

---

# Ring Attention with Blockwise Transformers for Near-Infinite Context

---

Anonymous Author(s)

Affiliation  
Address  
email

## Abstract

1 Transformers have emerged as the architecture of choice for many state-of-the-art AI models,  
2 showcasing exceptional performance across a wide range of AI applications. However, the  
3 memory demands imposed by Transformers limit their ability to handle long sequences, thereby  
4 creating challenges for tasks involving extended sequences or long-term dependencies. We  
5 present a distinct approach, Ring Attention, which leverages blockwise computation of self-  
6 attention to distribute long sequences across multiple devices while concurrently overlapping the  
7 communication of key-value blocks with the computation of blockwise attention. By processing  
8 longer input sequences while maintaining memory efficiency, Ring Attention enables training  
9 and inference of sequences that are device count times longer than those of prior memory-  
10 efficient Transformers, effectively eliminating the memory constraints imposed by individual  
11 devices. Extensive experiments on language modeling tasks demonstrate the effectiveness of  
12 Ring Attention in allowing large sequence input size and improving performance.

## 13 1 Introduction

14 Transformers [35] have become the backbone of many state-of-the-art AI systems that have demon-  
15 strated impressive performance across a wide range of AI problems. Transformers achieve this success  
16 through their architecture design that uses self-attention and position-wise feedforward mechanisms.  
17 These components facilitate the efficient capture of long-range dependencies between input  
18 tokens, and enable scalability through highly parallel computations.

21 However, scaling up the context length of Trans-  
22 formers is a challenge [26], since the inherited  
23 architecture design of Transformers, *i.e.* the self-  
24 attention has memory cost quadratic in the input  
25 sequence length, which makes it challenging to  
26 scale to longer input sequences. Large context  
27 Transformers are essential for tackling a diverse  
28 array of AI challenges, ranging from processing  
29 books and high-resolution images to analyzing  
30 long videos and complex codebases. They excel  
31 at extracting information from the intercon-  
32 nected web and hyperlinked content, and are  
33 crucial for handling complex scientific experi-  
34 ment data. There have been emerging use cases  
35 of language models with significantly expanded  
36 context than before: GPT-3.5 [29] with context

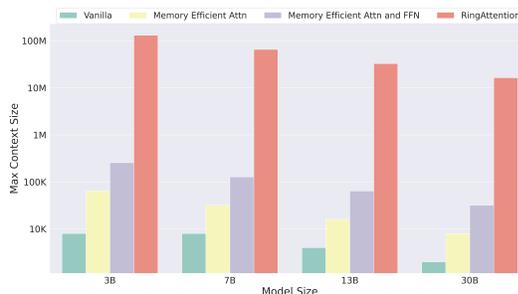


Figure 1: Maximum context length on TPUv4-512 (32GB memory on each TPUv4). Baselines are vanilla transformers [35], transformers with memory efficient attention [27], and memory efficient attention and feedforward (blockwise parallel transformers) [22]. Our proposed approach Ring Attention allows training 512 times longer sequence than prior SOTAs and enables the training of sequences that exceed 100 million in length without making approximations to attention.

37 length 16K, GPT-4 [26] with context length 32k, MosaicML’s MPT [24] with context length 65k,  
38 and Anthropic’s Claude [1] with context length 100k.

39 Driven by the significance, there has been surging research interests in reducing memory cost. One  
40 line of research leverages the observation that the softmax matrix in self-attention can be computed  
41 without materializing the full matrix [23] which has led to the development of blockwise computation  
42 of self-attention and feedforward [27, 9, 22] without making approximations. Despite the reduced  
43 memory, a significant challenge still arises from storing the output of each layer. This necessity arises  
44 from self-attention’s inherent nature, involving interactions among all elements (n to n interactions).  
45 The subsequent layer’s self-attention relies on accessing all of the prior layer’s outputs. Failing to  
46 do so would increase computational costs cubically, as every output must be recomputed for each  
47 sequence element, rendering it impractical for longer sequences. To put the memory demand in  
48 perspective, even when dealing with a batch size of 1, processing 100 million tokens requires over  
49 10,000GB of memory for a modest model with a hidden size of 1024. This is much greater than the  
50 capacity contemporary GPUs, which typically have less than 100GB of high-bandwidth memory  
51 (HBM).

52 To tackle this challenge, we make a key observation: by performing self-attention and feedforward  
53 network computations in a blockwise fashion [22], we can distribute sequence dimensions across  
54 multiple devices, allowing concurrent computation and communication. This insight stems from  
55 the fact that when we compute the attention on a block-by-block basis, the results are invariant to  
56 the ordering of these blockwise computations. Our method distributes the outer loop of computing  
57 blockwise attention among hosts, with each device managing its respective input block. For the inner  
58 loop, every device computes blockwise attention and feedforward operations specific to its designated  
59 input block. Host devices form a conceptual ring, where during the inner loop, each device sends  
60 a copy of its key-value blocks being used for blockwise computation to the next device in the ring,  
61 while simultaneously receiving key-value blocks from the previous one. Because block computations  
62 take longer than block transfers, overlapping these processes results in no added overhead compared  
63 to standard transformers. By doing so, each device requires memory only proportional to the block  
64 size, which is independent of the original input sequence length. This effectively eliminates the  
65 memory constraints imposed by individual devices. Since our approach overlaps the communication  
66 of key-value blocks between hosts in a ring with blockwise computation, we name it Ring Attention.

67 We evaluate the effectiveness of our approach on language modeling benchmarks. Our experiments  
68 show that Ring Attention can reduce the memory requirements of Transformers, enabling us to  
69 train more than 500 times longer sequence than prior memory efficient state-of-the-arts and enables  
70 the training of sequences that exceed 100 million in length without making approximations to  
71 attention. Importantly, Ring Attention eliminates the memory constraints imposed by individual  
72 devices, empowering the training and inference of sequences with lengths that scale in proportion to  
73 the number of devices, essentially achieving near-infinite context size.

74 Our contributions are twofold: (a) proposing a memory efficient transformers architecture that allows  
75 the context length to scale linearly with the number of devices while maintaining performance, elimi-  
76 nating the memory bottleneck imposed by individual devices, and (b) demonstrating the effectiveness  
77 of our approach through extensive experiments.

## 78 2 Large Context Memory Constraint

79 Given input sequences  $Q, K, V \in \mathbb{R}^{s \times d}$  where  $s$  is the sequence length and  $d$  is the head dimension.  
80 We compute the matrix of outputs as:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d}}\right)V,$$

81 where softmax is applied row-wise. Each self-attention sub-layer is accompanied with a feedforward  
82 network, which is applied to each position separately and identically. This consists of two linear  
83 transformations with a ReLU activation in between.

$$\text{FFN}(x) = \max(0, xW_1 + b_1)W_2 + b_2.$$

84 **Blockwise Parallel Transformers.** Prior state-of-the-arts have led to substantial reductions in mem-  
85 ory utilization, achieved through innovative techniques that enable attention computation without full

86 materialization by computing attention in a block by block manner [27, 9, 22]. These advancements  
87 lowered the memory overhead of attention to  $2bsh$  Bytes per layer, where  $b$  represents the batch  
88 size,  $s$  denotes the sequence length, and  $h$  stands for the hidden size of the model. To further reduce  
89 memory usage, blockwise parallel transformer (BPT) [22] introduced a strategy where the feedfor-  
90 ward network associated with each self-attention sub-layer is computed in a block-wise fashion. This  
91 approach effectively limits the maximum activation size of feedforward network from  $8bsh$  to  $2bsh$ .  
92 For a more detailed analysis of memory efficiency, please refer to the discussion provided therein. In  
93 summary, the state-of-the-art transformer layer’s memory cost of activation is  $2bsh$ .

94 **Large Output of Each Layer.** While BPT significantly reduces memory demand in Transformers, it  
95 still presents a major challenge for scaling up context length because it requires storing the output  
96 of each layer. This storage is crucial due to the inherent nature of self-attention, which involves  
97 interactions among all elements ( $n$  to  $n$  interactions). Without these stored outputs, the subsequent  
98 layer’s self-attention becomes computationally impractical, necessitating recomputation for each  
99 sequence element. To put it simply, processing 100 million tokens with a batch size of 1 requires  
100 over 10,000GB of memory even for a modest model with a hidden size of 1024. In contrast, modern  
101 GPUs typically provide less than 100GB of high-bandwidth memory (HBM), and the prospects for  
102 significant GPU HBM expansion are hindered by physical limitations and high manufacturing costs.

### 103 3 Ring Attention

104 Our primary objective is to eliminate the memory constraints imposed by individual devices by  
105 efficiently distribute long sequences across multiple hosts without adding overhead. To achieve this  
106 goal, we propose an enhancement to the blockwise parallel transformers (BPT) framework [22].  
107 When distributing an input sequence across different hosts, each host is responsible for running one  
108 element of the outer loop of blockwise attention corresponding to its designated block, as well as the  
109 feedforward network specific to that block. These operations do not necessitate communication with  
110 other hosts. However, a challenge arises in the inner loop, which involves key-value block interactions  
111 that require fetching blocks from other hosts. Since each host possesses only one key-value block,  
112 the naive approach of fetching blocks from other hosts results in two significant issues. Firstly,  
113 it introduces a computation delay as the system waits to receive the necessary key-value blocks.  
114 Secondly, the accumulation of key-value blocks leads to increased memory usage, which defeats the  
115 purpose of reducing memory cost.

116 **Ring-Based Blockwise Attention.** To tackle the aforementioned challenges, we leverage the per-  
117 mutation invariance property of the inner loop’s key-value block operations. This property stems  
118 from the fact that the self-attention between a query block and a group of key-value blocks can be  
119 computed in any order, as long as the statistics of each block are combined correctly for rescaling.  
120 We leverage this property by conceptualizing all hosts as forming a ring structure: host-1, host-2, ...,  
121 host- $N$ . As we compute blockwise attention and feedforward, each host efficiently coordinates by  
122 concurrently sending key-value blocks being used for attention computation to the next host while  
123 receiving key-value blocks from the preceding host, effectively overlapping transferring of blocks  
124 with blockwise computation. Concretely, for any host- $i$ , during the computation of attention between  
125 its query block and a key-value block, it concurrently sends key-value blocks to the next host- $(i + 1)$   
126 while receiving key-value blocks from the preceding host- $(i - 1)$ . If the computation time exceeds  
127 the time required for transferring key-value blocks, this results in no additional communication cost.  
128 This overlapping mechanism applies to both forward and backward passes of our approach since the  
129 same operations and techniques can be used.

130 **Arithmetic Intensity Between Hosts.** In order to determine the minimal required block size to  
131 overlap transferring with computation, assume that each host has  $F$  FLOPS and that the bandwidth  
132 between hosts is denoted as  $B$ . It’s worth noting that our approach involves interactions only with  
133 the immediately previous and next hosts in a circular configuration, thus our analysis applies to both  
134 GPU all-to-all topology and TPU torus topology. Let’s consider the variables: block size denoted  
135 as  $c$  and hidden size as  $d$ . When computing blockwise self-attention, we require  $2dc^2$  FLOPs for  
136 calculating attention scores using queries and keys, and an additional  $2dc^2$  FLOPs for multiplying  
137 these attention scores by values. In total, the computation demands amount to  $4dc^2$  FLOPs. We  
138 exclude the projection of queries, keys, and values, as well as blockwise feedforward operations,  
139 since they only add compute complexity without any communication costs between hosts. This  
140 simplification leads to more stringent condition and does not compromise the validity of our approach.



Figure 2: **Top (a):** In the framework of Ring Attention, key-value blocks traverse through hosts to facilitate attention and feedforward computations in a block-by-block fashion. As we compute attention, each host concurrently sends key-value blocks to the next host while receive key-value blocks from the preceding host, effectively overlapping communication with computation. **Bottom (b):** Ring Attention is the same as the original Transformer but with a different way of organizing the compute. In the diagram, we explain this by showing that the current device holds the left column first query block; then we iterate over the same key-value blocks sequence positioned horizontally. The query block, and the bottle middle key-value blocks, are used to compute self-attention (yellow box), whose output is pass to feedforward network (cyan box).

141 On the communication front, both key and value blocks require a total of  $2cd$  bytes. Thus, the  
 142 combined communication demand is  $4cd$  bytes. To achieve an overlap between communication and  
 143 computation, the following condition must hold:  $4dc^2/F \geq 4cd/B$ . This implies that the block size,  
 144 denoted as  $c$ , should be greater than or equal to  $F/B$ . Effectively, this means that the block size needs  
 145 to be larger than the ratio of FLOPs over bandwidth.

146 **Memory Requirement.** A host needs to store multiple blocks, including one block size to store  
 147 the current query block, two block sizes for the current key and value blocks, and two block sizes  
 148 for receiving key and value blocks. Furthermore, storing the output of blockwise attention and  
 149 feedforward necessitates one block size, as the output retains the shape of the query block. Therefore,  
 150 a total of six blocks are required, which translates to  $6bch$  bytes of memory. It's worth noting that  
 151 the blockwise feedforward network has a maximum activation size of  $2bch$  [22]. Consequently, the

Table 1: Comparison of maximum activation sizes among different Transformer architectures. Here,  $b$  is batch size,  $h$  is hidden dimension,  $n$  is number of head,  $s$  is sequence length,  $c$  is block size, the block size ( $c$ ) is independent of the input sequence length ( $s$ ). The comparison is between vanilla Transformer [35], memory efficient attention [27], memory efficient attention and feedforward [22], and our proposed approach Ring Attention. Numbers are shown in Bytes per layer, assuming *bfloat16* precision.

Layer Type	Self-Attention	FeedForward	Total
Vanilla	$2bns^2$	$8bsh$	$2bhs^2$
Memory efficient attention	$2bsh + 4bch$	$8bsh$	$8bsh$
Memory efficient attention and feedforward	$2bsh$	$2bsh$	$2bsh$
Ring Attention	$6bch$	$2bch$	$6bch$

Table 2: Minimal sequence length needed on each device. Interconnect Bandwidth is the uni-directional bandwidth between hosts, *i.e.*, NVLink / InfiniBand bandwidth between GPUs and ICI bandwidth between TPUs. Minimal sequence length  $s = 6c$  and minimal block size  $c = \text{FLOPS}/\text{Bandwidth}$ .

Spec Per Host	FLOPS	HBM	Interconnect Bandwidth	Minimal Blocksize	Minimal Sequence Len
	(TF)	(GB)	(GB/s)	( $\times 1e3$ )	( $\times 1e3$ )
A100 NVLink	312	80	300	1.0	6.2
A100 InfiniBand	312	80	100	3.1	18.7
TPU v3	123	16	112	1.1	6.6
TPU v4	275	32	268	1.0	6.2
TPU v5e	196	16	186	1.1	6.3

152 total maximum activation size remains at  $6bch$  bytes. Table 1 provides a detailed comparison of the  
 153 memory costs between our method and other approaches. Notably, our method exhibits the advantage  
 154 of linear memory scaling with respect to the block size  $c$ , and is independent of the input sequence  
 155 length  $s$ .

156 Our analysis shows that the model needs to fit in  $s = 6c$  sequence length, *i.e.*, six times of minimal  
 157 block size. Requirements on popular computing servers as shown in Table 3, the required minimal  
 158 sequence length to be fit in each host is between 6K to 20K. This requirement is easy to meet  
 159 using blockwise computation of attention and feedforward [22], which we will show in experiment  
 160 section 5.

161 **Algorithm and Implementation.** Algorithm 1 provides the pseudocode of the algorithm. Ring  
 162 Attention is compatible with existing code for memory efficient transformers: Ring Attention just  
 163 needs to call whatever available memory efficient computation locally on each host, and overlap the  
 164 communication of key-value blocks between hosts with blockwise computation. We use collective  
 165 operation `jax.lax.ppermute` to send and receive key value blocks between nearby hosts. A Jax  
 166 implementation is provided in Appendix A.

## 167 4 Setting

168 We evaluate the impact of using Ring Attention in improving Transformer models by benchmarking  
 169 maximum sequence length and model flops utilization.

170 **Model Configuration.** Our study is built upon the LLaMA architecture, we consider 3B, 7B, 13B,  
 171 and 30B model sizes in our experiments.

172 **Baselines.** We evaluate our method by comparing it with vanilla transformers [35] which computes  
 173 self-attention by materializing the attention matrix and computes the feedforward network normally,  
 174 transformers with memory efficient attention [27] and its efficient CUDA implementation [9], and  
 175 transformers with both memory efficient attention and feedforward [22].

---

**Algorithm 1** Reducing Transformers Memory Cost with Ring Attention.

---

**Required:** Input sequence  $x$ . Number of hosts  $N_h$ .  
Initialize  
Split input sequence into  $N_h$  blocks that each host has one input block.  
Compute query, key, and value for its input block on each host.  
**for** Each transformer layer **do**  
  **for**  $count = 1$  to  $N_h - 1$  **do**  
    **for** For each host concurrently. **do**  
      Compute memory efficient attention incrementally using local query, key, value blocks.  
      Send key and value blocks to next host and receive key and value blocks from previous host.  
    **end for**  
  **end for**  
  **for** For each host concurrently. **do**  
    Compute memory efficient feedforward using local attention output.  
  **end for**  
**end for**

---

176 **Training Configuration.** For all methods, we apply full gradient checkpointing [5] to both attention  
177 and feedforward, following prior works [27, 22]. The experiments are on both GPUs and TPUs.  
178 For GPUs, we consider both single DGX A100 server with 8 GPUs and distributed 32 A100 GPUs.  
179 We also experiment with TPUs, from older generations TPUv3 to newer generations of TPUv4 and  
180 TPUv5e. We note that all of our results are obtained using full precision instead of mixed precision.

## 181 5 Results

182 In our experiments, our primary objective is to comprehensively evaluate the performance of Ring  
183 Attention across multiple key metrics, including maximum supported sequence length within acceler-  
184 ator memory, model flops utilization, and throughput. We compare Ring Attention’s performance  
185 with several baseline models, including the vanilla transformers [35], transformers with memory  
186 efficient attention [27], and transformers with both memory efficient attention and feedforward [22],  
187 across different model sizes and accelerator configurations.

### 188 5.1 Evaluating Max Context Size

189 We evaluate maximum supported context length using tensor parallelism and batch size 1 in sequences.  
190 Following prior works [22, 31], we note that no data parallelism is considered in our evaluations  
191 since our approach is independent of data parallelism. As a result, the batch sizes used in our analysis  
192 are much lower than the ones used for the end-to-end training. Practitioners can combine our method  
193 with data parallelism to scale up batch size, which we will show in Section 5.2. Table 3 summarizes  
194 the results of our experiments.

195 Our Ring Attention model consistently surpasses baselines, delivering superior scalability across  
196 diverse hardware setups. For example, with 32 A100 GPUs, we achieve over 32 million tokens in  
197 context size, a significant improvement over baselines. Furthermore, when utilizing larger accelerators  
198 like TPUv4-512, Ring Attention enables a 512x increase in context size, allows training sequences of  
199 over 100 million tokens. Furthermore, our Ring Attention model scales linearly with the number of  
200 devices, as demonstrated by the 8x improvement over BPT on 8 A100 and the 512x improvement on  
201 TPUv4-512. If a model can be trained with context size  $s$  on  $n$  GPUs using the blockwise attention  
202 and feedforward, with our Ring Attention approach, it becomes possible to train a model with a  
203 context size of  $ns$ .

---

<sup>1</sup>Unlike TPUv4-256 and TPUv5-256 where the number 256 represents the count of TPUv4 (v5) hosts, TPUv3 uses a doubled host count notation. So, TPUv3-512 means there are 256 hosts. See [https://cloud.google.com/tpu/docs/system-architecture-tpu-vm#tpu\\_v3](https://cloud.google.com/tpu/docs/system-architecture-tpu-vm#tpu_v3) for more details.

Table 3: Maximum context length supported in device memory on different model sizes and clusters of accelerators. Baselines are vanilla transformer [35], transformer with memory efficient attention [27], and transformer with memory efficient attention and feedforward [22]. The context size is reported in tokens (1e3). Our Ring Attention substantially outperforms baselines and scales linearly with number of devices, achieving over 100M context size.

	Max context size supported ( $\times 1e3$ )				Ours vs SOTA
	Vanilla	Memory Efficient Attn	Memory Efficient Attn and FFN	Ring Attention (Ours)	
8x A100 NVLink					
3B	16	256	512	4096 (4M)	8x
7B	16	256	512	4096 (4M)	8x
13B	8	128	256	2048 (2M)	8x
30B	8	64	256	2048 (2M)	8x
32x A100 InfiniBand					
7B	32	512	1024	32768 (32M)	32x
30B	16	128	512	16384 (16M)	32x
TPUv3-512 <sup>1</sup>					
7B	4	16	64	16384 (16M)	256x
13B	2	8	32	8192 (8M)	256x
30B	1	4	16	4096 (4M)	256x
TPUv4-512					
3B	8	64	256	131072 (131M)	512x
7B	8	32	128	65536 (65M)	512x
13B	4	16	64	32768 (32M)	512x
30B	2	8	32	16384 (16M)	512x
TPUv5e-256					
7B	4	16	64	16384 (16M)	256x
30B	1	4	16	4096 (4M)	256x

## 204 5.2 Evaluating Model Flops Utilization

205 We evaluate the model flops utilization (MFU) of Ring Attention in standard training settings  
 206 using fully sharded data parallelism(FSDP) [10] and tensor parallelism following LLaMA and  
 207 OpenLLaMA [34, 11]. The batch size in tokens are 2M on 8/32x A100 and 4M on TPUv4-256. Our  
 208 goal is investigating the impact of model size and context length on MFU, a critical performance  
 209 metrics while highlighting the benefits of our approach. Table 5.1 presents the results of our  
 210 experiments on MFU for different model sizes and context lengths. We present the achieved MFU  
 211 using state-of-the-art memory efficient transformers BPT [22], compare it to our anticipated MFU  
 212 based on these results, and demonstrate the actual MFU obtained with our approach (Ring Attention).  
 213 For fair comparison, both BPT and our approach are based on the same BPT implementation<sup>2</sup> on  
 214 both GPUs and TPUs. It’s worth noting that on GPUs our approach Ring Attention can be also  
 215 integrated with the more compute efficient Triton code [16] or CUDA code [9] of memory efficient  
 216 attention [27], similarly on TPUs it is also compatible with Pallas [33]. Combing these low level  
 217 kernels implementations with our approach can maximize MFU, we leave that to future work.

218 Ring Attention trains much longer context sizes for self-attention, resulting in higher self-attention  
 219 FLOPs compared to baseline models. Since self-attention has a lower MFU than feedforward, Ring  
 220 Attention is expected to have a lower MFU than the baseline models. Our method offers a clear  
 221 advantage in terms of maintaining MFU while enabling training with significantly longer context  
 222 lengths. As shown in Table 5.1, when comparing our approach to prior state-of-the-arts, it is evident  
 223 that we can train very large context models without compromising MFU or throughput.

<sup>2</sup>[https://github.com/lhao499/llm\\_large\\_context](https://github.com/lhao499/llm_large_context)

Table 4: Model flops utilization (MFU) with different training configurations: model sizes, compute, and context lengths. Ring Attention enables training large models (30B-65B) for over 1M context size with negligible overheads.

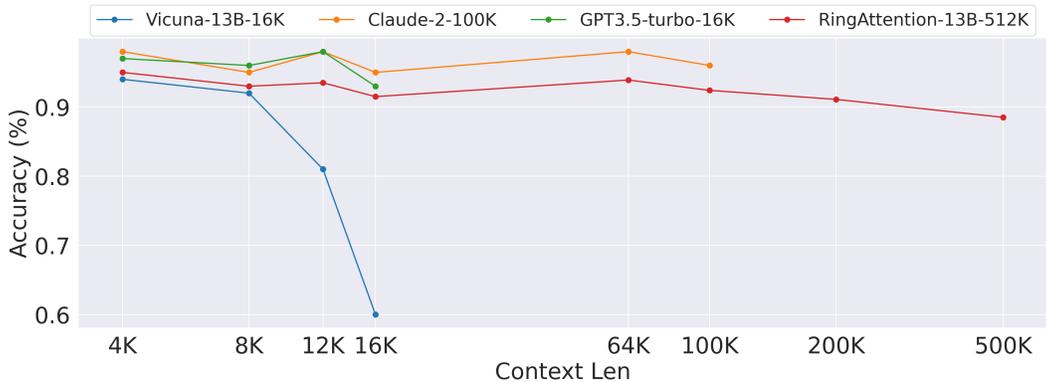
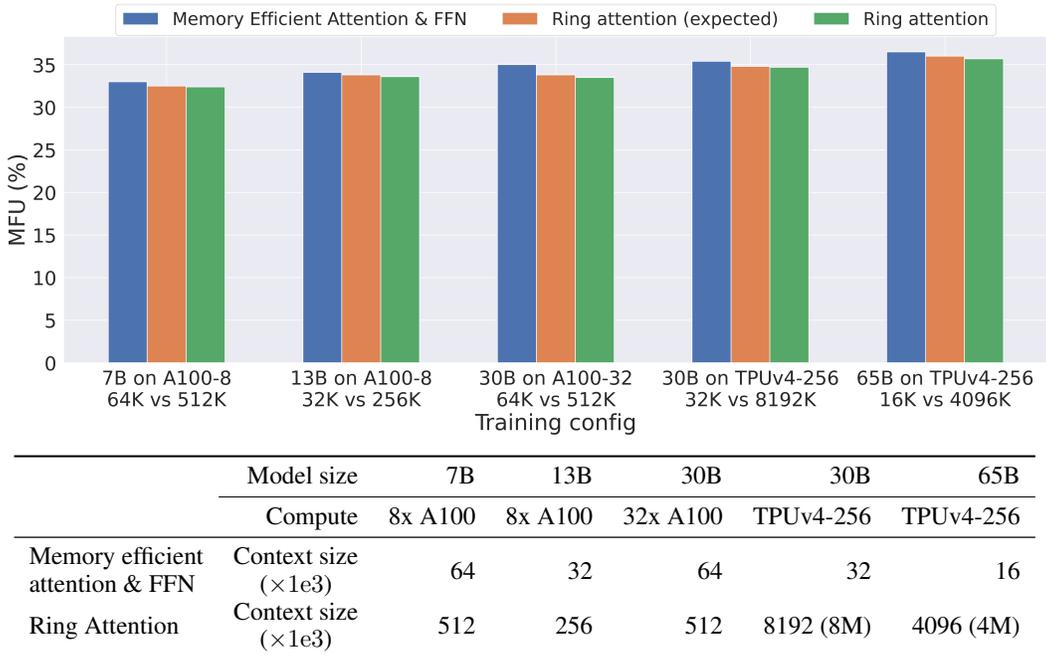


Figure 3: Comparison of different models on the long-range line retrieval task.

### 224 5.3 Impact on LLM Performance

225 We evaluate Ring Attention by applying our method to finetune LLaMA model to longer context. In  
 226 this experiment, while our approach enables training with millions of context tokens, we conducted  
 227 finetuning on the LLaMA-13B model, limiting the context length to 512K tokens due to constraints  
 228 on our cloud compute budget. This finetuning was carried out on 8 A100 GPUs, using the ShareGPT  
 229 dataset, following methodologies as outlined in prior works [6, 12]. We then evaluated our finetuned  
 230 model on the line retrieval test [19]. In this test, the model needs to precisely retrieve a number  
 231 from a long document, the task can effectively capture the abilities of text generation, retrieval, and  
 232 information association at long context, reflected by the retrieving accuracy. Figure 3 presents the  
 233 accuracy results for different models across varying context lengths (measured in tokens). Notably,  
 234 our model, Ring Attention-13B-512K, stands out as it maintains high accuracy levels even with  
 235 long contexts. GPT3.5-turbo-16K, Vicuna-16B-16K, and Claude-2-100K demonstrate competitive  
 236 accuracy within short context lengths. However, they cannot handle extended context lengths.

## 237 6 Related Work

238 Transformers have garnered significant attention in the field of AI and have become the backbone  
239 for numerous state-of-the-art models. Several works have explored memory-efficient techniques  
240 to address the memory limitations of Transformers and enable their application to a wider range  
241 of problems. Computing exact self-attention in a blockwise manner using the tiling technique [23]  
242 has led to the development of memory efficient attention mechanisms [27] and its efficient CUDA  
243 implementation [9], and blockwise parallel transformer [22] that proposes computing both feedfor-  
244 ward and self-attention block-by-block, resulting in a significant reduction in memory requirements.  
245 In line with these advancements, our work falls into the category of memory efficient computation  
246 for Transformers. Other works have investigated the approximation of attention mechanisms, yet  
247 these efforts have often yielded sub-optimal results or encountered challenges during scaling up.  
248 For an in-depth review of these techniques, we recommend referring to the surveys by Narang et al.  
249 [25], Tay et al. [32]. Another avenue of research explores various parallelism methods, including  
250 tensor parallelism [31], pipeline parallelism [14], sequence parallelism [20, 17], and FSDP [10, 28].  
251 The activations of self-attention take a substantial amount of memory for large context models and  
252 tensor parallelism can only reduce parts of activations memory. Sequence parallelism of self-attention  
253 introduces a significant communication overhead that cannot be overlapped with computation, our  
254 work leverages on blockwise parallel transformers to distribute blockwise computation across devices  
255 and concurrently overlaps the communication of key-value blocks between hosts with blockwise  
256 computation. Overlapping communication with computation has been studied in high performance  
257 computing literature [7, 36, 8, *inter alia*]. While ring communication has found applications in other  
258 parallel computing scenarios [2, 15, 13, 30], our work stands out as the first work to show that it can  
259 be applied to self-attention as used in Transformers and to make it fit efficiently into Transformer  
260 training and inference without adding significant overhead by overlapping blockwise computation  
261 and communication.

## 262 7 Conclusion

263 In conclusion, we propose a memory efficient approach to reduce the memory requirements of  
264 Transformers, the backbone of state-of-the-art AI models. Our approach enables the context length to  
265 scale linearly with the number of devices while maintaining performance, eliminating the memory  
266 bottleneck imposed by individual devices. Through extensive experiments, we demonstrate its  
267 effectiveness, achieving up to 512x memory reduction than prior memory efficient Transformers. Our  
268 contributions include a practical method for large context sizes in large Transformer models.

269 **Limitations and Future Work.** Although our method achieves state-of-the-art low memory usage  
270 for Transformer models, it does have some limitations that need to be addressed:

- 271 • *Scaled up training:* due to compute budget constraint, our experiments focus on evaluation the  
272 effectiveness of the proposed approach without large scale training models.
- 273 • *Optimal compute efficiency:* While Ring Attention enables the context length to scale linearly  
274 with the number of devices while maintaining performance, optimizing low-level operations is  
275 crucial for achieving optimal compute efficiency. In future works, we suggest considering porting  
276 our method to CUDA / OpenAI Triton /Jax Pallas to achieve maximum sequence length and  
277 compute performance.

278 In terms of future prospects, the possibility of near-infinite context introduces a vast array of exciting  
279 opportunities, such as large video-language models, decision making and tool use transformers on ex-  
280 tended trial-and-error experience, understanding and generating large code projects, and transforming  
281 language models into a versatile AI scientist for helping understand science experimental data.

## 282 References

- 283 [1] Anthropic. Introducing claude, 2023. URL [https://www.anthropic.com/index/  
284 introducing-claude](https://www.anthropic.com/index/introducing-claude).
- 285 [2] Christian Bischof. *Parallel computing: Architectures, algorithms, and applications*, volume 15.  
286 IOS Press, 2008.

- 287 [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal,  
288 Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are  
289 few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- 290 [4] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Misha Laskin, Pieter  
291 Abbeel, Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning  
292 via sequence modeling. *Advances in neural information processing systems*, 34:15084–15097,  
293 2021.
- 294 [5] Tianqi Chen, Bing Xu, Chiyuan Zhang, and Carlos Guestrin. Training deep nets with sublinear  
295 memory cost. *arXiv preprint arXiv:1604.06174*, 2016.
- 296 [6] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng,  
297 Siyuan Zhuang, Yonghao Zhuang, Joseph E Gonzalez, et al. Vicuna: An open-source chatbot  
298 impressing gpt-4 with 90%\* chatgpt quality. See <https://vicuna.lmsys.org>, 2023.
- 299 [7] Anthony Danalis, Ki-Yong Kim, Lori Pollock, and Martin Swany. Transformations to parallel  
300 codes for communication-computation overlap. In *SC’05: Proceedings of the 2005 ACM/IEEE  
301 conference on Supercomputing*, pages 58–58. IEEE, 2005.
- 302 [8] Anthony Danalis, Lori Pollock, Martin Swany, and John Cavazos. Mpi-aware compiler opt-  
303 imizations for improving communication-computation overlap. In *Proceedings of the 23rd  
304 international conference on Supercomputing*, pages 316–325, 2009.
- 305 [9] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and  
306 memory-efficient exact attention with io-awareness. *Advances in Neural Information Processing  
307 Systems*, 35:16344–16359, 2022.
- 308 [10] Facebook. Fully Sharded Data Parallel: faster AI training with fewer GPUs — engineer-  
309 ing.fb.com. <https://engineering.fb.com/2021/07/15/open-source/fSDP/>, 2023.
- 310 [11] Xinyang Geng and Hao Liu. Openllama: An open reproduction of llama, may 2023. URL  
311 [https://github.com/openlm-research/open\\_llama](https://github.com/openlm-research/open_llama), 2023.
- 312 [12] Xinyang Geng, Arnav Gudibande, Hao Liu, Eric Wallace, Pieter Abbeel, Sergey Levine, and  
313 Dawn Song. Koala: A dialogue model for academic research. *Blog post, April, 1, 2023*.
- 314 [13] Andrew Gibiansky. Bringing hpc techniques to deep learning. *Baidu Research, Tech. Rep.*,  
315 2017.
- 316 [14] Yanping Huang, Youlong Cheng, Ankur Bapna, Orhan Firat, Dehao Chen, Mia Chen, Hy-  
317 oukJoong Lee, Jiquan Ngiam, Quoc V Le, Yonghui Wu, et al. Gpipe: Efficient training of giant  
318 neural networks using pipeline parallelism. *Advances in neural information processing systems*,  
319 32, 2019.
- 320 [15] Joshua Hursey and Richard L Graham. Building a fault tolerant mpi application: A ring  
321 communication example. In *2011 IEEE International Symposium on Parallel and Distributed  
322 Processing Workshops and Phd Forum*, pages 1549–1556. IEEE, 2011.
- 323 [16] OpenAI kernel team. Openai triton fused attention, 2023. URL [https://github.com/  
324 openai/triton/blob/main/python/tutorials/06-fused-attention.py](https://github.com/openai/triton/blob/main/python/tutorials/06-fused-attention.py).
- 325 [17] Vijay Korthikanti, Jared Casper, Sangkug Lym, Lawrence McAfee, Michael Andersch, Moham-  
326 mad Shoeybi, and Bryan Catanzaro. Reducing activation recomputation in large transformer  
327 models. *arXiv preprint arXiv:2205.05198*, 2022.
- 328 [18] Michael Laskin, Denis Yarats, Hao Liu, Kimin Lee, Albert Zhan, Kevin Lu, Catherine Cang,  
329 Lerrel Pinto, and Pieter Abbeel. Urlb: Unsupervised reinforcement learning benchmark. *arXiv  
330 preprint arXiv:2110.15191*, 2021.
- 331 [19] Dacheng Li, Rulin Shao, Anze Xie, Ying Sheng, Lianmin Zheng, Joseph E. Gonzalez, Ion  
332 Stoica, Xuezhe Ma, and Hao Zhang. How long can open-source llms truly promise on context  
333 length?, June 2023. URL <https://lmsys.org/blog/2023-06-29-longchat>.

- 334 [20] Shenggui Li, Fuzhao Xue, Yongbin Li, and Yang You. Sequence parallelism: Making 4d  
335 parallelism possible. *arXiv preprint arXiv:2105.13120*, 2021.
- 336 [21] Hao Liu and Pieter Abbeel. Emergent agentic transformer from chain of hindsight experience.  
337 *International Conference on Machine Learning*, 2023.
- 338 [22] Hao Liu and Pieter Abbeel. Blockwise parallel transformer for large context models. *Advances*  
339 *in neural information processing systems*, 2023.
- 340 [23] Maxim Milakov and Natalia Gimelshein. Online normalizer calculation for softmax. *arXiv*  
341 *preprint arXiv:1805.02867*, 2018.
- 342 [24] MosaicML. Introducing mpt-7b: A new standard for open-source, commercially usable llms,  
343 2023. URL <https://www.mosaicml.com/blog/mpt-7b>.
- 344 [25] Sharan Narang, Hyung Won Chung, Yi Tay, William Fedus, Thibault Fevry, Michael Matena,  
345 Karishma Malkan, Noah Fiedel, Noam Shazeer, Zhenzhong Lan, et al. Do transformer modifi-  
346 cations transfer across implementations and applications? *arXiv preprint arXiv:2102.11972*,  
347 2021.
- 348 [26] OpenAI. Gpt-4 technical report, 2023.
- 349 [27] Markus N Rabe and Charles Staats. Self-attention does not need  $o(n^2)$  memory. *arXiv preprint*  
350 *arXiv:2112.05682*, 2021.
- 351 [28] Samyam Rajbhandari, Jeff Rasley, Olatunji Ruwase, and Yuxiong He. Zero: Memory optimiza-  
352 tions toward training trillion parameter models. In *SC20: International Conference for High*  
353 *Performance Computing, Networking, Storage and Analysis*, pages 1–16. IEEE, 2020.
- 354 [29] J. Schulman, B. Zoph, C. Kim, J. Hilton, J. Menick, J. Weng, J. F. C. Uribe, L. Fedus, L. Metz,  
355 M. Pokorny, R. G. Lopes, S. Zhao, A. Vijayvergiya, E. Sigler, A. Perelman, C. Voss, M. Heaton,  
356 J. Parish, D. Cummings, R. Nayak, V. Balcom, D. Schnurr, T. Kaftan, C. Hallacy, N. Turley,  
357 N. Deutsch, and V. Goel. Chatgpt: Optimizing language models for dialogue. *OpenAI Blog*,  
358 2022. URL <https://openai.com/blog/chatgpt>.
- 359 [30] Alexander Sergeev and Mike Del Balso. Horovod: fast and easy distributed deep learning in  
360 tensorflow. *arXiv preprint arXiv:1802.05799*, 2018.
- 361 [31] Mohammad Shoeybi, Mostofa Patwary, Raul Puri, Patrick LeGresley, Jared Casper, and Bryan  
362 Catanzaro. Megatron-lm: Training multi-billion parameter language models using model  
363 parallelism. *arXiv preprint arXiv:1909.08053*, 2019.
- 364 [32] Yi Tay, Mostafa Dehghani, Samira Abnar, Hyung Won Chung, William Fedus, Jinfeng Rao,  
365 Sharan Narang, Vinh Q Tran, Dani Yogatama, and Donald Metzler. Scaling laws vs model  
366 architectures: How does inductive bias influence scaling? *arXiv preprint arXiv:2207.10551*,  
367 2022.
- 368 [33] Jax team. Jax pallas fused attention, 2023. URL [https://github.com/google/jax/blob/main/jax/experimental/pallas/ops/tpu/flash\\_attention.py](https://github.com/google/jax/blob/main/jax/experimental/pallas/ops/tpu/flash_attention.py).
- 370 [34] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timo-  
371 thée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open  
372 and efficient foundation language models. *arXiv preprint arXiv:2302.13971*, 2023.
- 373 [35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,  
374 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*  
375 *processing systems*, 30, 2017.
- 376 [36] Shibo Wang, Jinliang Wei, Amit Sabne, Andy Davis, Berkin Ilbeyi, Blake Hechtman, Dehao  
377 Chen, Karthik Srinivasa Murthy, Marcello Maggioni, Qiao Zhang, et al. Overlap communication  
378 with dependent computation via decomposition in large deep learning models. In *Proceedings of*  
379 *the 28th ACM International Conference on Architectural Support for Programming Languages*  
380 *and Operating Systems, Volume 1*, pages 93–106, 2022.
- 381 [37] Denis Yarats, David Brandfonbrener, Hao Liu, Michael Laskin, Pieter Abbeel, Alessandro  
382 Lazaric, and Lerrel Pinto. Don’t change the algorithm, change the data: Exploratory data for  
383 offline reinforcement learning. *arXiv preprint arXiv:2201.13425*, 2022.

## 384 A Code

385 The implementation of Ring Attention in Jax is provided in Figure 4. We use `defvjp` function  
386 to define both the forward and backward passes, and use collective operation `jax.lax.ppermute`  
387 to facilitate the exchange of key-value blocks among a ring of hosts. The provided code snippet  
388 highlights essential components of Ring Attention. The complete implementation with  
389 maximum memory efficient just needs to replace the local blockwise computation, specifically  
390 `jnp.einsum("bshd,btd->bhst", q, k)` and `jnp.einsum("bhst,btd->bshd", s, v)`  
391 as well as the local blockwise feedforward computation with BPT’s Jax based blockwise attention  
392 and FFN computation. For maximum compute efficiency our Ring Attention can be integrated  
393 with exiting kernel-level fused-attention implementations, such as on GPUs Ring Attention can be  
394 integrated with Triton code [16] or CUDA code [9], similarly on TPUs it is also compatible with  
395 Pallas code [33] of the memory efficient attention [27].

## 396 B Experiment Details

### 397 B.1 Evaluation of context length

398 In the experimental results presented in Section 5.1, we used tensor parallelism to partition the model  
399 across GPUs or TPU units. Our evaluation focused on determining the maximum achievable sequence  
400 length, using a sequence number of one. For TPUs, we utilized its default training configuration,  
401 which involved performing matmul operations in `bf16` format with weight accumulation in  
402 `float32`. On the other hand, for GPUs, we adopted the default setup, where all operations were  
403 performed in `float32`.

### 404 B.2 Evaluation of MFU

405 In the evaluation presented in Section 5.2, the training was conducted using FSDP [10] with no  
406 gradient accumulation. The batch size in tokens is 2 million per batch on GPU and 4 million per batch  
407 on TPU. For gradient checkpointing [5], we used `nothing_saveable` as checkpointing policies for  
408 attention and feedforward network (FFN). For more details, please refer to Jax documentation.

### 409 B.3 Evaluation on line retrieval

410 In the evaluation presented in Section 5.3, the training was conducted using FSDP on 8x A100 80GB  
411 Cloud GPUs. We finetuned the LLaMA-13B model [34], limiting context length to 512K tokens due  
412 to constraints on our cloud compute budget, though our approach enables training with millions of  
413 context tokens. We use user-shared conversations gathered from ShareGPT.com with its public APIs  
414 for finetuning, following methodologies as outlined in prior works [6, 12]. ShareGPT is a website  
415 where users can share their ChatGPT conversations. To ensure data quality, we convert the HTML  
416 back to markdown and filter out some inappropriate or low-quality samples, which results in 125K  
417 conversations after data cleaning.

```

1  @partial(jax.custom_vjp, nondiff_argnums=[3, 4, 5])
2  def _ring_attention_fwd(q, k, v, mask, axis_name, float32_logits):
3      if float32_logits:
4          q, k = q.astype(jnp.float32), k.astype(jnp.float32)
5          batch, q_len, num_heads, _ = q.shape
6          batch, kv_len, dim_per_head = k.shape
7          numerator = jnp.zeros((batch, q_len, num_heads, dim_per_head)).astype(q.dtype)
8          denominator = jnp.zeros((batch, num_heads, q_len)).astype(q.dtype)
9          axis_size = lax.psum(1, axis_name)
10         scale = jnp.sqrt(q.shape[-1])
11         def scan_kv_block(carry, idx):
12             prev_max_score, numerator, denominator, k, v = carry
13             mask = lax.dynamic_slice_in_dim(mask,
14                 (lax.axis_index(axis_name) - idx) % axis_size * kv_len, kv_len, axis=-1)
15             attn_weights = jnp.einsum("bqhd,bkd->bhqk", q, k) / scale
16             attn_weights = jnp.where(mask, -jnp.inf, attn_weights)
17             max_score = jnp.maximum(prev_max_score, jnp.max(attn_weights, axis=-1))
18             exp_weights = jnp.exp(attn_weights - max_score[... , None])
19             correction = rearrange(jnp.exp(prev_max_score - max_score), 'b h q -> b q h')[... , None]
20             numerator = numerator * correction + jnp.einsum("bhqk,bkd->bqhd", exp_weights, v)
21             denominator = denominator * jnp.exp(prev_max_score - max_score) + jnp.sum(exp_weights, axis=-1)
22             k, v = map(lambda x: lax.ppermute(x, axis_name, perm=[[i,
23                 (i + 1) % axis_size] for i in range(axis_size)]), (k, v))
24             return (max_score, numerator, denominator, k, v), None
25         prev_max_score = jnp.full((batch, num_heads, q_len), -jnp.inf).astype(q.dtype)
26         (numerator, max_score, denominator, _, _) = lax.scan(scan_kv_block,
27             init=(prev_max_score, numerator, denominator, k, v), xs=jnp.arange(0, axis_size))
28         output = numerator / rearrange(denominator, 'b h q -> b q h')[... , None]
29         return output.astype(v.dtype), (output, q, k, v, numerator, denominator, max_score)
30
31 def _ring_attention_bwd(mask, axis_name, float32_logits, res, g):
32     del float32_logits
33     axis_size = lax.psum(1, axis_name)
34     output, q, k, v, numerator, denominator, max_score = res
35     dq = jnp.zeros_like(q, dtype=jnp.float32)
36     dk = jnp.zeros_like(k, dtype=jnp.float32)
37     dv = jnp.zeros_like(v, dtype=jnp.float32)
38     batch, kv_len, dim_per_head = k.shape
39     scale = jnp.sqrt(q.shape[-1])
40     def scan_kv_block(carry, idx):
41         dq, dk, dv, k, v = carry
42         mask = lax.dynamic_slice_in_dim(mask,
43             (lax.axis_index(axis_name) - idx) % axis_size * kv_len, kv_len, axis=-1)
44         attn_weights = jnp.einsum("bqhd,bkd->bhqk", q, k) / scale
45         attn_weights = jnp.where(mask, -jnp.inf, attn_weights)
46         exp_weights = jnp.exp(attn_weights - max_score[... , None]) / denominator[... , None]
47         ds = jnp.einsum("bqhd,bkd->bhqk", g, v)
48         dl = (ds - jnp.einsum("bqhd,bqhd->bhs", g, output))[... , None] * exp_weights
49         dq = dq + jnp.einsum("bhqk,bkd->bqhd", dl, k) / scale
50         dk = dk + jnp.einsum("bqhd,bhqk->bkd", q, dl) / scale
51         dv = dv + jnp.einsum("bhqk,bqhd->bkd", exp_weights, g)
52         k, v, dk, dv = map(lambda x: lax.ppermute(x, axis_name, perm=[[i,
53             (i + 1) % axis_size] for i in range(axis_size)]), (k, v, dk, dv))
54         return (dq, dk, dv, k, v), None
55     (dq, dk, dv, k, v) = lax.scan(scan_kv_block, init=(dq, dk, dv, k, v), xs=jnp.arange(0, axis_size))
56     dq, dk, dv = dq.astype(q.dtype), dk.astype(k.dtype), dv.astype(v.dtype)
57     return dq, dk, dv
58
59 @partial(jax.custom_vjp, nondiff_argnums=[3, 4, 5])
60 def ring_attention(q, k, v, mask, axis_name, float32_logits=True):
61     y, _ = _ring_attention_fwd(q, k, v, mask, axis_name, float32_logits)
62     return y
63
64 ring_attention.defvjp(_ring_attention_fwd, _ring_attention_bwd)

```

Figure 4: Key parts of the implementation of Ring Attention in Jax. We use collective operation `lax.ppermute` to send and receive key value blocks between previous and next hosts.

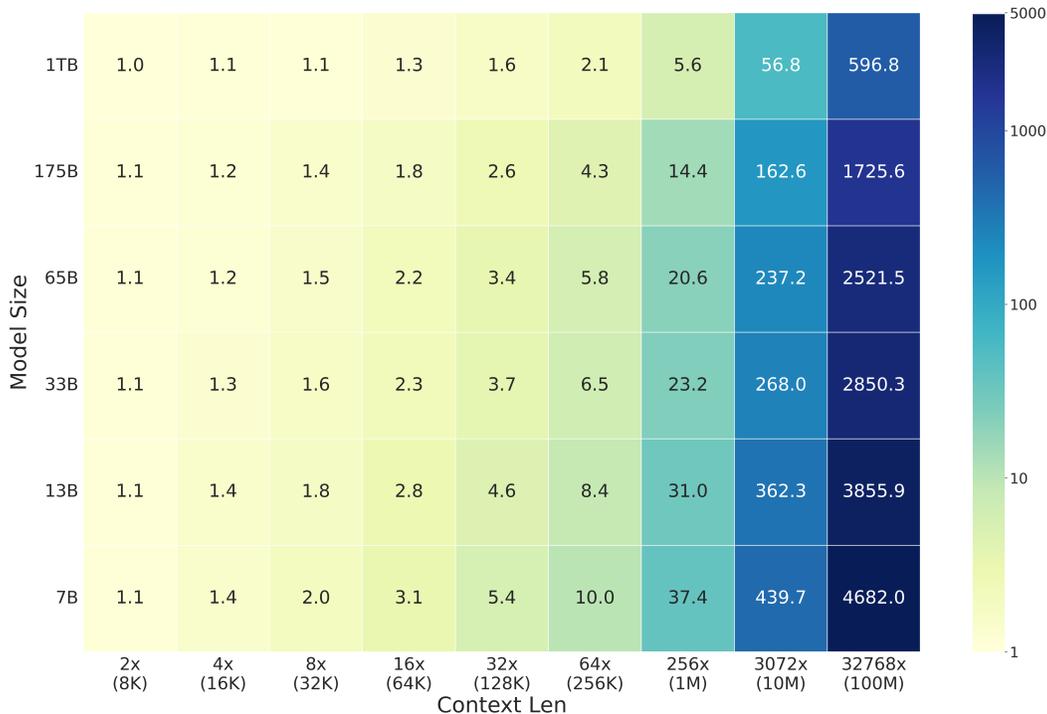


Figure 5: The per dataset training FLOPs cost ratio relative to a 4k context size, considering different model dimensions. On the x-axis, you’ll find the context length, where, for example, 32x(128k) denotes a context length of 128k, 32x the size of the same model’s 4k context length.

## 418 C Training FLOPs Scaling of Context Size

419 Given that our proposed approach unlocks the possibility of training with a context size exceeding 100  
 420 million tokens and allows for linear scaling of the context size based on the number of devices, it is  
 421 essential to understand how the training FLOPs per dataset scale with the context size. While a larger  
 422 context size results in a higher number of FLOPs, the increased ratio does not scale quadratically  
 423 because the number of tokens remains fixed. We present these results in Figure 5, which showcases  
 424 various model sizes and context lengths, representing different computational budgets. The figure  
 425 illustrates the ratio of FLOPs for larger context lengths compared to the same model with a shorter  
 426 4K context size. We calculated the per sequence FLOPs using  $(24bs_2h^2 + 4bs_2^2h)n$  where  $h$  is  
 427 model hidden dimension,  $b$  is batch size,  $s$  is total sequence length, and  $n$  is number of layers. The  
 428 per dataset FLOPs ratio is then given by  $((24bs_2h^2 + 4bs_2^2h)/(24bs_1h^2 + 4bs_1^2h))/(s_2/s_1) =$   
 429  $(6h + s_2)/(6h + s_1)$ , where  $s_2$  and  $s_1$  are new and old context lengths. Model sizes and their  
 430 hidden dimensions are as follows: LLaMA-7B (4096), LLaMA-13B (5140), LLaMA-33B (7168),  
 431 LLaMA-65B (8192), GPT3-175B (12288), and 1TB (36864). These model configurations are from  
 432 LLaMA [34] and GPT-3 [3] papers, except the 1TB model size and dimension were defined by us.

433 As depicted in Figure 5, scaling up small models to a 1M context size results in approximately 20-40  
 434 times more FLOPs, and even more for 10M and 100M token context sizes. However, as the model  
 435 sizes increase, the cost ratio decreases. For instance, scaling up the 170B model from 4K to 10M  
 436 incurs 162.6x higher per dataset FLOPs, despite the context size being 3072 times longer.

## 437 D Impact on In Context RL Performance

438 In addition to show the application of Ring Attention to finetune LLM in Section 5.3, we present  
 439 additional results of applying Ring Attention for learning trial-and-error RL experience using Trans-  
 440 formers. We report our results in Table 5, where we evaluate our proposed model on the ExoRL  
 441 benchmark across six different tasks. On ExoRL, we report the cumulative return, as per ExoRL [37].  
 442 We compare BC, DT [4], AT [21], and AT with memory efficient attention [27] (AT+ME), AT with  
 443 blockwise parallel transformers [22] (AT+BPT), and AT with our Ring Attention (AT+Ring Attention).

Table 5: Application of Ring Attention on improving Transformer in RL. BC and DT use vanilla attention. AT + ME denotes using memory efficient attention, AT + BPT denotes using blockwise parallel transformer. AT + RA denotes using Ring Attention.

ExoRL	BC-10%	DT	AT + ME	AT + BPT	AT + BPT	AT + RA
Task			N Trajs = 32	N Trajs = 32	N Trajs = 128	N Trajs = 128
Walker Stand	52.91	34.54	oom	95.45	oom	98.23
Walker Run	34.81	49.82	oom	105.88	oom	110.45
Walker Walk	13.53	34.94	oom	78.56	oom	78.95
Cheetah Run	34.66	67.53	oom	178.75	oom	181.34
Jaco Reach	23.95	18.64	oom	87.56	oom	89.51
Cartpole Swingup	56.82	67.56	oom	120.56	oom	123.45
<b>Total Average</b>	36.11	45.51	oom	111.13	oom	113.66

444 The numbers of BC, DT, AT are from the ExoRL and AT paper. AT + Ring Attention numbers are  
445 run by ourselves. Since the ExoRL data is highly diverse, having been collected using unsupervised  
446 RL [18], it has been found that TD learning performs best, while behavior cloning struggles [37].  
447 AT [21] shows that conditioning Transformer on multiple trajectories with relabeled target return can  
448 achieve competitive results with TD learning. For more details, please refer to their papers. We are  
449 interested in applying Ring Attention to improve the performance of AT by conditioning on a larger  
450 number of trajectories rather than 32 trajectories in prior works. It is worth noting that each trajectory  
451 has  $1000 \times 4$  length where 1000 is sequence length while 4 is return-state-action-reward, making  
452 training 128 trajectories with modest 350M size model infeasible for prior state-of-the-art blockwise  
453 parallel transformers. Results in Table 5 show that, by scaling up the sequence length (number of  
454 trajectories), AT + Ring Attention consistently outperforms original AT with BPT across all six tasks,  
455 achieving a total average return of 113.66 compared to the AT with BPT model’s total average return  
456 of 111.13. The results show that the advantage of Ring Attention for training and inference with long  
457 sequences.