52(4):65–76, 2009.
- [90] Michael E Wolf and Monica S Lam. A data locality optimizing algorithm. In *Proceedings of the ACM SIGPLAN 1991 conference on Programming language design and implementation*, pages 30–44, 1991.
- [91] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander M. Rush. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online, October 2020. Association for Computational Linguistics. URL https://www.aclweb.org/anthology/2020.emnlp-demos.6.
- [92] David P Woodruff. Optimal space lower bounds for all frequency moments. In *SODA*, volume 4, pages 167–175. Citeseer, 2004.
- [93] Felix Wu, Angela Fan, Alexei Baevski, Yann N Dauphin, and Michael Auli. Pay less attention with lightweight and dynamic convolutions. In *The International Conference on Learning Representations* (ICLR), 2019.
- [94] Yunyang Xiong, Zhanpeng Zeng, Rudrasis Chakraborty, Mingxing Tan, Glenn Fung, Yin Li, and Vikas Singh. Nyströmformer: A nystöm-based algorithm for approximating self-attention. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, page 14138, 2021.
- [95] Li Yuan, Yunpeng Chen, Tao Wang, Weihao Yu, Yujun Shi, Zi-Hang Jiang, Francis EH Tay, Jiashi Feng, and Shuicheng Yan. Tokens-to-token vit: Training vision transformers from scratch on imagenet. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 558–567, 2021.
- [96] Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, et al. Big bird: Transformers for longer sequences. Advances in Neural Information Processing Systems, 33, 2020.
- [97] Shuangfei Zhai, Walter Talbott, Nitish Srivastava, Chen Huang, Hanlin Goh, Ruixiang Zhang, and Josh Susskind. An attention free transformer. arXiv preprint arXiv:2105.14103, 2021.
- [98] Chen Zhu, Wei Ping, Chaowei Xiao, Mohammad Shoeybi, Tom Goldstein, Anima Anandkumar, and Bryan Catanzaro. Long-short transformer: Efficient transformers for language and vision. *Advances in Neural Information Processing Systems*, 34, 2021.

A Related Work

IO-Aware Runtime Optimization. The broad concept of optimizing for reading and writing to fast/slow memory has a long history in computer science and has been known by many names. We draw the most direct connection to the literature of analyzing I/O complexity in this work [1], but concepts of memory hierarchies are fundamental and has appeared in many forms, from the working set model [23], to data locality [90], to the Roofline model of arithmetic intensity [89], to analyses of scalability [61], to standard textbook treatments of computer architecture [42]. We hope that this work encourages the community to adopt these ideas in more parts of the deep learning stack.

Efficient ML Models with Structured Matrices. Matrix multiply is the core computational bottleneck of most machine learning models. To reduce the computational complexity, there have been numerous approaches to learn over a more efficient set of matrices. These matrices are called *structured matrices*, which have subquadratic $(o(n^2))$ for dimension $n \times n$ number of parameters and runtime. Most common examples of structured matrices are sparse and low-rank matrices, along with fast transforms commonly encountered in signal processing (Fourier, Chebyshev, sine/cosine, orthogonal polynomials). There have been several more general classes of structured matrices proposed in machine learning: Toeplitz-like [81], low-displacement rank [51], quasi-separable [27]). The butterfly pattern we use for our block-sparse attention is motivated by the fact that butterfly matrices [16, 67] and their products have been shown to be able to express any structured matrices with almost optimal runtime and number of parameters [17, 21]. However, even though structured matrices are efficient in theory, they have not seen wide adoption since it is hard to translate their efficiency to wall-clock speedup since dense unconstrained matrix multiply has very optimize implementation, a phenomenon known as the hardware lottery [43]. Extensions of butterfly matrices [18, 19] aimed to make butterfly matrices more hardware-friendly.

Sparse Training. Our block-sparse FLASHATTENTION can be seen as a step towards making sparse model training more efficient. Sparse models have seen success in compressing models for inference (pruning) by sparsifying the weight matrices [25, 40, 41, 57, 79]. For model training, the lottery tickets hypothesis [30, 31, 32] suggests that there are a set of small sub-networks derived from a larger dense network that performs as well as the original dense network. Out block-sparse FLASHATTENTION can also be seen as a fixed lottery ticket in the context of attention: we fix the sparsity pattern to be the butterfly pattern through training, and observe that it performs almost as well as the (dense) FLASHATTENTION on the Long-range Arena tasks.

Efficient Transformer. Transformer-based models have become the most widely-used architecture in natural language processing [24] and computer vision [26, 95]. However, one of their computational bottlenecks is that their time and memory scales quadratic in the sequence length. There are numerous approaches to overcome this bottleneck, including approximation with hashing (i.e., sparse) such as Reformer [53] and Smyrf [20] and with low-rank approximation such as Performer [13, 56]. One can even combine sparse and low-rank approximation for better accuracy (e.g., Longformer [3], BigBird [96], Scatterbrain [9], Long-short transformer [98], Combiner [76]). Other approaches include compressing along the sequence dimension to attend to multiple tokens at once [54, 59, 82, 93]. One can also attend over the states from previous sequences to help lengthen the context (e.g., Transformer-XL [15] and Compressive Transformer [72]). We recommend the survey [84] for more details.

There are several lines of work on developing other modules instead of attention to model longer context. HiPPO [37] and its extensions, most notably S4 [33, 38, 39] projects the history on a polynomial basis, allowing accurate reconstruction of the history through state-space models. They combine the strengths of CNNs (efficient training), RNNs (efficient inference), and continuous models (robust to change in sampling rates). LambdaNetworks [2], AFT [97] and FLASH [44] are other attempts at replacing attention in the context of image classification and language modeling.

B Algorithm Details

We first derive the forward and backward passes of attention and show that they can be computed in a memory-efficient manner (requiring extra memory linear instead of quadratic in the sequence length). Though they reduce the amount of extra memory required, naively they still incur quadratic HBM accesses, resulting in slower execution speed. We describe the FLASHATTENTION algorithm to implement both the forward

and the backward passes on GPUs that reduces HBM accesses, leading to both faster runtime and smaller memory footprint.

B.1 Memory-efficient forward pass

The main challenge in making attention memory-efficient is the softmax that couples the columns of \mathbf{K} (and columns of \mathbf{V}). Our approach is to compute the softmax normalization constant separately to decouple the columns. This technique [62] has been used in the literature [53, 69] to show that attention computation does not need quadratic *extra* memory (though the number of HBM accesses is still quadratic, resulting in slow run-time).

For simplicity, we omit here the max-shifting step during softmax. The full algorithm in Appendix B.3 contains all the steps.

Recall that given input sequences $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$, we want to compute the attention output $\mathbf{O} \in \mathbb{R}^{N \times d}$:

$$\mathbf{S} = \mathbf{Q}\mathbf{K}^{\top} \in \mathbb{R}^{N \times N}, \quad \mathbf{P} = \operatorname{softmax}(\mathbf{S}) \in \mathbb{R}^{N \times N}, \quad \mathbf{O} = \mathbf{P}\mathbf{V} \in \mathbb{R}^{N \times d}$$

We have that $S_{ij} = q_i^T k_j$ where q_i and k_j are the *i*-th and *j*-th columns of **Q** and **K** respectively. Define the normalization constants of softmax:

$$L_i = \sum_j e^{q_i^T k_j}. (1)$$

Let v_i be the j-th column of \mathbf{V} , then the i-th columns of the output is

$$o_i = P_{i:} \mathbf{V} = \sum_{i} P_{ij} v_j = \sum_{i} \frac{e^{q_i^T k_j}}{L_i} v_j.$$
 (2)

We see that once L_i is computed, we can compute o_i without extra memory by repeatedly summing $\frac{e^{q_i^T k_j}}{L} v_j$. Therefore the forward pass can be computed with O(n) extra memory:

- 1. Compute L_i for all i according to Eq. (1), which takes O(n) extra memory.
- 2. Compute o_i for all i according to Eq. (2), which takes O(d) extra memory.

B.2 Memory-efficient backward pass

We derive the backward pass of attention and show that it can also be computed with linear memory. Rabe and Staats [69] suggests that the backward pass can be done without quadratic extra memory by applying gradient checkpointing to the memory-efficient forward pass. We instead derive the backward pass explicitly and show how it can be computed in a memory-efficient manner.

Suppose that there is a scalar loss function ϕ , and let the output gradient be $\mathbf{dO} \in \mathbb{R}^{n \times d}$ (where \mathbf{dO} denotes $\frac{\partial \phi}{\partial \mathbf{O}}$). We want to compute the input gradients $\mathbf{dQ}, \mathbf{dK}, \mathbf{dV} \in \mathbb{R}^{n \times d}$ (where $\mathbf{dQ}, \mathbf{dK}, \mathbf{dV}$ denote $\frac{\partial \phi}{\partial \mathbf{Q}}, \frac{\partial \phi}{\partial \mathbf{K}}, \frac{\partial \phi}{\partial \mathbf{V}}$ respectively).

The gradient dV is easy to see. Applying reverse-mode autodiff by hand (aka the chain rule), we obtain (in matrix notation) $dV = P^T dO$. Thus:

$$dv_j = \sum_i P_{ij} do_i = \sum_i \frac{e^{q_i^T k_j}}{L_i} do_i.$$
 (3)

Since we already computed L_i , dv_i can be computed without extra memory by repeated summing.

The gradients $d\mathbf{Q}$ and $d\mathbf{K}$ are a little more complicated. We go through the gradients $d\mathbf{P}$ and $d\mathbf{S}$ first. From Eq. (2), we have that $d\mathbf{P} = d\mathbf{O}\mathbf{V}^T$, and so:

$$dP_{ij} = do_i^T v_j.$$

Recall that $P_{i:} = \text{softmax}(S_{i:})$. Using the fact that the Jacobian of y = softmax(x) is $\text{diag}(y) - yy^T$, we have that

$$dS_{i:} = (\operatorname{diag}(P_{i:}) - P_{i:}P_{i:}^T)dP_{i:} = P_{i:} \circ dP_{i:} - (P_{i:}^TdP_{i:})P_{i:},$$

where • denotes pointwise multiplication.

Define

$$D_{i} = P_{i:}^{T} dP_{i:} = \sum_{j} \frac{e^{q_{i}^{T} k_{j}}}{L_{i}} do_{i}^{T} v_{j} = do_{i}^{T} \sum_{j} \frac{e^{q_{i}^{T} k_{j}}}{L_{i}} v_{j} = do_{i}^{T} o_{i},$$

$$(4)$$

then

$$dS_{i:} = P_{i:} \circ dP_{i:} - D_i P_{i:}.$$

Hence

$$dS_{ij} = P_{ij}dP_{ij} - D_iP_{ij} = P_{ij}(dP_{ij} - D_i).$$

Now we can get the gradients **dQ** and **dK**. Recall that $S_{ij} = q_i^T k_j$, so

$$dq_{i} = \sum_{j} dS_{ij}k_{j} = \sum_{j} P_{ij}(dP_{ij} - D_{i})k_{j} = \sum_{j} \frac{e^{q_{i}^{T}k_{j}}}{L_{i}}(do_{i}^{T}v_{j} - D_{i})k_{j}.$$
 (5)

Similarly,

$$dk_{j} = \sum_{i} dS_{ij} q_{i} = \sum_{i} P_{ij} (dP_{ij} - D_{i}) q_{i} = \sum_{i} \frac{e^{q_{i}^{T} k_{j}}}{L_{i}} (do_{i}^{T} v_{j} - D_{i}) q_{i}.$$
 (6)

Therefore the backward pass can also be computed with O(n) extra memory:

- 1. Compute dv_j for all j according to Eq. (3), which takes O(d) extra memory.
- 2. Compute D_i for all i according to Eq. (4), which takes O(n) extra memory.
- 3. Compute dq_i for all i according to Eq. (5), which takes O(d) extra memory.
- 4. Compute dk_i for all j according to Eq. (6), which takes O(d) extra memory.

B.3 FLASHATTENTION: Forward Pass

We describe the full details of FLASHATTENTION forward pass. Given input sequences $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$, we want to compute the attention output $\mathbf{O} \in \mathbb{R}^{N \times d}$:

$$\begin{split} \mathbf{S} &= \tau \mathbf{Q} \mathbf{K}^{\top} \in \mathbb{R}^{N \times N}, \quad \mathbf{S}^{\text{masked}} &= \text{MASK}(S) \in \mathbb{R}^{N \times N}, \quad \mathbf{P} = \text{softmax}(\mathbf{S}^{\text{masked}}) \in \mathbb{R}^{N \times N}, \\ \mathbf{P}^{\text{dropped}} &= \text{dropout}(\mathbf{P}, p_{\text{drop}}), \quad \mathbf{O} &= \mathbf{P}^{\text{dropped}} \mathbf{V} \in \mathbb{R}^{N \times d}, \end{split}$$

where $\tau \in \mathbb{R}$ is some softmax scaling (typically $\frac{1}{\sqrt{d}}$), MASK is some masking function that sets some entries of the input to $-\infty$ and keep other entries the same (e.g., key padding mask when sequences in the batch don't have the same lengths and are padded), and dropout(x, p) applies dropout to x elementwise (i.e., output $\frac{x}{1-p}$ with probability 1-p and output 0 with probability p for each element x).

The full algorithm is in Algorithm 2. We save the output $\mathbf{0}$, the softmax statistics ℓ and m, and the pseudo-random number generator state \mathcal{R} for the backward pass.

Algorithm 2 FlashAttention Forward Pass

```
Require: Matrices \mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d} in HBM, on-chip SRAM of size M, softmax scaling constant \tau \in \mathbb{R},
       masking function MASK, dropout probability p_{\text{drop}}.
  1: Initialize the pseudo-random number generator state \mathcal R and save to HBM.
  2: Set block sizes B_c = \left\lceil \frac{M}{4d} \right\rceil, B_r = \min\left(\left\lceil \frac{M}{4d} \right\rceil, d\right).

3: Initialize \mathbf{O} = (0)_{N \times d} \in \mathbb{R}^{N \times d}, \ell = (0)_N \in \mathbb{R}^N, m = (-\infty)_N \in \mathbb{R}^N in HBM.
  4: Divide Q into T_r = \begin{bmatrix} \frac{N}{B_r} \end{bmatrix} blocks \mathbf{Q}_1, \dots, \mathbf{Q}_{T_r} of size B_r \times d each, and divide \mathbf{K}, \mathbf{V} in to T_c = \begin{bmatrix} \frac{N}{B_c} \end{bmatrix} blocks
       \mathbf{K}_1, \ldots, \mathbf{K}_{T_c} and \mathbf{V}_1, \ldots, \mathbf{V}_{T_c}, of size B_c \times d each.
  5: Divide \mathbf{0} into T_r blocks \mathbf{0}_i, \dots, \mathbf{0}_{T_r} of size B_r \times d each, divide \ell into T_r blocks \ell_i, \dots, \ell_{T_r} of size B_r each,
       divide m into T_r blocks m_1, \ldots, m_{T_r} of size B_r each.
  6: for 1 \le j \le T_c do
            Load \mathbf{K}_i, \mathbf{V}_i from HBM to on-chip SRAM.
            for 1 \le i \le T_r do
                 Load \mathbf{Q}_i, \mathbf{O}_i, \ell_i, m_i from HBM to on-chip SRAM.
 9:
                 On chip, compute \mathbf{S}_{ij} = \tau \mathbf{Q}_i \mathbf{K}_j^T \in \mathbb{R}^{B_r \times \overline{B}_c}.
 10:
                 On chip, compute \mathbf{S}_{ij}^{\text{masked}} = \text{MASK}(\mathbf{S}_{ij}).
11:
                 On chip, compute \tilde{m}_{ij} = \text{rowmax}(\mathbf{S}_{ij}^{\text{masked}}) \in \mathbb{R}^{B_r}, \tilde{\mathbf{P}}_{ij} = \exp(\mathbf{S}_{ij}^{\text{masked}} - \tilde{m}_{ij}) \in \mathbb{R}^{B_r \times B_c} (pointwise),
 12:
                 \tilde{\ell}_{ij} = \text{rowsum}(\tilde{\mathbf{P}}_{ij}) \in \mathbb{R}^{B_r}.
                 On chip, compute \tilde{\mathbf{m}}_{i}^{\text{new}} = \max(m_{i}, \tilde{m}_{ij}) \in \mathbb{R}^{B_{r}}, \, \ell_{i}^{\text{new}} = e^{m_{i} - m_{i}^{\text{new}}} \ell_{i} + e^{\tilde{m}_{ij} - m_{i}^{\text{new}}} \tilde{\ell}_{ij} \in \mathbb{R}^{B_{r}}.
On chip, compute \tilde{\mathbf{P}}_{ij}^{\text{dropped}} = \text{dropout}(\tilde{\mathbf{P}}_{ij}, p_{\text{drop}}).
13:
```

Write $\mathbf{O}_{i} \leftarrow \operatorname{diag}(\ell_{i}^{\operatorname{new}})^{-1}(\operatorname{diag}(\ell_{i})e^{m_{i}-m_{i}^{\operatorname{new}}}\mathbf{O}_{i} + e^{\tilde{m}_{ij}-m_{i}^{\operatorname{new}}}\tilde{\mathbf{P}}_{ij}^{\operatorname{dropped}}\mathbf{V}_{j})$ to HBM. Write $\ell_{i} \leftarrow \ell_{i}^{\operatorname{new}}$, $m_{i} \leftarrow m_{i}^{\operatorname{new}}$ to HBM. 16: end for 17:

18: end for

14:

15:

19: Return $\mathbf{0}, \ell, m, \mathcal{R}$.

B.4 FLASHATTENTION: Backward Pass

We describe the full details of FlashAttention backward pass. Given input sequences $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$, the output $\mathbf{O} \in \mathbb{R}^{N \times d}$, and the output gradient $d\mathbf{O}$, we want to compute the input gradients $d\mathbf{Q}, d\mathbf{K}, d\mathbf{V} \in \mathbb{R}^{N \times d}$. We first describe the standard attention backward pass in Algorithm 3 for completeness.

Algorithm 3 Standard Attention Backward Pass

```
Require: Matrices \mathbf{Q}, \mathbf{K}, \mathbf{V}, \mathbf{dO} \in \mathbb{R}^{N \times d}, \mathbf{P} \in \mathbb{R}^{N \times N} in HBM.
```

- 1: Load \mathbf{P}, \mathbf{dO} by blocks from HBM, compute $\mathbf{dV} = \mathbf{P}^{\mathsf{T}} \mathbf{dO} \in \mathbb{R}^{N \times d}$, write \mathbf{dV} to HBM.
- 2: Load $d\mathbf{O}, \mathbf{V}$ by blocks from HBM, compute $d\mathbf{P} = d\mathbf{O}\mathbf{V}^{\top} \in \mathbb{R}^{N \times N}$, write $d\mathbf{P}$ to HBM.
- 3: Read **P**, **dP** from HBM, compute $\mathbf{dS} \in \mathbb{R}^{N \times N}$ where $dS_{ij} = P_{ij}(dP_{ij} \sum_{l} P_{il}dP_{il})$, write \mathbf{dS} to HBM.
- 4: Load dS and K by blocks from HBM, compute dQ = dSK, write dQ to HBM.
- 5: Load dS and Q by blocks from HBM, compute $dK = dS^TQ$, write dK to HBM.
- 6: Return dQ, dK, dV.

We now make two observations about FlashAttention backward pass:

- 1. We do not need to store the dropout mask of size $O(N^2)$ from the forward pass. Instead, we can save the pseudo-random number generator states from the forward pass and re-generate the dropout mask in the backward pass. This allows us to only use O(N) extra memory.
- 2. When computing the softmax gradient, we use Eq. (4) to compute $D_i = P_{i}^{\mathsf{T}} dP_i$: without reducing over $P_{i:}$ and $dP_{i:}$ of size N (they might not fit into SRAM). Instead we can rewrite $D_i = do_i^{\mathsf{T}} o_i$ and compute the dot product between vectors of size d.

The full FlashAttention backward pass algorithm is in Algorithm 4. Conceptually it is just a block version of the derivation in Appendix B.2.

Algorithm 4 FlashAttention Backward Pass

Require: Matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V}, \mathbf{O}, \mathbf{dO} \in \mathbb{R}^{N \times d}$ in HBM, vectors $\ell, m \in \mathbb{R}^N$ in HBM, on-chip SRAM of size M, softmax scaling constant $\tau \in \mathbb{R}$, masking function MASK, dropout probability p_{drop} , pseudo-random number generator state \mathcal{R} from the forward pass.

```
1: Set the pseudo-random number generator state to \mathcal{R}.
```

```
2: Set block sizes B_c = \left\lceil \frac{M}{4d} \right\rceil, B_r = \min\left(\left\lceil \frac{M}{4d} \right\rceil, d\right).
```

- 3: Divide **Q** into $T_r = \left\lceil \frac{N}{B_r} \right\rceil$ blocks $\mathbf{Q}_1, \dots, \mathbf{Q}_{T_r}$ of size $B_r \times d$ each, and divide \mathbf{K}, \mathbf{V} in to $T_c = \left\lceil \frac{N}{B_c} \right\rceil$ blocks $\mathbf{K}_1, \ldots, \mathbf{K}_{T_c}$ and $\mathbf{V}_1, \ldots, \mathbf{V}_{T_c}$, of size $B_c \times d$ each.
- 4: Divide **O** into T_r blocks $\mathbf{O}_i, \dots, \mathbf{O}_{T_r}$ of size $B_r \times d$ each, divide \mathbf{dO} into T_r blocks $\mathbf{dO}_i, \dots, \mathbf{dO}_{T_r}$ of size $B_r \times d$ each, divide ℓ into T_r blocks $\ell_i, \ldots, \ell_{T_r}$ of size B_r each, divide m into T_r blocks m_1, \ldots, m_{T_r} of size
- 5: Initialize $\mathbf{dQ} = (0)_{N \times d}$ in HBM and divide it into T_r blocks $\mathbf{dQ}_1, \dots, \mathbf{dQ}_{T_r}$ of size $B_r \times d$ each. Initialize $\mathbf{dK} = (0)_{N \times d}, \mathbf{dV} = (0)_{N \times d}$ in HBM and divide \mathbf{dK}, \mathbf{dV} in to T_c blocks $\mathbf{dK}_1, \dots, \mathbf{dK}_{T_c}$ and $\mathbf{dV}_1, \dots, \mathbf{dV}_{T_c}$, of size $B_c \times d$ each.

```
6: for 1 \le j \le T_c do
```

- Load \mathbf{K}_j , \mathbf{V}_j from HBM to on-chip SRAM.
- Initialize $\mathbf{dK}_j = (0)_{B_c \times d}, \mathbf{dV}_j = (0)_{B_c \times d}$ on SRAM. 8:
- 9: for $1 \le i \le T_r$ do
- 10: Load $\mathbf{Q}_i, \mathbf{O}_i, \mathbf{dO}_i, \mathbf{dQ}_i, \ell_i, m_i$ from HBM to on-chip SRAM.
- On chip, compute $\mathbf{S}_{ij} = \tau \mathbf{Q}_i \mathbf{K}_i^T \in \mathbb{R}^{B_r \times B_c}$. 11:
- On chip, compute $\mathbf{S}_{ij}^{\text{masked}} = \text{MASK}(\mathbf{S}_{ij})$. 12:
- On chip, compute $\mathbf{P}_{ij} = \operatorname{diag}(l_i)^{-1} \exp(\mathbf{S}_{ij}^{\text{masked}} m_i) \in \mathbb{R}^{B_r \times B_c}$. 13:
- On chip, compute dropout mask $\mathbf{Z}_{ij} \in \mathbb{R}^{B_r \times B_c}$ where each entry has value $\frac{1}{1-p_{\text{drop}}}$ with probability 14: $1-p_{\rm drop}$ and value 0 with probability $p_{\rm drop}.$
- On chip, compute $\mathbf{P}_{ij}^{\text{dropped}} = \mathbf{P}_{ij} \circ \mathbf{Z}_{ij}$ (pointwise multiply). 15:
- On chip, compute $\mathbf{d}\tilde{\mathbf{V}}_{j} \leftarrow \mathbf{d}\tilde{\mathbf{V}}_{j} + (\mathbf{P}_{ij}^{\text{dropped}})^{\top} \mathbf{d}\mathbf{O}_{i} \in \mathbb{R}^{B_{c} \times d}$. 16:
- On chip, compute $\mathbf{dP}_{ij}^{\text{dropped}} = \mathbf{dO}_i \mathbf{V}_j^{\top} \in \mathbb{R}^{B_r \times B_c}$. 17:
- On chip, compute $\mathbf{dP}_{ij} = \mathbf{dP}_{ij}^{\text{dropped}} \circ \mathbf{Z}_{ij}$ (pointwise multiply). 18:
- On chip, compute $D_i = \text{rowsum}(\mathbf{dO}_i \circ \mathbf{O}_i) \in \mathbb{R}^{B_r}$. 19:
- On chip, compute $\mathbf{dS}_{ij} = \mathbf{P}_{ij} \circ (\mathbf{dP}_{ij} D_i) \in \mathbb{R}^{B_r \times B_c}$. Write $\mathbf{dQ}_i \leftarrow \mathbf{dQ}_i + \tau \mathbf{dS}_{ij} \mathbf{K}_j \in \mathbb{R}^{B_r \times d}$ to HBM. 20:
- 21:
- On chip, compute $\mathbf{d}\mathbf{K}_i \leftarrow \mathbf{d}\mathbf{K}_i + \tau \mathbf{d}\mathbf{S}_{ii}^{\top}\mathbf{Q}_i \in \mathbb{R}^{B_c \times d}$. 22:
- 23:
- Write $d\mathbf{K}_i \leftarrow d\tilde{\mathbf{K}}_i, d\mathbf{V}_i \leftarrow d\tilde{\mathbf{V}}_i$ to HBM. 24:
- end for 25:
- 26: Return dQ, dK, dV.

We see that similar to the forward pass, the backward pass performs $O(N^2)$ FLOPs and only requires O(N) extra memory beyond inputs, output, output gradient, and input gradients.

We analyze the IO-complexity of the backward pass, similar to the forward pass (Theorem 2).

Theorem 5. Let N be the sequence length, d be the head dimension, and M be size of SRAM with $d \le M \le Nd$. Standard attention (Algorithm 0) backward pass requires $\Theta(Nd+N^2)$ HBM accesses, while FlashAttention backward pass (Algorithm 4) requires $\Theta(N^2d^2M^{-1})$ HBM accesses.

The proof is in Appendix C.

B.5 Comparison with Rabe and Staats [69]

We describe here some similarities and differences between our FlashAttention algorithm and the algorithm of Rabe and Staats [69].

Conceptually, both FLASHATTENTION and Rabe and Staats [69] operate on blocks of the attention matrix using the well-established technique of tiling (or softmax scaling) [53, 62]. To reduce the memory footprint, both methods avoid storing the large attention matrix in the forward pass and recompute it in the backward pass.

The first major difference is that Rabe and Staats [69] focuses on the reducing the total memory footprint (maximum amount of GPU memory required) while FLASHATTENTION focuses on reducing memory accesses (the number of memory reads/writes). As mentioned in Section 2, the amount of memory access is the primary determining factor of runtime. Reducing memory accesses also necessarily reduces the total amount of memory required (e.g., if an operation incurs A memory accesses, then its total memory requirement is at most A). As a result, FLASHATTENTION is faster than standard attention (2-4×) while Rabe and Staats [69] is around the same speed or slightly slower than standard attention. In terms of total memory required, both methods offer substantial memory saving.

The second difference between the two methods is the way information is summarized from each block to pass to the next block. Rabe and Staats [69] summarizes each block with its temporary output along with the softmax normalization statistics. At the end of the forward pass, the temporary outputs of all the blocks are combined using the statistics to produce the final output. FLASHATTENTION instead incrementally updates the output (Algorithm 1 line 12) after processing each block, so only one copy of the output is needed (instead of K copies for K blocks). This means that FLASHATTENTION has smaller total memory requirement compared to Rabe and Staats [69].

The final major difference is the way the backward pass is computed. Rabe and Staats [69] uses gradient checkpointing to recompute the attention matrix and the temporary output of each block. FlashAttention instead simplifies the backward pass analytically (Appendices B.2 and B.4). It only recomputes the attention matrix and does not recompute the temporary output of each block. This reduces the memory requirement for the backward pass and yields speedup.

C Proofs

Proof of Theorem 1. We first count the number of FLOPs and extra memory required.

The dominating FLOPs are from matrix multiplication. In the inner loop, (Algorithm 1 line 9), we compute $\mathbf{Q}_i \mathbf{K}_j^{\top} \in \mathbb{R}^{B_r \times B_c}$ for $\mathbf{Q}_i \in \mathbb{R}^{B_r \times d}$ and $\mathbf{K}_j \in \mathbb{R}^{B_c \times d}$, which takes $O(B_r B_c d)$ FLOPs. We also compute (Algorithm 1 line 12) $\tilde{\mathbf{P}}_{ij} \mathbf{V}_j \in \mathbb{R}^{B_r \times d}$ for $\tilde{\mathbf{P}}_{ij} \in \mathbb{R}^{B_r \times B_c}$ and $\mathbf{V}_j \in \mathbb{R}^{B_c \times d}$, which takes $O(B_r B_c d)$ FLOPs. We execute the inner loops $T_c T_r = \left\lceil \frac{N}{B_c} \right\rceil \left\lceil \frac{N}{B_r} \right\rceil$ times. Therefore the total number of FLOPs is

$$O\left(\frac{N^2}{B_c B_r} B_r B_c d\right) = O(N^2 d).$$

In terms of extra memory required, we see that we need O(N) memory to store the statistics (ℓ, m) .

We now prove the algorithm's correctness by induction on j for $0 \le j \le T_c$. Let $\mathbf{K}_{:j} \in \mathbb{R}^{jB_c \times d}$ be the first jB_c rows of \mathbf{K} , and similarly $\mathbf{V}_{:j} \in \mathbb{R}^{jB_c \times d}$ the the first jB_c rows of \mathbf{V} . Let $\mathbf{S}_{:,:j} = \mathbf{Q}\mathbf{K}_{:j}^{\top} \in \mathbb{R}^{N \times jB_c}$, and $\mathbf{P}_{:,:j} = \text{softmax}(\mathbf{S}_{:,:j}) \in \mathbb{R}^{N \times jB_c}$ (softmax applied row-wise). Let $m^j, \ell^{(j)}, \mathbf{O}^{(j)}$ be the values of m, ℓ, \mathbf{O} in HBM after the j-th iteration of the outer loop (Algorithm 1 line 5). (Note that these values of m, ℓ, \mathbf{O} are updated after each iteration of the outer loop.) We want to show that after the j-th iteration of the outer loop, we have computed in HBM:

$$m^{(j)} = \operatorname{rowmax}(\mathbf{S}_{:,:j}) \in \mathbb{R}^N, \quad \ell^{(j)} = \operatorname{rowsum}(\exp(\mathbf{S}_{:,:j} - m^{(j)})) \in \mathbb{R}^N, \quad \mathbf{O}^{(j)} = \mathbf{P}_{:,:j} \mathbf{V}_{:j} \in \mathbb{R}^{N \times d}.$$

Based on our initialization (Algorithm 1 line 2), this claim is true for j = 0 (i.e., before the any iteration of the outer loop is executed). Suppose that the claim holds for some $j = 0, ..., T_c - 1$. We want to show that the claim also holds for j + 1. Indeed, when we update the statistics in the inner loop (Algorithm 1 line 10)

on the (j+1)-th iteration of the outer loop, we update $m^{(j+1)} = \max(m^{(j)}, \tilde{m})$ where $\tilde{m} \in \mathbb{R}^N$ is the row-max of $\mathbf{S}_{:,j:j+1}$, the slice of \mathbf{S} from column jB_c to column $(j+1)B_c - 1$. This implies that

$$m^{(j+1)} = \operatorname{rowmax}(\mathbf{S}_{:,:j+1}) \in \mathbb{R}^N.$$

Similarly, we update

$$\ell^{(j+1)} = e^{m^{(j)} - m^{(j+1)}} \ell^{(j)} + e^{\tilde{m} - m^{(j+1)}} \tilde{\ell}.$$

where $\tilde{\ell} = \text{rowsum}(\exp(\mathbf{S}_{:,j:j+1} - \tilde{m})) \in \mathbb{R}^N$. By the same algebraic manipulation in Section 3.1, we obtain:

$$\ell^{(j+1)} = \operatorname{rowsum}(\exp(\mathbf{S}_{\dots j+1} - m^{(j+1)})) \in \mathbb{R}^{N}.$$

Let $V_{j:j+1}$ be the slice of V from column jB_c to column $(j+1)B_c-1$, we also update:

$$\begin{split} \mathbf{O}^{(j+1)} &= \mathrm{diag}(\ell^{(j+1)})^{-1}(\mathrm{diag}(\ell^{(j)})e^{m^{(j)}-m^{(j+1)}}\mathbf{O}^{(j)} + e^{\tilde{m}-m^{(j+1)}}\exp(\mathbf{S}_{j:j+1}-\tilde{m})\mathbf{V}_{j:j+1}) \\ &= \mathrm{diag}(\ell^{(j+1)})^{-1}(\mathrm{diag}(\ell^{(j)})e^{m^{(j)}-m^{(j+1)}}\mathbf{P}_{:,:j}\mathbf{V}_{:j} + e^{-m^{(j+1)}}\exp(\mathbf{S}_{j:j+1})\mathbf{V}_{j:j+1}) \\ &= \mathrm{diag}(\ell^{(j+1)})^{-1}(\mathrm{diag}(\ell^{(j)})e^{m^{(j)}-m^{(j+1)}}\mathrm{diag}(\ell^{(j)})^{-1}\exp(\mathbf{S}_{:,:j}-m^{(j)})\mathbf{V}_{:j} + e^{-m^{(j+1)}}\exp(\mathbf{S}_{j:j+1})\mathbf{V}_{j:j+1}) \\ &= \mathrm{diag}(\ell^{(j+1)})^{-1}(e^{-m^{(j+1)}}\exp(\mathbf{S}_{:,:j})\mathbf{V}_{:j} + e^{-m^{(j+1)}}\exp(\mathbf{S}_{j:j+1})\mathbf{V}_{j:j+1}) \\ &= \mathrm{diag}(\ell^{(j+1)})^{-1}(\exp(\mathbf{S}_{:,:j}-m^{(j+1)})\mathbf{V}_{:j} + \exp(\mathbf{S}_{j:j+1}-m^{(j+1)})\mathbf{V}_{j:j+1}) \\ &= \mathrm{diag}(\ell^{(j+1)})^{-1}\left(\exp\left(\left[\mathbf{S}_{:,:j}-\mathbf{S}_{j:j+1}\right]-m^{(j+1)}\right)\right)\left[\begin{array}{c} \mathbf{V}_{:j} \\ \mathbf{V}_{j:j+1} \end{array}\right] \\ &= \operatorname{softmax}(\mathbf{S}_{:j+1})\mathbf{V}_{:j+1}. \end{split}$$

We then see that the claim is also true for j+1. By induction, the claim is true for all $j=0,\ldots,T_c$. When $j=T_c$, we conclude that the final value of $\mathbf{0}$ in HBM is softmax(\mathbf{S}) \mathbf{V} = softmax($\mathbf{Q}\mathbf{K}^{\mathsf{T}}$) \mathbf{V} .

Proof of Theorem 2. We first analyze the IO complexity of standard attention implementation. The inputs $\mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d}$ reside in HBM, and the at the end of the algorithm the output $\mathbf{O} \in \mathbb{R}^{N \times d}$ is written to HBM.

In the first step of computing the matrix multiply $\mathbf{S} = \mathbf{Q}\mathbf{K}^{\top}$, the inputs \mathbf{Q}, \mathbf{K} are read from HBM and the output $\mathbf{S} \in \mathbb{R}^{N \times N}$ is written to HBM (Algorithm 0 line 1). This incurs $\Theta(Nd + N^2)$ HBM accesses.

In the second step of computing $\mathbf{P} = \operatorname{softmax}(\mathbf{S})$, the input \mathbf{S} is read from HBM and the output \mathbf{P} is written to HBM (Algorithm 0 line 2). This incurs $\Theta(N^2)$ HBM accesses.

In the last step of computing $\mathbf{O} = \mathbf{PV}$, the inputs \mathbf{P}, \mathbf{V} are read from global memory and the output \mathbf{O} is written to HBM (Algorithm 0 line 3). This incurs $\Theta(Nd+N^2)$ HBM accesses.

Overall, standard attention implementation requires $\Theta(Nd + N^2)$ global memory accesses.

We now analyze the IO complexity of streaming attention.

Following Algorithm 1, we see that each element of **K** and **V** is loaded from HBM once (Algorithm 1 line 6). We make T_c passes over **Q** and **O**, each pass loading all of **Q** and all of **O** to HBM (Algorithm 1 line 8). Therefore the number of HBM accesses is $\Theta(Nd + NdT_c) = \Theta(NdT_c)$.

We derive the conditions on the block sizes B_c and B_r . We need the blocks \mathbf{K}_j and \mathbf{V}_j of size $B_c \times d$ to fit into on-chip memory, which translates to:

$$B_c d = O(M) \Leftrightarrow B_c = O\left(\frac{M}{d}\right).$$

Similarly, we need the blocks \mathbf{Q}_i , \mathbf{O}_i of size $B_r \times d$ to fit into on-chip memory, which translates to:

$$B_r d = O(M) \Leftrightarrow B_r = O\left(\frac{M}{d}\right).$$

Finally, we need the block S_{ij} of size $B_r \times B_c$ to fit into on-chip memory, which translates to:

$$B_rB_c=O(M).$$

We therefore set:

$$B_c = \Theta\left(\frac{M}{d}\right), \qquad B_r = \Theta\left(\min\left(\frac{M}{d}, \frac{M}{B_c}\right)\right) = \Theta\left(\min\left(\frac{M}{d}, d\right)\right).$$

We then have:

$$T_c = \frac{N}{B_c} = \Theta\left(\frac{Nd}{M}\right).$$

As a result, the number of HBM accesses is:

$$\Theta\left(NdT_c\right) = \Theta\left(\frac{N^2d^2}{M}\right).$$

Proof of Proposition 3. For contradiction, suppose that there exists an algorithm that computes exact attention where the number for HBM access for all $M \in [d, Nd]$ is

$$o\left(\frac{N^2d^2}{M}\right)$$
.

In the regime of $M = \Theta(Nd)$, this results in the number of HBM accesses:

$$o\left(\frac{N^2d^2}{Nd}\right) = o(Nd).$$

However, the input to attention (matrices $\mathbf{Q}, \mathbf{K}, \mathbf{V}$) and the output \mathbf{O} have size Nd and they start out being in HBM, so if the algorithm computes exact attention it must incur at least $\Omega(Nd)$ HBM accesses. This is a contradiction.

Proof of Theorem 5. The IO complexity of the attention backward is very similar to the IO complexity of the attention forward (Theorem 2). Here we provide a sketch of the proof.

We first analyze the IO complexity of standard attention backward pass. The inputs $\mathbf{Q}, \mathbf{K}, \mathbf{V}, \mathbf{dO} \in \mathbb{R}^{N \times d}$ reside in HBM, and the at the end of the algorithm the outputs $\mathbf{dQ}, \mathbf{dK}, \mathbf{dV} \in \mathbb{R}^{N \times d}$ are written to HBM.

At each step of the standard attention backward pass, one needs to load inputs of size Nd or N^2 from HBM, and needs to write the outputs of size N^2 or Nd to HBM. This incurs $\Theta(Nd+N^2)$ HBM accesses.

We now analyze the IO complexity of FlashAttention backward pass.

Similar to Theorem 2, we see that each element of **K** and **V** is loaded from HBM once. Each element of **dK** and **dV** is only written to HBM once. We make T_c passes over **Q**, **O**, **dO**, each pass loading all of **Q**, **O**, **dO** to HBM. We also make T_c passes over **dQ**, each pass reading/writing all of **dQ** from/to HBM. Therefore the number of HBM accesses is $\Theta(Nd + NdT_c) = \Theta(NdT_c)$.

As in the proof of Theorem 2, the constraints on the block sizes are that:

$$B_c = \Theta\left(\frac{M}{d}\right), \qquad B_r = \Theta\left(\min\left(\frac{M}{d}, d\right)\right).$$

We then have:

$$T_c = \frac{N}{B_c} = \Theta\left(\frac{Nd}{M}\right).$$

As a result, the number of HBM accesses is:

$$\Theta\left(NdT_{c}\right)=\Theta\left(\frac{N^{2}d^{2}}{M}\right).$$

```
Require: Matrices \mathbf{Q}, \mathbf{K}, \mathbf{V} \in \mathbb{R}^{N \times d} in HBM, on-chip SRAM of size M, softmax scaling constant \tau \in \mathbb{R},
        masking function MASK, dropout probability p_{\text{drop}}, block sizes B_c = \left[\frac{M}{4d}\right], B_r = \min\left(\left[\frac{M}{4d}\right], d\right), block
        sparsity mask M \in \{0,1\}^{N/B_r \times N/B_c}..
   1: Initialize the pseudo-random number generator state \mathcal R and save to HBM.
  2: Initialize \mathbf{O} = (0)_{N \times d} \in \mathbb{R}^{N \times d}, \ell = (0)_N \in \mathbb{R}^N, m = (-\infty)_N \in \mathbb{R}^N \text{ in HBM.}

3: Divide \mathbf{Q} into T_r = \left\lceil \frac{N}{B_r} \right\rceil blocks \mathbf{Q}_1, \dots, \mathbf{Q}_{T_r} of size B_r \times d each, and divide \mathbf{K}, \mathbf{V} in to T_c = \left\lceil \frac{N}{B_c} \right\rceil blocks
        \mathbf{K}_1, \dots, \mathbf{K}_{T_c} and \mathbf{V}_1, \dots, \mathbf{V}_{T_c}, of size B_c \times d each.
   4: Divide \mathbf{0} into T_r blocks \mathbf{0}_i, \dots, \mathbf{0}_{T_r} of size B_r \times d each, divide \ell into T_r blocks \ell_i, \dots, \ell_{T_r} of size B_r each,
        divide m into T_r blocks m_1, \ldots, m_{T_r} of size B_r each.
        for 1 \le j \le T_c do
             Load \mathbf{K}_i, \mathbf{V}_i from HBM to on-chip SRAM.
             for 1 \le i \le T_r do
   7:
                   if M_{ij} \neq 0 then
   8:
   9:
                        Load \mathbf{Q}_i, \mathbf{O}_i, \ell_i, m_i from HBM to on-chip SRAM.
                        On chip, compute \mathbf{S}_{ij} = \tau \mathbf{Q}_i \mathbf{K}_i^T \in \mathbb{R}^{B_r \times \overline{B}_c}.
 10:
                        On chip, compute \mathbf{S}_{ij}^{\text{masked}} = \text{MASK}(\mathbf{S}_{ij}).
 11:
                        On chip, compute \tilde{m}_{ij} = \text{rowmax}(\mathbf{S}_{ij}^{\text{masked}}) \in \mathbb{R}^{B_r}, \tilde{\mathbf{P}}_{ij} = \exp(\mathbf{S}_{ij}^{\text{masked}} - \tilde{m}_{ij}) \in \mathbb{R}^{B_r \times B_c} (pointwise).
 12:
                        \tilde{\ell}_{ii} = \text{rowsum}(\tilde{\mathbf{P}}_{ii}) \in \mathbb{R}^{B_r}.
                       On chip, compute m_i^{\text{new}} = \max(m_i, \tilde{m}_{ij}) \in \mathbb{R}^{B_r}, \ell_i^{\text{new}} = e^{m_i - m_i^{\text{new}}} \ell_i + e^{\tilde{m}_{ij} - m_i^{\text{new}}} \tilde{\ell}_{ij} \in \mathbb{R}^{B_r}.
On chip, compute \tilde{\mathbf{P}}_{ij}^{\text{dropped}} = \text{dropout}(\tilde{\mathbf{P}}_{ij}, p_{\text{drop}}).
 13:
 14:
                       Write \mathbf{O}_{i} \leftarrow \operatorname{diag}(\ell_{i}^{\operatorname{new}})^{-1}(\operatorname{diag}(\ell_{i})e^{m_{i}-m_{i}^{\operatorname{new}}}\mathbf{O}_{i} + e^{\tilde{m}_{ij}-m_{i}^{\operatorname{new}}}\tilde{\mathbf{P}}_{ij}^{\operatorname{dropped}}\mathbf{V}_{j}) to HBM. Write \ell_{i} \leftarrow \ell_{i}^{\operatorname{new}}, m_{i} \leftarrow m_{i}^{\operatorname{new}} to HBM.
 15:
 16:
                   end if
 17:
             end for
 18:
 19: end for
 20: Return \mathbf{0}, \ell, m, \mathcal{R}.
```

D Extension Details

D.1 Block-sparse FlashAttention

We describe the full block-sparse FlashAttention algorithm in Algorithm 5. The algorithm is identical to Algorithm 2, except that we skip zero blocks.

We prove the IO-complexity of block-sparse FlashAttention.

Proof of Proposition 4. The proof is very similar to the proof of Theorem 2. For the block-sparse case, notice that we only need to load blocks corresponding to nonzero blocks. As a result, the number of HBM accesses are scaled by s, the fraction of nonzero blocks in the block-sparsity mask. However, for small values of s, we would still need to write the result $\mathbf{O} \in \mathbb{R}^{N \times d}$. Therefore the number of HBM accesses is

$$\Theta\left(Nd + \frac{N^2d^2}{M}s\right).$$

D.2 Potential Extensions

We discuss here a few potential extensions of the IO-aware approach to speed up deep learning training.

Multi-GPU Attention. Large language models are trained on hundreds or thousands of GPUs, and one typically splits the attention computation between 4-8 GPUs on the same node [80]. This introduces another level of memory hierarchy: beside GPU SRAM and GPU HBM, we also have the HBM of other GPUs. For

25

very long sequences, the different GPUs on the same node can cooperate to compute attention by taking into account the asymmetry of different levels of memory hierarchy.

Sparse MLP layers. Typical dense MLP layers are compute-bound and not memory-bound. To improve their efficiency, MLP layers with sparse weight matrices can be used [18]. However, many sparse MLP layers are instead memory-bound, and their speedup is often not proportional to the sparsity. We believe that an IO-aware implementation can alleviate this issue and realize the benefits of sparsity. We are excited about future work in this direction, to reduce the computational requirement of large models and improve their wall-block runtime.

Kernel machine learning. Our approach in FLASHATTENTION relies on the fact that the $N \times N$ attention matrix is a function of a low-rank matrix $\mathbf{Q}\mathbf{K}^{\top}$ (of rank $d \ll N$). As a result, we can repeatedly load the inputs \mathbf{Q} , \mathbf{K} and recompute the block of the attention matrix that we need, significantly reducing HBM access. As similar scenario happens in kernel machine learning: each element K_{ij} of the $N \times N$ kernel matrix \mathbf{K} is a function of two vectors of size $d \ll N$, as it measures the similarity between two datapoints x_i and x_j . The KeOps library [8, 28] is a successful example of how reducing memory reads/writes can speed up kernel operations. We hope that this will motivate kernel methods that focus more on reducing IOs instead of just FLOPs.

E Full Experimental Results

E.1 BERT

We train BERT-large following the training procedure and hyperparameters of the reference MLPerf 1.1 implementation. In particular, we use the LAMB optimizer with learning rate 3.75e-3, with batch size 448, trained for at most 7100 steps. The training is stopped once the validation accuracy (for masked language modeling) reaches the target 72.0%, and the wall-clock run-time is measured. We train with FP16 precision using Apex AMP (with O2 optimization level).

We compare our results with the reported training speed from Nvidia that was submitted to MLPerf 1.1 (Table 1).

We use the same train / validation data split provided by MLPerf 1.1 reference implementation. In particular, we evaluate on the same 10000 validation examples as the baseline from Nvidia.

We train the model on $8\times A100$ -80GB GPUs. Each training run takes between 16 and 19 minutes, and we average the results of 10 runs.

We see a memory saving of 1.8times (from 58GB to 32GB) for the same batch size.

In Table 7, we additionally compare against the commonly used Huggingface implementation. FlashAttention is 3.2× faster than this implementation.

Table 7: Training time of BERT-large, starting from the same initialization provided by the MLPerf benchmark, to reach the target accuracy of 72.0% on masked language modeling. Averaged over 10 runs on 8×A100 GPUs.

BERT Implementation	Training time (minutes)
Huggingface [91]	55.6 ± 3.9
Nvidia MLPerf 1.1 [63]	20.0 ± 1.5
FLASHATTENTION (ours)	17.4 ± 1.4

E.2 GPT-2

We use the standard implementations of GPT-2 [70] from Huggingface transformers library and from Nvidia's Megatron-LM repo. We follow the training recipe of the Megatron-LM repo.

We use an effective batch size of 512, and use gradient accumulation to fit into available GPU memory. We use the AdamW optimizer, with learning rate 6e-4 for GPT-2 small and 1.5e-4 for GPT-2 medium, and weight decay of 0.1. All models are trained with the same hyperparameters for 400K steps. We run all implementations with mixed-precision training (PyTorch AMP).

We use the Openwebtext dataset, with the GPT-2 BPE tokenizer. We randomly select 0.5% of the dataset as the validation set, with the rest being used as training set. This random selection of validation set is done once, and all models are evaluated on the same validation set.

We train the model on 8×A100-40GB GPUs, and we measure the wall-clock training time. Training GPT-2 small takes between 2.7-9.5 days, and training GPT-2 medium takes between 6.9-21.0 days (Table 2).

For GPT-2 small, we see a memory saving of 3.5times (from 39GB to 11GB) for the same batch size of 16 (which means we could run FlashAttention with 4× larger device batch size while keeping the global batch size of 512 the same).

In Fig. 4, we plot of the validation perplexity throughout training of GPT-2 small/medium, using either HuggingFace implementation or our FLASHATTENTION implementation. We see that FLASHATTENTION behaves the same as the baseline implementation and the validation perplexity curves of the two implementations almost lie on top of each other.

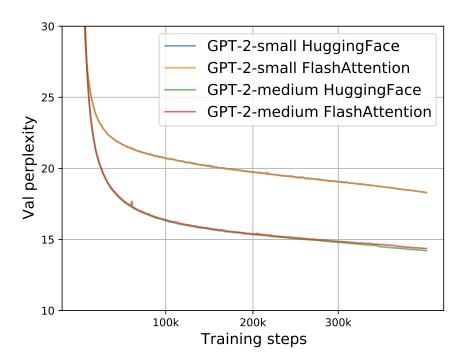


Figure 4: Validation perplexity of GPT-2 small/medium using two implementations. We confirm that FLASHATTENTION yields the same validation curves as the baseline implementation from HuggingFace.

We additionally compare the speedup of FlashAttention as we scale the number of GPUs from 1 to 8, for GPT-2 Small. All the training hyperparameters are kept the same. As we change the number of GPUs, we change the number of gradient accumulations to keep the global batch size the same (512). We use PyTorch DistributedDataParallel when there are more than 1 GPUs. Table 8 shows that the speedup is consistent, around $3.5-3.7\times$.

Table 8: Training speedup (in wallclock-time) of FlashAttention compared to Huggingface implementation on GPT-2 small as we vary the number of GPUs, measured on A100-SXM4-40GB GPUs. The speedup varies from $3.7 \times$ to $3.5 \times$.

FLASHATTENTION vs. Huggingface on GPT-2 Small	1 GPU	2 GPUs	$4 \mathrm{GPUs}$	8 GPUs
Wallclock-time speedup	3.7×	$3.6 \times$	$3.6 \times$	3.5×

Long Document Classification. For MIMIC-III and ECtHR, we follow the hyperparameters of Dai et al. [14].

E.3 LRA details

We follow the hyperparameters from the Long-range arena paper [83], the Long-range arena repo (https://github.com/google-research/long-range-arena), and the Nyströmformer reproduction [94]. To be generous to the baseline methods, if we are unable to reproduce the performance of any baseline for any of

the five tasks, we report the better performance from Tay et al. [83] or Xiong et al. [94] for that baseline on that task.

After hyperparameter tuning, almost all of the attention methods achieve similar accuracy on all of the five LRA tasks.

We run all methods with mixed-precision training, except for Performer (not stable with mixed precision) and Local Attention (implementation does not support FP16).

To calculate the overall wallclock-time speedup, we take the geometric mean of the wallclock-time speedup of each of the five tasks.

Path-X For Path-X and Path-256, we follow the hyperparameters from the PathFinder-32 experiments from the long-range arena paper [83]. For both, we first pretrain a model on Path-64. We take the checkpoint after 200 epochs, upsample its positional embedding (we duplicate the positional embeddings gridwise in space), and fine-tune it on the downstream task for 200 epochs with one epoch of linear warmup, and cosine decay of the learning rate. For Path-X, we take the best performing checkpoint (according to val accuracy), and additionally fine-tune it for 200 epochs with the same warmup and learning rate (this adds roughly 4 points of accuracy to FlashAttention for Path-X, but the model starts overfitting afterwards).

E.4 Faster Vision Transformer with FlashAttention on ImageNet

On the popular vision benchmark, ImageNet [22], we show that FlashAttention can also speedup Vision Transformers (ViT) [26] by 1.5times, where the sequence length is 196 (patch size 16×16 for 224×224 images). For longer sequences, FlashAttention yields up to 3.5× speedup.

We use the ViT-base implementation from the widely-used library timm, and replace the standard attention implementation with FLASHATTENTION. We follow the same training recipe as that of DeiT [85], which improves on the original training recipe of ViT. We measure accuracy and training time of both models (for 300 epochs) on 8×A100s. Table 9 shows that FLASHATTENTION achieves up 1.5× speed-up compared to standard attention.

Table 9: Training time of ViT-base on ImageNet for 300 epochs, on 8×A100 GPUs. Even with relatively small sequence length (196), FlashAttention still yields 1.5x speedup.

ViT-base implementation	ImageNet top-1 val accuracy	Training time (hours)
timm	81.8%	29.1
FLASHATTENTION (ours)	81.8%	19.5

We also compare ViT-Large with smaller patch sizes (i.e., longer sequence lengths). Table 10 shows that FLASHATTENTION yields $3.5\times$ speedup and saves up to 3.6x memory, compared to standard attention. Table 10: Forward + Backward time of ViT-Large on a batch of 224×224 images an A100 GPU. With longer sequence length, FLASHATTENTION yields 3.5x speedup and up to $3.6\times$ memory saving.

ViT-Large implementation	Sequence length	Batch size	Fwd + bwd time	Memory
ViT-Large (timm) patch size 8	784	32	1400ms	36GB
ViT-Large (FlashAttention) patch size 8	784	32	405ms $(3.5 \times)$	22GB
ViT-Large (timm) patch size 4	3136	2	$1200 \mathrm{ms}$	22GB
ViT-Large (FlashAttention) patch size 4	3136	2	350ms $(3.4\times)$	6GB

E.5 Comparison with Automatic Fusion

We compare FlashAttention with automatic fusion methods: NVFuser from Pytorch 1.12 (newest version at the time of writing), AOT compiler from Functorch, and TVM [11]. For context, we also include the runtime of standard implementation from Pytorch and the more optimized implementation from Megatron-LM [80].

We benchmark for batch size 16, 32, and 64, sequence length 1024, 16 heads, head dimension 64, with key-padding mask (and no dropout). The runtime is measured on an A100-SXM4-40GB GPU. Table 11 shows that FLASHATTENTION is about 2-3× faster than these methods.

Table 11: Runtime (ms) of FLASHATTENTION compared to automatic fusion methods by sequence length, with key padding masking, measured on an A100-SXM4-40GB GPU. Batch size 16, 32, and 64, sequence length 1024, 16 heads, head dimension 64. FLASHATTENTION is 2-3× faster than these methods.

Method	Batch size 16		Batch size 32			Batch size 64			
	Fwd	Bwd	Total	Fwd	Bwd	Total	Fwd	Bwd	Total
Pytorch eager mode	4.1	5.0	9.1	8.1	9.5	17.6	16.1	19.0	35.1
Pytorch JIT (NVFuser)	2.8	4.8	7.6	5.5	9.5	15.0	11.0	18.7	29.7
AOT compiler (Functorch)	2.7	4.9	7.6	5.4	9.8	15.2	10.8	19.6	30.4
$ ext{TVM}$	2.7	4.9	7.6	5.5	9.7	15.2	11.0	19.0	30.0
Megatron-LM	2.9	3.8	6.9	5.5	7.1	12.6	11.6	14.3	25.9
FLASHATTENTION	1.0	2.6	3.6	1.8	4.1	5.9	3.3	8.4	11.7

Table 12: Runtime (ms) of FLASHATTENTION compared to FMHA by sequence length, with masking and dropout, measured on an A100-SXM4-40GB GPU. Batch size 64, 16 heads, head dimension 64 (i.e., BERT-large size).

Attention Method	128	256	512
Apex FMHA forward	0.10	0.29	1.14
FLASHATTENTION forward	0.08	0.22	0.81
Apex FMHA backward	0.17	0.52	1.81
FLASHATTENTION backward	0.20	0.53	2.00
	0.27	0.81	2.95
FLASHATTENTION forward $+$ backward	0.28	0.75	2.81

E.6 Comparison with Apex FMHA

We compare our method/implementation with Apex FMHA (https://github.com/NVIDIA/apex/tree/master/apex/contrib/csrc/fmha).

When we started this project, Apex FMHA was the fastest implementation of attention (that we knew of), tailored for short sequences of length at most 512. In fact, almost all MLPerf submissions for BERT training benchmark running on Nvidia GPUs use FMHA for their model code, as of MLPerf 1.1 [60]. Since FMHA targets BERT models, it only supports head dimension 64, and only runs on A100 GPUs. FMHA fuses the attention computation dropout(softmax(MASK(\mathbf{QK}^{T}))) \mathbf{V} into one CUDA kernel. In the forward pass, it stores the attention matrix softmax(MASK(\mathbf{QK}^{T})) to HBM to be used in gradient computation. As a result, it does not offer substantial memory saving (though for shorter sequences memory footprint is often not a primary concern).

We use FMHA code as a starting point, and apply two well-established techniques (tiling and recomputation) to deal with long sequences and to save memory as mentioned in Section 3. As a result, we can support much longer sequences (e.g., up to length 64K). We also support more head dimensions (16, 32, 64, 128) and broader GPU types (all Turing and Ampere GPUs at the time of writing).

In Table 12, we compare the performance of FlashAttention and Apex FMHA for short sequences (as FMHA only supports sequence length at most 512). Generally FlashAttention is slightly faster than FMHA in the forward pass and slightly slower than FMHA in the backward pass. This is because we do not store the attention matrix in the forward pass and recompute it in the backward pass. Compared to FMHA, the overall runtime of FlashAttention is about 4% slower for sequence length 128, 8% faster for sequence length 256, and 5% faster for sequence length 512.

E.7 Roofline analysis

In Fig. 5, we include a roofline analysis of the FlashAttention forward pass, taken from Nvidia Nsight Compute (batch size 16, seglen 512, 16 heads, head dimension 64) on an A100-SXM4-40GB GPU.

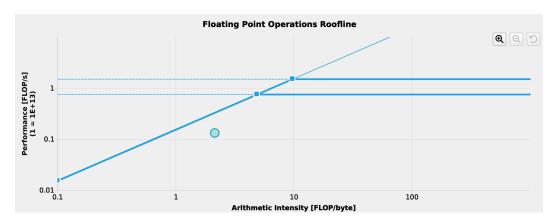


Figure 5: Roofline analysis of FlashAttention forward pass. While FlashAttention substantially speeds up attention, there is still some potential headroom to gain further speedup.

E.8 Speedup On Different Hardware and Configurations

Speedup varies between different types of GPU types and generations depending on HBM bandwidth and SRAM size. In this section, we profile FlashAttention speedup on different GPUs and configurations.

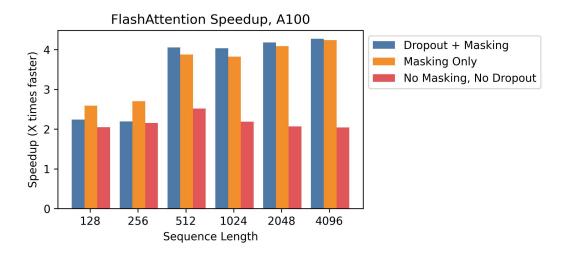


Figure 6: Speedup over standard PyTorch attention at different sequence lengths, on A100.

A100 Figure 6 shows speedup on an A100 GPU with batch size 8, head dimension 64, and 12 attention heads, across different sequence lengths. We generally see $2-4\times$ speedup, and we see more speedup when using dropout and masking due to kernel fusion.

A100, Head Dimension 128 Speedup also changes when we increase the head dimension. Each block requires more memory, so we need to use smaller block sizes to fit into SRAM. Figure 7 shows speedup with head dimension 128 on an A100 (batch size 16, 12 heads). We see less speedup overall—but we can still see significant speedup (up to $3\times$) with a causal mask, where half the blocks are masked out.

RTX 3090 Figure 8 shows speedup on an RTX 3090 GPU. Here, we use batch size 12 with 12 attention heads. We observe slightly higher speedups on the RTX 3090 (between 2.5-4.5×), since the memory bandwidth on an RTX 3090 is lower than on an A100 (roughly 900 GB/s vs. 1.5 TB/s).

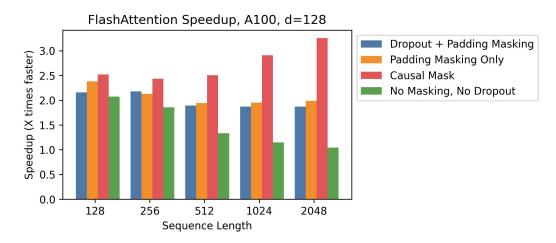


Figure 7: Speedup over standard PyTorch attention at different sequence lengths, on A100, with head dimension 128.

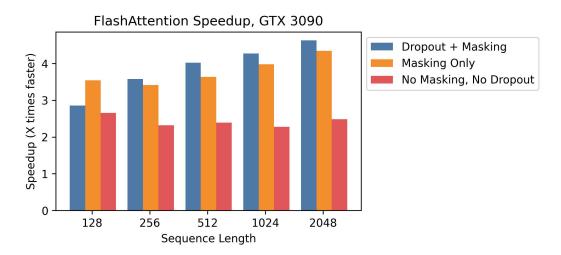


Figure 8: Speedup over standard PyTorch attention at different sequence lengths, on RTX 3090.

T4 Figure 9 shows speedup on a T4 GPU. T4 SRAM is smaller than A100, so we need to make the block sizes smaller in FlashAttention. As a result, we observe less speedup on T4, which matches the IO complexity analysis in Section 3.2. T4 GPUs are commonly used for inference, so we also report speedup on the forward pass only.

E.9 Full Benchmarking Results

We report the full benchmarking results and experimental details on A100.

Baselines We compare against reference implementations for exact attention from PyTorch/HuggingFace and Megatron, approximate attention, and sparse attention. For approximate attention, we compare against reference implementations of Reformer [53], Local Attention [71], Linformer Attention [88], Smyrf [20], and LongShortFormer (LSFormer) [98]. For sparse attention, we compare against reference implementations of Block-Sparse Attention form OpenAI [12], Longformer[3], and BigBird Attention [96]. For the approximate and sparse attention, we use a compression ratio of 1/8, or a compressed sequence length of 256, whichever is smaller.

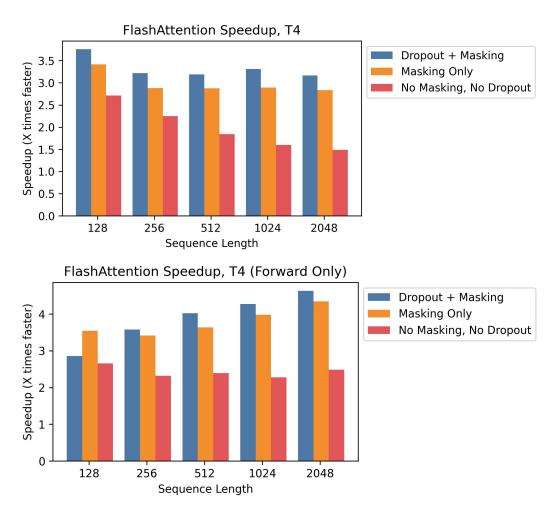


Figure 9: Speedup over standard PyTorch attention at different sequence lengths, on T4. **Top:** Combined forward pass + backward pass. **Bottom:** Forward pass only.

Setup We measure runtime and memory usage of the attention computation with 8 heads of dimension 64, and batch size 16 on a machine with one A100 GPU with 40 GB of GPU HBM. We vary sequence length in our experiments. We compute attention on random vectors for \mathbf{Q} , \mathbf{K} , and \mathbf{V} (we do not measure the projection from the hidden layer). For dropout, we use dropout 0.1; for masking, we use a padding mask with uniformly-random mask lengths between the total sequence length and the total sequence length minus 20. To measure runtime, we take the average of 100 measurements of the attention call. We only measure memory footprint once, since it does not vary between runs.

We report timing results on the forward pass, backward pass, and combined forward + backward pass. We measure each method with and without dropout, masking, or both—except for Block Sparse, Longformer, and BigBird. These methods did not successfully run the backward pass with masking due to a bug in external libraries, so we measured them without masking to be generous. We use FP16 for all measurements, except for Local Attention, whose implementation only supports FP32.

For each baseline, we increase sequence length until it runs out of memory on the GPU, except for the following exceptions: The Megatron implementation does not support sequence lengths longer than 2048. Block-Sparse (OpenAI) does not support sequence lengths longer than 4096. Longformer and BigBird do not support sequence lengths longer than 8092.

We measure memory usage on the combined forward + backward pass, without dropout or masking.

Results Table 13 summarizes all the experimental configurations and contains pointers to the results tables.

Table 13: Pointers to results tables.

Dropout	Masking	Pass	Table
Yes	Yes	Forward	Table 14
Yes	Yes	Backward	Table 15
Yes	Yes	Combined	Table 16
No	Yes	Forward	Table 17
No	Yes	Backward	Table 18
No	Yes	Combined	Table 19
Yes	No	Forward	Table 20
Yes	No	Backward	Table 21
Yes	No	Combined	Table 22
No	No	Forward	Table 23
No	No	Backward	Table 24
No	No	Combined	Table 25
No	No	Memory Usage (Combined)	Table 26

Table 14: Forward pass runtime (ms) of various exact/approximate/sparse attention mechanisms by sequence length, with dropout and masking. Best in **bold**, second best <u>underlined</u>.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
PyTorch Attention	0.36	0.34	0.78	2.54	9.33	36.33	-	-	-	
Megatron	0.40	0.40	1.10	3.65	16.19	-	-	-	-	-
Reformer	2.03	3.15	5.67	11.02	22.59	46.14	97.38	212.13	-	
Local Attention	0.83	0.86	1.01	2.20	7.13	14.32	28.60	57.79	117.67	-
Linformer	0.67	0.52	0.69	0.71	1.65	3.18	6.15	12.16	24.17	52.39
Smyrf	2.27	2.34	3.91	$\overline{7.44}$	14.71	29.22	58.27	116.41		
LSformer	1.18	1.27	1.34	3.38	11.40	22.55	44.95	89.76	179.66	-
Block Sparse	1.12	1.11	2.13	2.77	6.95	20.91	-	-	-	
Longformer	1.22	1.14	1.08	1.95	5.72	12.98	-	-	-	-
$_{ m BigBird}$	1.13	1.12	1.12	1.77	6.03	13.68	-	-	-	-
FLASHATTENTION	0.04	0.06	0.21	0.82	2.85	10.41	41.74	167.19	670.76	2682.35
Block-Sparse FlashAttention	0.06	0.06	0.06	0.12	0.44	0.86	1.70	3.29	6.55	13.34

Table 15: Backward pass runtime (ms) of various exact/approximate/sparse attention mechanisms by sequence length, with dropout and masking. Best in **bold**, second best <u>underlined</u>.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
PyTorch Attention	0.37	0.49	1.66	5.81	22.32	87.67	-	-	-	-
Megatron	0.35	0.32	0.77	2.42	8.43	-	-	-	-	-
Reformer	2.37	4.59	8.91	17.68	35.13	70.05	140.01	-	-	-
Local Attention	0.55	0.62	1.49	4.03	13.78	27.61	55.20	110.27	221.40	-
Linformer	0.89	0.80	0.81	0.93	2.48	4.75	9.29	18.27	36.53	-
Smyrf	1.41	2.83	5.43	10.72	21.25	42.31	84.48	168.95	-	-
LSformer	1.75	1.76	3.01	7.50	20.07	39.08	76.39	150.82	-	-
Block Sparse	1.29	1.28	2.18	3.04	7.27	21.16	-	-	-	-
Longformer	1.27	1.31	1.29	2.04	5.24	10.74	25.95	-	-	-
$_{ m BigBird}$	1.33	1.28	1.32	1.81	5.55	11.44	27.45	-	-	-
FLASHATTENTION	0.30	0.26	0.68	2.02	6.84	26.89	105.70	418.96	1666.89	6660.44
Block-Sparse FlashAttention	0.30	0.27	0.29	0.59	1.50	2.94	5.82	11.85	23.98	47.61

Table 16: Forward pass + backward pass runtime (ms) of various exact/approximate/sparse attention mechanisms by sequence length, with dropout and masking. Best in **bold**, second best underlined.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
PyTorch Attention	0.84	0.86	2.35	8.29	31.75	124.19	-	-	-	
Megatron	0.87	0.89	1.33	4.21	16.50	-	-	-	-	-
Reformer	4.30	7.76	14.60	28.74	57.79	116.34	237.57	-	-	-
Local Attention	1.40	1.60	2.06	6.06	20.94	42.01	84.08	168.48	339.45	-
Linformer	1.57	1.49	1.55	1.60	4.19	8.04	15.71	30.92	61.47	-
Smyrf	3.41	5.08	9.35	18.18	36.03	71.68	143.04	285.87	-	-
LSformer	3.08	3.10	4.26	10.90	31.59	61.72	121.51	241.18	-	-
Block Sparse	2.54	2.52	3.71	5.44	13.29	39.19	-	-	-	-
Longformer	2.47	2.49	2.51	3.10	10.39	22.49	60.44	-	-	-
$_{ m BigBird}$	2.51	2.49	2.52	3.40	10.97	23.89	63.28	-	-	-
FLASHATTENTION	0.43	0.41	0.95	2.55	9.56	37.49	147.75	586.61	2339.11	9341.30
Block-Sparse FlashAttention	0.44	0.44	$\overline{0.45}$	0.89	1.95	4.12	7.64	16.60	32.73	64.11

Table 17: Forward pass runtime (ms) of various exact/approximate/sparse attention mechanisms by sequence length, with masking. Best in bold, second best <u>underlined</u>.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
PyTorch Attention	0.30	0.30	0.63	1.93	7.08	27.45	112.90	-	-	-
${f Megatron}$	0.45	0.41	0.43	1.52	5.80	-	-	-	-	-
Reformer	1.87	3.00	5.37	10.43	21.40	43.83	92.80	203.24	-	-
Local Attention	0.70	0.81	1.02	2.09	6.64	13.34	26.77	54.02	110.11	-
Linformer	0.63	0.50	0.67	0.65	1.36	2.60	5.04	9.92	19.69	43.47
Smyrf	2.38	2.32	3.76	7.16	14.14	28.09	55.98	$1\overline{11.73}$	-	_
LSformer	1.22	1.29	1.44	3.28	10.99	21.72	43.29	86.32	172.76	-
Block Sparse	0.96	1.04	1.66	2.16	5.41	16.15	-	-	-	-
Longformer	0.99	0.98	0.99	1.56	4.79	11.07	32.98	-	-	-
BigBird	0.96	1.02	1.02	1.48	5.05	11.59	34.16	-	-	-
FLASHATTENTION	0.03	0.04	0.17	0.68	2.28	8.40	33.55	134.14	537.50	2150.88
Block-Sparse FlashAttention	0.05	0.04	$\overline{0.05}$	0.11	0.35	0.68	1.33	2.54	5.34	10.73

Table 18: Backward pass runtime (ms) of various exact/approximate/sparse attention mechanisms by sequence length, with masking. Best in bold, second best <u>underlined</u>.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
PyTorch Attention	0.44	0.46	1.53	5.33	20.34	79.87		-	_	
Megatron	0.29	0.31	0.65	1.95	6.49	-	-	-	-	-
Reformer	2.31	4.47	8.68	17.20	34.14	68.09	136.02	-	-	-
Local Attention	0.51	0.62	1.30	3.81	13.33	26.72	53.41	106.82	214.15	-
Linformer	0.76	0.81	0.94	0.87	2.24	4.25	8.35	16.38	32.67	72.11
Smyrf	1.34	2.77	5.30	10.46	20.73	41.27	82.41	164.86	-	-
LSformer	1.66	1.61	3.09	7.42	19.68	38.35	74.92	147.86	-	-
Block Sparse	1.24	1.25	2.04	2.91	6.78	19.67	-	-	-	-
Longformer	1.27	1.23	1.24	1.85	4.99	10.21	24.89	-	-	-
$\operatorname{BigBird}$	1.43	1.50	1.44	1.69	5.25	10.86	26.26	-	-	-
FLASHATTENTION	0.21	0.22	0.62	1.84	5.77	22.25	86.21	338.91	1343.91	5361.09
Block-Sparse FlashAttention	0.22	0.22	0.26	0.57	1.55	3.13	5.98	12.21	23.49	47.85

Table 19: Forward pass + backward pass runtime (ms) of various exact/approximate/sparse attention mechanisms by sequence length, with masking. Best in **bold**, second best <u>underlined</u>.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
PyTorch Attention	0.80	0.81	2.08	7.23	27.51	107.58	-	-	-	-
Megatron	0.81	0.83	1.09	3.36	12.39	-	-	-	-	-
Reformer	4.16	7.46	14.06	27.68	55.66	112.15	229.37	-	-	-
Local Attention	1.39	1.68	2.08	5.83	20.04	40.16	80.44	161.35	325.11	-
Linformer	1.51	1.42	1.56	1.67	3.67	6.99	13.63	26.77	53.36	117.56
Smyrf	3.38	4.93	9.07	17.66	34.94	69.55	138.72	277.41	-	-
LSformer	3.08	3.10	4.26	10.90	31.59	61.72	121.51	241.18	-	-
Block Sparse	2.39	2.40	3.31	5.02	12.25	35.94	-	-	-	-
Longformer	2.36	2.34	2.38	2.94	9.83	21.35	58.12	-	-	-
$_{ m BigBird}$	2.35	2.35	2.37	3.25	10.36	22.57	60.63	-	-	-
FLASHATTENTION	0.32	0.30	0.83	2.37	7.95	30.77	119.98	473.65	1883.43	7513.01
Block-Sparse FlashAttention	0.34	0.34	0.36	0.69	1.85	3.89	7.16	14.85	30.46	60.03

Table 20: Forward pass runtime (ms) of various exact/approximate/sparse attention mechanisms by sequence length, with dropout. Best in **bold**, second best underlined.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
PyTorch Attention	0.26	0.24	0.57	1.80	6.56	25.34	-	-	-	-
Megatron	0.27	$\overline{0.27}$	0.56	1.88	6.56	-	-	-	-	-
Reformer	1.83	2.96	5.31	10.33	21.19	43.42	91.96	201.34	-	-
Local Attention	0.51	0.60	0.78	2.01	6.23	12.52	25.07	50.50	102.18	-
Linformer	0.47	0.37	0.49	0.52	1.37	2.65	5.12	10.13	20.25	44.16
Smyrf	2.12	2.01	3.15	5.97	11.83	23.36	46.48	92.72	-	-
LSformer	1.28	1.33	1.51	3.39	11.40	22.54	44.96	89.85	179.73	-
Block Sparse	1.03	1.00	1.72	2.39	5.96	17.88	-	-	-	-
Longformer	1.02	1.03	1.03	1.73	5.10	11.63	34.22	-	-	-
$_{ m BigBird}$	0.99	1.03	1.01	1.58	5.36	12.27	35.56	-	-	-
FLASHATTENTION	0.10	0.10	0.22	0.83	2.81	10.38	41.63	167.01	668.74	2678.11
Block-Sparse FlashAttention	0.54	0.51	0.68	0.61	0.67	1.10	1.89	3.71	7.18	14.41

Table 21: Backward pass runtime (ms) of various exact/approximate/sparse attention mechanisms by sequence length, with dropout. Best in bold, second best <u>underlined</u>.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
PyTorch Attention	0.44	0.35	0.90	2.94	10.77	41.67	-	-	-	-
Megatron	0.28	0.33	0.92	2.94	10.80	-	-	-	-	-
Reformer	2.24	4.34	8.39	16.62	33.02	65.77	131.52	-	-	-
Local Attention	0.51	0.58	1.41	3.71	12.96	25.98	51.94	103.72	207.78	-
Linformer	0.84	0.74	0.79	0.85	2.28	4.37	8.66	17.02	33.78	-
Smyrf	1.27	2.56	4.90	9.66	19.16	38.13	76.17	152.39		-
LSformer	1.67	1.77	3.03	7.52	20.10	39.13	76.35	150.83	-	-
Block Sparse	1.27	1.36	2.15	3.04	7.27	21.18	-	-	-	-
Longformer	1.28	1.34	1.38	1.98	5.24	10.74	25.95	-	-	-
$\overline{ ext{BigBird}}$	1.48	1.47	1.50	1.81	5.57	11.38	27.43	-	-	-
FLASHATTENTION	0.15	0.18	0.58	1.86	6.50	26.21	104.27	416.10	1661.92	6643.01
Block-Sparse FlashAttention	0.17	0.17	0.17	0.40	1.10	2.04	4.43	9.33	18.28	37.31

Table 22: Forward pass + backward pass runtime (ms) of various exact/approximate/sparse attention mechanisms by sequence length, with dropout. Best in bold, second best <u>underlined</u>.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
PyTorch Attention	0.66	0.67	1.43	4.82	17.47	67.29	-	-	-	-
Megatron	0.88	0.90	1.49	4.73	17.41	-	-	-	-	-
Reformer	4.06	7.28	13.68	26.98	54.27	109.39	223.80	-	-	-
Local Attention	1.09	1.40	1.99	5.61	19.23	38.62	77.30	154.63	311.12	-
Linformer	1.31	1.21	1.30	1.39	3.73	7.15	14.05	27.69	55.00	-
Smyrf	3.00	4.37	8.05	15.66	31.04	61.64	123.04	245.65	-	-
LSformer	3.07	3.17	4.31	10.89	31.54	61.78	121.56	240.94	-	-
Block Sparse	2.54	2.52	3.71	5.44	13.29	39.19	-	-	-	-
Longformer	2.47	2.49	2.51	3.10	10.39	22.49	60.44	-	-	-
$\operatorname{BigBird}$	2.51	2.49	2.52	3.40	10.97	23.89	63.28	-	-	-
FLASHATTENTION	0.35	0.36	0.80	2.52	9.16	36.70	146.13	583.45	2332.01	9323.63
Block-Sparse FlashAttention	0.91	0.83	0.94	0.92	1.83	3.50	7.02	13.56	26.71	53.92

Table 23: Forward pass runtime (ms) of various exact/approximate/sparse attention mechanisms by sequence length. Best in **bold**, second best <u>underlined</u>.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
PyTorch Attention	0.21	0.22	0.43	1.27	4.32	16.47	67.77	-	-	-
Megatron	0.24	0.26	0.42	1.33	4.28	-	-	-	-	-
Reformer	1.77	2.82	5.01	9.74	20.03	41.11	87.39	192.40	-	-
Local Attention	0.48	0.57	0.80	1.90	5.76	11.56	23.13	46.65	94.74	-
Linformer	0.46	0.36	0.45	0.50	1.09	2.09	4.01	7.90	15.70	35.40
Smyrf	1.94	1.96	3.01	5.69	11.26	22.23	44.21	88.22	-	-
LSformer	1.21	1.34	1.34	3.31	11.01	21.71	43.27	86.32	172.85	-
Block Sparse	0.96	1.04	1.66	2.16	5.41	16.15	-	-	-	-
Longformer	0.99	0.98	0.99	1.56	4.79	11.07	32.98	-	-	-
$\operatorname{BigBird}$	0.96	1.02	1.02	1.48	5.05	11.59	34.16	-	-	-
FLASHATTENTION	0.08	0.09	0.18	0.68	2.40	8.42	33.54	134.03	535.95	2147.05
Block-Sparse FlashAttention	0.56	0.52	0.63	0.65	0.61	0.96	1.69	3.02	5.69	11.77

Table 24: Backward pass runtime (ms) of various exact/approximate/sparse attention mechanisms by sequence length. Best in **bold**, second best <u>underlined</u>.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
PyTorch Attention	0.26	0.29	0.78	2.44	8.82	33.87	-	-	-	-
Megatron	0.29	0.30	0.80	2.59	8.86	-	-	-	-	-
Reformer	2.18	4.21	8.14	16.12	32.02	63.84	127.60	-	-	-
Local Attention	0.51	0.64	1.28	3.60	12.52	25.08	50.22	100.23	200.66	-
Linformer	0.69	0.76	0.69	0.80	2.04	3.88	7.67	15.04	30.11	63.15
Smyrf	1.24	2.49	4.77	9.42	18.65	37.12	74.15	148.35	-	-
LSformer	1.68	1.61	3.02	7.40	19.72	38.27	74.89	147.99	-	-
Block Sparse	1.24	1.25	2.04	2.91	6.78	19.67	-	-	-	-
Longformer	1.27	1.23	1.24	1.85	4.99	10.21	24.89	-	-	-
BigBird	1.43	1.50	1.44	1.69	5.25	10.86	26.26	-	-	-
FLASHATTENTION	0.11	0.16	0.52	1.62	5.45	21.57	84.75	336.00	1338.56	5343.19
Block-Sparse FlashAttention	0.11	$\overline{0.12}$	$\overline{0.16}$	0.38	1.20	2.34	4.69	9.10	18.74	37.04

Table 25: Forward pass + backward pass runtime (ms) of various exact/approximate/sparse attention mechanisms by sequence length. Best in **bold**, second best <u>underlined</u>.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
PyTorch Attention	0.67	0.70	1.18	3.67	13.22	50.44	-	-	-	_
Megatron	0.74	0.65	1.23	3.80	13.21	-	-	-	-	-
Reformer	3.93	7.01	13.15	25.89	52.09	105.00	215.13	-	-	-
Local Attention	1.09	1.27	1.99	5.38	18.32	36.77	73.67	147.29	296.35	-
Linformer	1.31	1.25	1.30	1.29	3.20	6.10	11.93	23.39	46.72	100.52
Smyrf	2.98	4.23	7.78	15.12	29.96	59.45	118.60	$2\overline{37.02}$		
LSformer	3.03	3.05	4.26	10.70	30.77	60.15	118.33	234.94	-	-
Block Sparse	2.39	2.40	3.31	5.02	12.25	35.94	-	-	-	-
Longformer	2.36	2.34	2.38	2.94	9.83	21.35	58.12	-	-	-
$_{ m BigBird}$	2.35	2.35	2.37	3.25	10.36	22.57	60.63	-	-	-
FLASHATTENTION	0.31	0.31	0.73	2.29	7.64	30.09	118.50	470.51	1876.08	7492.85
Block-Sparse FlashAttention	0.74	0.77	0.82	0.88	1.71	3.21	6.56	12.60	24.93	50.39

Table 26: Memory usage (MB) of various exact/approximate/sparse attention mechanisms by sequence length. Best in **bold**, second best $\underline{\text{underlined}}$.

Attention Method	128	256	512	1024	2048	4096	8192	16384	32768	65536
PyTorch Attention	36	104	336	1184	4416	17024	-	-	-	-
Megatron	36	104	336	1184	4416	-	-	-	-	-
Reformer	377	754	1508	3016	6033	12067	24134	-	-	-
Local Attention	53	110	232	592	1696	3392	6784	13568	27136	-
Linformer	25	52	114	287	832	1652	3292	6572	13132	26252
Smyrf	217	434	868	1737	3474	6947	13894	27788	-	-
LSformer	72	152	333	796	2540	5068	10125	20240	-	-
Block Sparse	33	82	228	408	910	2401	-	-	-	-
Longformer	30	61	124	277	681	1370	2748	-	-	-
BigBird	33	66	131	294	708	1431	2872	-	-	-
FLASHATTENTION	22	44	104	209	418	836	1672	3344	6688	13376
Block-Sparse FlashAttention	<u>22</u>	$\underline{44}$	104	209	418	836	1672	3344	6690	13384