
SeCOv2: Highly Memory-Efficient Training for LLMs with Million-Token Contexts

Anonymous Author(s)

Affiliation
Address
email

Abstract

1 Training Large Language Models (LLMs) on long contexts is hindered by the pro-
2hibitive GPU memory cost of activations. We introduce SeCOv2, a block-recurrent
3 framework that overcomes this limitation by processing sequences in blocks, akin
4 to an RNN. By discarding and recomputing block activations on-the-fly, our method
5 maintains a constant activation memory footprint, irrespective of sequence length.
6 To address the subsequent bottleneck from the linearly growing KV cache (a bot-
7 tleneck ignored by v1), we introduce three co-designed optimizations: (1) a paged
8 memory manager for the KV cache and its gradients that eliminates redundant mem-
9 ory operations and fragmentation; (2) Op-level cache management independent
10 of the PyTorch autograd system; and (3) an asynchronous mechanism to offload
11 inactive KV caches to the CPU, reducing their GPU memory consumption by over
12 10x. The synergy of these techniques yields a remarkably low memory overhead,
13 adding only 70-160 MB per 10,240 additional context tokens for 7-8B models. This
14 efficiency allows us to train a Qwen2.5-7B model with a 4M token context on a
15 single 96GB GPU, a substantial advance over existing open-source frameworks. To
16 further enhance time efficiency, we also combined SeCOv2 with tensor parallelism
17 and page-level sparse attention, accelerating the training process.

18

1 Introduction

19 Long-context modeling remains a formidable challenge in the training of Large Language Models
20 (LLMs), a problem rooted not in training time, but in the prohibitive GPU memory overhead
21 associated with long sequences [28, 38]. Research indicates that significant long-context capabilities
22 can be achieved with a surprisingly small number of training steps (as few as one thousand) [30, 10]
23 and that reaching a commercial standard may only require fine-tuning on approximately 5% of the
24 original pretraining token count [14, 4].

25 Despite this modest data requirement, the number of tokens processed in a single training iteration
26 grows dramatically, causing GPU memory consumption to scale linearly with context length and
27 rapidly exhaust available resources. For instance, a model with a 4× Grouped-Query Attention
28 (GQA) [2] ratio requires 64GB of HBM for the KV cache alone when processing a 256K context,
29 a demand that single-handedly overwhelms a 96GB H20 GPU. This figure does not even account
30 for the substantial memory footprint of other network activations, which often makes 128K tokens
31 a practical ceiling for many models [15, 40]. The immense memory pressure from the KV cache
32 and activations dwarfs the optimizer state, rendering popular memory-saving techniques like ZeRO3
33 offloading [33] and PyTorch FSDP [1] insufficient to overcome this fundamental bottleneck.

34 While various training-free methods for long-context extension have been proposed [20, 18, 3],
35 their practical utility is often overstated. As established in prior work [13], common evaluation
36 metrics such as Perplexity (PPL) and targeted retrieval tasks like Needle-in-a-Haystack (NIAH)

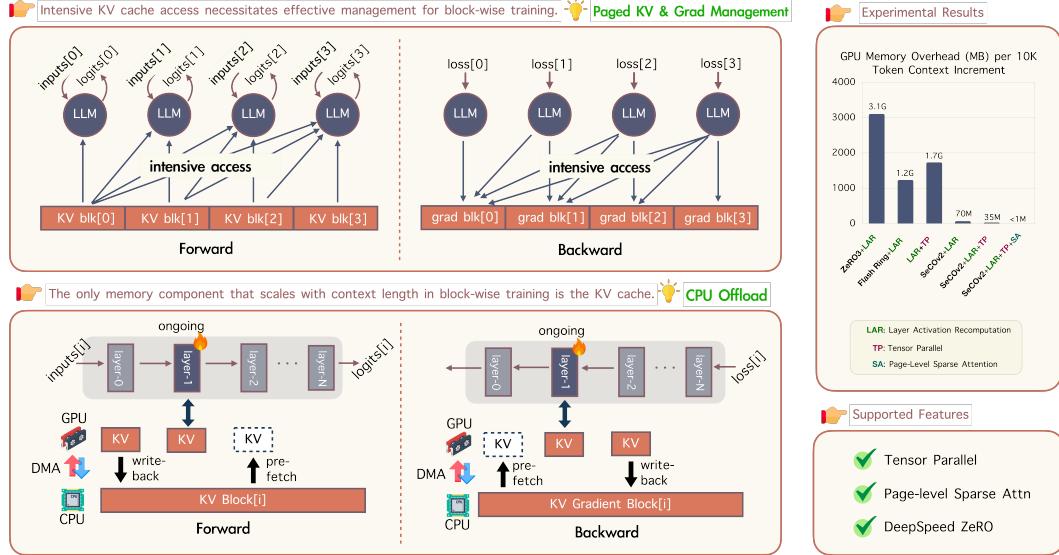


Figure 1: **Overview of the SeCoV2 Training Framework.** (*Top*) SeCoV2 processes long sequences in blocks, maintaining a constant activation memory footprint regardless of context length. (*Bottom*) This exposes the KV cache as the primary memory bottleneck in block-wise training. To solve this, we offload the cache to CPU memory and apply sparse attention, achieving near-constant GPU memory usage. (*Right*) The combined approach yields a $>1000\times$ reduction in memory consumption for every additional 10K context compared to strong baselines.

37 are unreliable indicators of true long-range reasoning abilities. Indeed, [13, 4] emphasizes that
 38 achieving commercial-grade performance requires dedicated fine-tuning with a meticulously curated
 39 data pipeline.

40 To this end, this paper directly confronts the central problem: how to train longer-context models
 41 with fewer resources, without compromising performance. We introduce a block-recurrent training
 42 framework that processes sequences in segments, analogous to an RNN. During the forward pass, each
 43 block’s activations are computed and then immediately discarded. When needed for the backward
 44 pass, they are recomputed on-the-fly. This strategy maintains a constant activation memory footprint
 45 equivalent to a single block, regardless of the total sequence length.

46 However, a critical distinction from standard RNNs complicates this picture: Transformer blocks are
 47 globally coupled via the KV cache. Processing any single block necessitates access to the KV caches
 48 from all preceding blocks. This global dependency introduces two primary bottlenecks: (1) **Memory**
 49 **Bottleneck**: The KV cache and its associated gradients must be retained throughout the training step,
 50 scaling linearly with the context length. (2) **Computational Overhead**: The forward and backward
 51 passes require frequent access to the entire KV cache and its gradient, introducing significant latency.

52 To overcome these challenges, we have engineered a highly integrated system centered on efficient
 53 KV cache management, synergistically combines several co-designed optimizations:

- 54 **1. Paged KV Cache and Gradient Management.** We introduce a paged memory manager
 55 (inspired by [21]) for both the KV cache and its gradients. This design eliminates
 56 expensive concatenation, reallocation and copy operations and mitigates memory
 57 fragmentation when appending new key-value pairs.
- 58 **2. Specialized Op and Kernels.** We execute all operations on the KV cache and its gradients
 59 directly within the Op, making them opaque to PyTorch’s autograd system. This avoids
 60 storing the KV cache as an activation and allows for direct gradient accumulation via CUDA
 61 `atomic_c_add`, improving both storage and computational efficiency.
- 62 **3. Asynchronous KV Cache CPU Offloading.** Since the KV cache and its gradients are the
 63 only components whose memory grows with context length, we designed an asynchronous
 64 mechanism to preemptively offload them to CPU memory. Experiments confirm that the
 65 resulting data transfer latency is effectively masked by the attention computation.

Table 1: **Comparison of Long-Context Training Methods.** SeCO series uniquely achieve $\mathcal{O}(1)$ activation memory complexity, where N denotes the context length. This key advantage allows SeCOv2, when combined with other techniques, to theoretically enable single-GPU training on contexts exceeding 10M tokens.

Method	GPUs	Act Memory	Exact Attn	Max Ctx Len.
ZeRO3 Offload [33]	8×H20	$\mathcal{O}(N)$	✓	64K
Ring Flash Attention [27]	8×H20	$\mathcal{O}(N/8)$	✓	256K
Tensor Parallel [19, 34]	8×H20	$\mathcal{O}(N/8)$	✓	256K
SeCO [23]	1×H20	$\mathcal{O}(1)$	✓	256K
SeCOv2	1×H20	$\mathcal{O}(1)$	✓	4M
SeCOv2 + Sparse Attn	1×H20	$\mathcal{O}(1)$	✗	10M+

66 The synergy of these techniques yields exceptional GPU memory efficiency. Our empirical results
 67 show that for every additional 10K tokens of context, the end-to-end training memory overhead
 68 increases by a mere 160MB for LLaMA3-8B [15] and 70MB for Qwen2.5-7B [40]. This stable and
 69 highly reproducible incremental cost is attributable solely to the paged KV cache and its gradients.

70 While context parallelism (CP) methods [27, 25, 34] are also highly memory-efficient, they do not
 71 extend the maximum training sequence length on a single GPU. Consequently, training a Qwen2.5-
 72 7B [40] model with a 4M context length requires a large cluster of 256 H100s [39]. In contrast, our
 73 approach enables training Qwen2.5-7B [40] with 4M-token context on a single 96GB GPU, repre-
 74 senting a substantial advance in resource efficiency. To showcase its flexibility, we also successfully
 75 integrated our framework with tensor parallel [34, 19] and page-level sparse attention [29, 42].

76 2 Related Work

77 **Efficient Long-Context LLM Training.** A core challenge in extending language models to long
 78 contexts is their failure to generalize to token distances and positional encodings not seen during
 79 pre-training [18]. While continued training can resolve these issues, the associated GPU memory
 80 costs are often prohibitive. This has motivated the development of training-free context extension
 81 methods, which typically modify positional encodings to accommodate longer sequences [7, 18, 20, 3].
 82 These approaches, however, introduce significant trade-offs. Methods based on interpolation or
 83 extrapolation of positional encodings [7] often degrade performance on short-text tasks by reducing
 84 the resolution of positional information [13]. Others that reuse position indices [18, 20, 3] achieve
 85 high scores on perplexity benchmarks but have been shown to fail on tasks requiring deep contextual
 86 understanding [13], indicating a superficial grasp of the extended context.

87 Given these limitations, state-of-the-art commercial and open-source models [15, 40, 26] still rely
 88 on fine-tuning to expand their context windows. However, the prevailing trend has shifted from
 89 algorithmic efficiency [8] towards overcoming memory barriers with massive computational re-
 90 sources [27, 22, 19, 34]. For instance, recent efforts [39] have required 256×H100 GPU cluster to
 91 extend LLaMA3.1 [15] from 128K to 4M context.

92 **Serial vs. Parallel Training Paradigms.** Language model training architectures can be broadly
 93 classified as parallel (Transformers [37], CNNs [31, 12], SSMs [11, 16]) or serial (RNNs [36]).
 94 Parallel training processes an entire sequence in a single forward and backward pass, maximizing
 95 GPU utilization and enabling scalability. However, this paradigm’s memory footprint, which scales
 96 linearly with sequence length, creates a significant bottleneck for long-context models. In contrast,
 97 serial training is highly memory-efficient, as it only activates a small part of the network at any given
 98 time, but at the cost of reduced parallelism.

99 For long-context fine-tuning where dataset sizes are often moderate, memory consumption, not raw
 100 throughput, is the primary constraint [30, 13]. This motivates a hybrid approach that processes
 101 sequences in blocks—parallel within each block and serial between them. This block-wise strategy
 102 is already standard in SOTA LLM inference engines (vLLM [21], FlashInfer [41]), where it incurs
 103 negligible latency. However, its application to training remains underdeveloped. An early exploration,

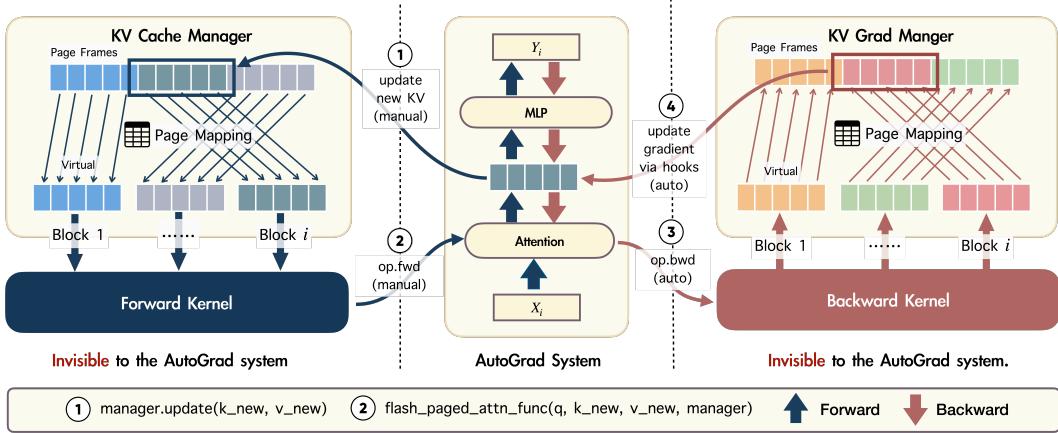


Figure 2: **An Overview of the Efficient KV Cache Management Framework.** Our framework optimizes frequent KV cache operations via two key strategies: (1) We bypass the autograd engine and manually manage KV cache values and their gradients. This ensures the portion of the computation graph handled by autograd has a constant memory footprint, resulting in $\mathcal{O}(1)$ activation memory. (2) We employ a paged memory scheme for the cache, which improves efficiency, reduces fragmentation, and natively supports page-level sparse attention [42, 29].

104 SeCO [23], introduced block-wise training but lacked critical optimizations at the operator and
 105 memory management levels. As a result, it could only train a 16K context model on a single RTX
 106 3090 GPU. Under identical conditions, our fully-optimized system extends this capability by more
 107 than 10 \times , demonstrating a dramatic improvement in efficiency.

108 **Context Parallelism.** Context Parallelism (CP) and our proposed method both segment long
 109 sequences but diverge significantly in their scalability and overhead profiles.

110 CP, inspired by Ring Attention [27] and implemented in frameworks like Megatron-LM [34], Deep-
 111 Speed Ulysses [19] and TorchTitan [24], distributes a single sequence across multiple GPUs. This
 112 design introduces substantial inter-GPU communication overhead that scales quadratically with GPU
 113 count. In contrast, our serial approach’s overhead is limited to repeated KV cache accesses and the
 114 latency from CPU offloading.

115 The trade-off is clear: CP excels with large GPU clusters but requires a high GPU count (e.g., at
 116 least 128 \times 80GB GPUs for a 4M context) and specialized, high-bandwidth networks to manage
 117 its communication demands. Our method, however, maintains moderate and predictable overhead,
 118 eliminating the need for large-scale GPU clusters and specialized interconnects.

119 3 Preliminary: SeCOv1

120 Activation recomputation is a well-established technique for reducing the memory footprint of training
 121 RNNs. In a standard RNN, the forward pass iteratively computes an output y_i and a new hidden state
 122 m_i at each timestep i based on the input x_i and the previous state m_{i-1} :

$$y_i, m_i, \mathbf{a}_i \leftarrow \mathcal{F}_{\text{RNN}}(x_i, m_{i-1}; \Theta) \quad (1)$$

123 where \mathbf{a}_i represents the intermediate activations required for the backward pass:

$$d\Theta, dm_{i-1} \leftarrow \text{backward}(dy_i, \mathbf{a}_i) \quad (2)$$

124 For long sequences, caching the activations \mathbf{a}_i for all timesteps results in a memory footprint that
 125 scales linearly with sequence length, quickly becoming a bottleneck. Activation recomputation [35, 5]
 126 circumvents this by discarding \mathbf{a}_i during the forward pass and only storing the much smaller hidden
 127 states m_i . The required activations are then regenerated on-the-fly just before their use in the
 128 backward pass. This trades a modest amount of recomputation for a substantial reduction in memory
 129 consumption.

130 **Hybrid Serial-Parallel Training of Transformers.** While Transformers are typically trained in a
 131 fully parallel fashion, the same principle of serial processing can be adapted to make their training

132 memory-efficient. We partition the input sequence into S blocks, $\{X_1, X_2, \dots, X_S\}$, and process
 133 them serially, akin to the timesteps of an RNN. The forward pass for block i is defined as:

$$Y_i, M_i, \mathbf{A}_i = \mathcal{F}_{\text{Transformer}}(X_i, \{M_1, M_2, \dots, M_{i-1}\}; \Theta) \quad (3)$$

134 Here, the model attends to the KV caches $\{M_1, \dots, M_{i-1}\}$ from preceding blocks, which serve a
 135 role analogous to the RNN’s hidden state. It produces the output Y_i , its own KV cache M_i , and the
 136 intermediate activations \mathbf{A}_i . The backward pass is likewise performed block-by-block:

$$d\Theta, \{dM_1, dM_2, \dots, dM_{i-1}\} \leftarrow \text{backward}(dY_i, dM_i, \mathbf{A}_i) \quad (4)$$

137 Following the recomputation strategy, the large activation tensor \mathbf{A}_i for each block is discarded
 138 immediately after the forward pass and recomputed just before its corresponding backward pass.

139 This block-wise recomputation fundamentally differs from the layer-wise activation checkpointing
 140 commonly used in existing training frameworks [17, 33, 19]. While layer-wise checkpointing reduces
 141 memory, the total activation memory still scales linearly with the sequence length. In contrast, our
 142 block-wise approach ensures that the activation memory footprint remains constant, regardless of the
 143 sequence length, making it uniquely suited for training on extremely long contexts.

144 4 SeCOv2

145 4.1 Efficient KV Cache Management

146 Existing attention operators require the KV cache and its gradients to be stored in contiguous memory.
 147 This constraint poses no performance penalty in standard parallel training, where the KV cache is
 148 computed once per iteration and remains static. Consequently, most mainstream training frameworks
 149 lack a dedicated memory management mechanism for the KV cache.

150 In our block-wise training setting, however, the KV cache grows dynamically with each forward
 151 step. Enforcing memory contiguity under these conditions severely degrades system efficiency and
 152 memory utilization. Maintaining a contiguous KV cache requires expensive memory operations
 153 with each new block, such as `concatenate`, `copy`, and `reallocate`. Furthermore, allocating large,
 154 contiguous tensors for long contexts leads to memory fragmentation, preventing smaller, unused
 155 memory segments from being repurposed. These combined issues result in a substantial waste of
 156 both memory and computational resources.

157 Although pre-allocating a single large, contiguous buffer might seem like a solution, this approach is
 158 incompatible with dynamic memory-saving techniques like sequence-level activation recomputation
 159 or KV cache offloading. In our system, the KV cache is highly dynamic; its constituent blocks
 160 may be discarded, recomputed, or offloaded to CPU memory. A large, static pre-allocation would
 161 unnecessarily reserve memory that could otherwise be used by other components.

162 To quantify the memory inefficiency in a sys-
 163 tem lacking a dynamic KV cache manager, we
 164 benchmarked SeCOv1 [23] with a block-wise
 165 regime. As shown in Figure 3, even when using
 166 Flash Attention [9] and KV cache offloading,
 167 the system’s memory consumption exhibits a
 168 sharp inflection point beyond a 50K token con-
 169 text. Past this point, the memory footprint for
 170 each additional 10K tokens increases at a rate
 171 nearly three times higher than theoretically ex-
 172 pected. After eliminating other potential causes,
 173 including the attention operator itself, we con-
 174 cluded that this anomalous memory growth is
 175 a direct consequence of the lack of a dynamic
 176 memory management mechanism.

177 **Paged Management for the KV Cache in Training.** To address these challenges and maximize
 178 trainable context length, we introduce a paged memory management system for both the KV cache
 179 and its gradients, tightly integrated with our block-wise training framework. Existing paged attention
 180 implementations [21, 41] are designed for inference and lack support for backpropagation. We extend

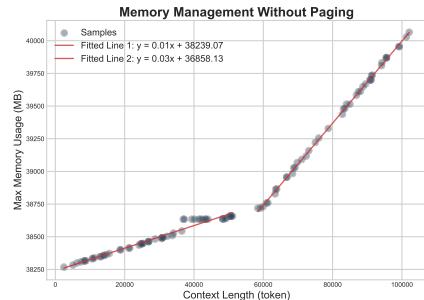


Figure 3: Without a paging mechanism, peak memory allocation for the long-context KV cache can be 3× higher than anticipated.

181 this approach by developing a paged management system for the KV cache gradients. Specifically,
182 as depicted in Figure 2, our mechanism partitions the KV cache tensors M_i and their corresponding
183 gradients dM_i (from Eq. 4) into fixed-size, non-contiguous pages. These pages are stored in a
184 pre-allocated memory pool, and a page table tracks their physical addresses.

185 **Optimizing for L2 Cache Locality.** A potential drawback of paged memory is reduced data locality,
186 which can lower the L2 cache hit rate. To overcome this, our memory allocator attempts to assign
187 pages from the same logical block to physically contiguous page frames whenever possible. This
188 simple heuristic yields a measurable improvement in system throughput.

189 **Decoupling the KV Cache from the Autograd System.** Managing the KV cache via PyTorch’s
190 autograd system, while straightforward, introduces significant overhead. Standard autograd execution
191 requires allocating intermediate buffers for gradients and performing separate read-add-write operations
192 on the cache’s .grad tensor. This process is inefficient for the large and frequently accessed KV
193 cache and complicates advanced memory management like offloading. To circumvent this, we bypass
194 autograd for KV cache gradients (see Figure 2). Our CUDA Op computes the gradient and uses a
195 single atomic_add operation to directly update the gradient tensor in place. This approach eliminates
196 intermediate memory allocation and redundant data movement, increasing system throughput
197 especially in large-context scenarios where the KV cache consumes tens of gigabytes.

198 **Eliminating Redundant KV Computation in Activation Recomputation.** Activation recomputation
199 saves memory by re-calculating intermediate results during the backward pass. However, a naive
200 implementation recomputes all activations, including the KV cache. This is unnecessary, as the KV
201 cache is intentionally preserved across forward passes and remains available for the backward pass.
202 We therefore implement a specialized recomputation pipeline that selectively regenerates only the
203 required intermediate results, entirely skipping the redundant KV projection. This targeted approach
204 avoids wasteful computation and improves system throughput.

205 4.2 Scaling Context Length via KV Cache Offloading

206 The KV cache and its gradient (which are isolated from the autograd system) are the sole components
207 whose memory demands scale linearly with the context length. This overhead presents a significant
208 bottleneck.

209 To overcome this limitation, we leverage the temporal locality inherent in the Transformer’s sequential,
210 layer-by-layer execution. We introduce a strategy of asynchronous CPU offloading and pre-fetching
211 that minimizes peak GPU memory by retaining only the KV cache for the currently active layer.

212 During computation for layer i , dedicated CUDA streams concurrently pre-fetch the KV cache
213 for layer $i + 1$ from the CPU and offload the cache for layer $i - 1$ to the CPU. This data flow is
214 reversed during the backward pass, as illustrated in Figure 1 (Bottom). When layer computations
215 are sufficiently large (*i.e.*, using a large block size), the latency from these asynchronous transfers is
216 effectively masked by the main computation.

217 Experiments on a single A100 GPU¹
218 show that for 4096 block size, our
219 method achieves a GPU-CPU transfer
220 throughput approximately double that
221 of FlashAttention-2.5.8 [9] (Table 2).
222 This result indicates that, in theory, the
223 overhead from CPU offloading can be
224 entirely masked by computation.

Table 2: Our CPU offload mechanism’s throughput is significantly higher than that of FlashAttention-2.5.8, ensuring the communication overhead is fully masked by computation.

Context	KV Size (GB)	Latency (s)		Throughput (GB/s)	
		flash attn	pre-fetch	flash attn	pre-fetch
100K	0.4	0.0324	0.0172	12.32	23.25
1M	4.0	0.3285	0.1668	12.17	24.17

225 4.3 Overcoming Training Time with Sparse Attention

226 With memory efficiency achieved, excessive training time becomes the primary bottleneck. To
227 address this, we incorporate Native Sparse Attention (NSA) [42], which holds CPU offloading
228 communication overhead constant and scales attention computation nearly linearly with context
229 length. This capability is crucial for processing extremely long contexts. Extensive experiments
230 confirm that training with NSA on ultra-long texts leads to significant efficiency gains.

¹The A100 features faster computation but a lower bandwidth compared to the H20.

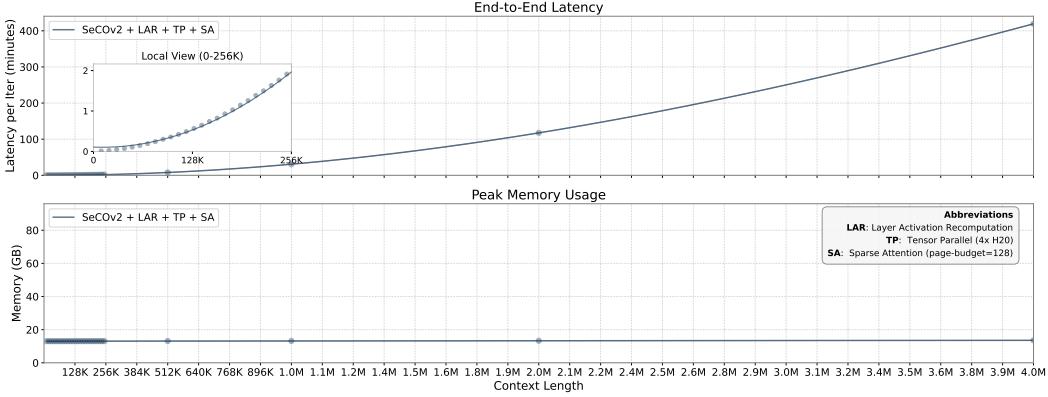


Figure 4: **Scalability of SeCOv2 for Training 7B Models.** SeCOv2, combined with tTensor parallelism (TP) and sparse attention (SA), achieves a nearly constant memory footprint while maintaining manageable latency growth. The limited gains from sparse attention are due to Triton kernel implementation details, with significant improvements anticipated in future CUDA releases.

231 5 Experiments

232 SeCOv2 is a general purpose framework for long context training, designed to be agnostic to specific
 233 training recipes and data. Consequently, our experiments (Section 5.1 and 5.2) benchmark the
 234 intrinsic performance of the framework itself, with a primary focus on its time and memory efficiency.

235 5.1 Main Results

236 **Experimental Setup.** To evaluate time and memory efficiency, we use the Qwen2.5-7B [40] model
 237 on the PG19 [32] dataset. We duplicated and concatenated the dataset as necessary to achieve the
 238 desired sequence lengths. All experiments were conducted using bf16 precision with the AdamW
 239 optimizer. Since our focus is on efficiency rather than model convergence, all other hyperparameters
 240 were kept at their default values.

241 All baseline methods utilize a GQA-compatible version of FlashAttention [9] with a per GPU batch
 242 size of 1. All experiments involve full-parameter fine-tuning. For SeCOv2, we always set the block
 243 size to 4096 and the page size to 128. For experiments involving sparse attention, we consistently
 244 use a page budget of 128. We measure peak memory usage using `max_memory_allocated()` and
 245 iteration latency using CUDA events. All reported results are the minimum from 3 runs.

246 *By default, methods do not employ acceleration techniques such as tensor parallelism (TP), layer
 247 activation recomputation (LAR), or sparse attention (SA) unless explicitly indicated by a suffix.*

248 **Efficiency and Scalability up to 4M Tokens.** We benchmark the per-iteration training time and
 249 memory of SeCOv2 on sequence lengths ranging from 2K to 4M tokens, comparing full and sparse
 250 attention implementations. As shown in Figure 4, SeCOv2 achieves high memory efficiency with
 251 both attention mechanisms. With sparse attention, the training time exhibits a more favorable growth.

252 **Comparison with ZeRO3-Offload.** Since both SeCOv2 and ZeRO3-Offload [33] utilize CPU
 253 offloading, we conduct a direct comparison to evaluate their relative time and memory efficiency.
 254 We evaluate three DeepSpeed ZeRO3 configurations: (1) pure ZeRO3, (2) ZeRO3 with parameter
 255 offloading, and (3) ZeRO3 with both parameter offloading and layer activation recomputation. The
 256 ZeRO3 experiments are run on 8×H20 GPUs with a per device batch size of 1, and we report statistics
 257 from GPU 0 (see Figure 5). SeCOv2 is run on a single H20 GPU for this comparison.

258 **Comparison with Context Parallelism Methods.** We benchmark SeCOv2 against CP methods,
 259 which also partition input sequences into blocks. We selected Ring Flash Attention (RFA) [25], a
 260 heavily optimized and representative open-source implementation of Ring Attention [27], as our
 261 primary baseline. Both methods were benchmarked for per device training throughput and memory

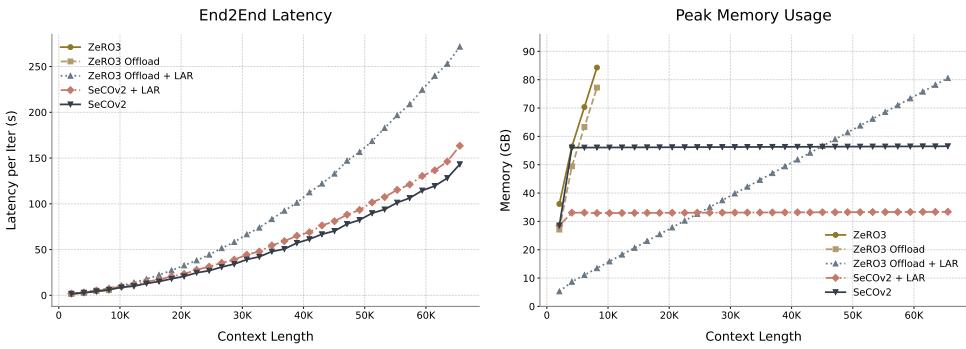


Figure 5: Compared to the widely-recognized DeepSpeed ZeRO3 [33], SeCOv2 demonstrates substantially higher per-GPU efficiency with respect to both computation time and memory usage.

Table 3: Throughput comparison with context parallel methods. We compared the per-GPU training throughput for 128K and 256K context lengths. Results shows that the SOTA implementation, Ring Flash Attention (RFA) [25], achieves a throughput comparable to that of SeCOv2. (The results highlighted in red are from [6].)

Context	Ring Attention [27]		RFA [25]		SeCOv2-TP	
	8×A100	16×TPUv4	8×H20	1×H20	4×H20	
128K	-	-	232 tok/s	240 tok/s	229 tok/s	
256K	50 tok/s	49 tok/s	125 tok/s	139 tok/s	121 tok/s	

usage. The results in Table 3 and Figure 6 demonstrate that SeCOv2 delivers throughput comparable to RFA but with substantially improved memory efficiency.

5.2 Ablation Results

Ablation on Block Size. We ablate the block size, doubling it from 512 to 4096. As shown in Figure 7, computation time saturates with increasing block size, yielding diminishing returns beyond 4096. Meanwhile, memory consumption shows minimal difference. Notably, as the block size approaches infinity, our method degenerates to standard parallel training.

Ablation on CPU Offload. To hide the data transfer latency with computation, we evaluated the effect of CPU offloading on system latency and memory footprint across various context lengths (Figure 8). The results show that our

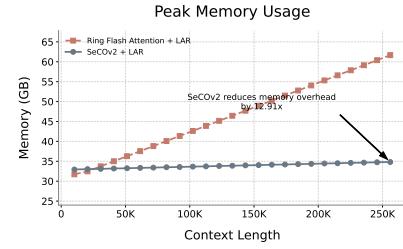


Figure 6: Memory usage compared to context parallel methods.

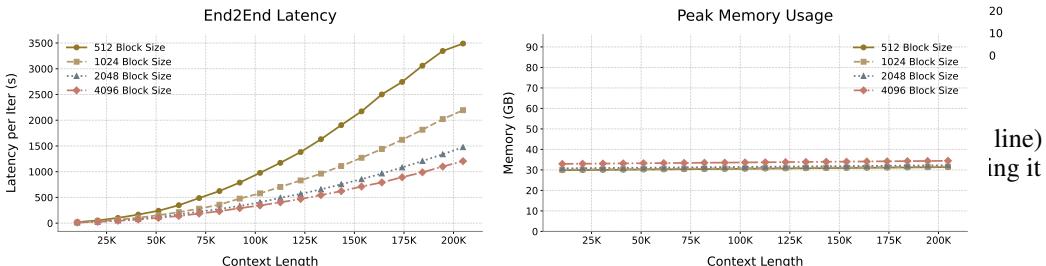


Figure 7: Increasing the block size from 512 to 4096 yields diminishing marginal returns, with performance eventually converging to the parallel training baseline (corresponds to an ∞ block size).

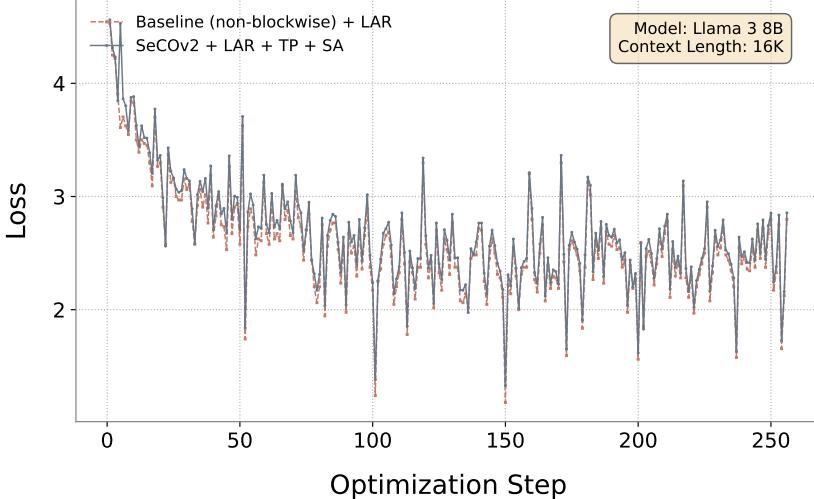


Figure 9: **Accuracy Validation.** Our SeCOv2 implementation, using 25% sparse attention and 4-GPU tensor parallelism, achieves nearly identical final accuracy to a baseline employing FlashAttention-2 (v2.5.8). This result confirms the correctness and reliability of our system.

274 method introduces a minor latency overhead of
 275 5%, primarily due to data transfer initiation. In
 276 return, it substantially reduces the memory foot-
 277 print, which remains nearly constant for con-
 278 texts up to 200K. These findings confirm that
 279 the KV cache is the primary memory bottleneck
 280 and demonstrate that offloading is an effective
 281 mitigation strategy.

282 5.3 Accuracy Validation

283 To evaluate SeCOv2, we quantify discrepancies
 284 arising from its custom sparse attention, tensor parallelism, and block-wise optimizations. We
 285 compare the loss reduction of SeCOv2 against a baseline using parallel training with FlashAttention
 286 (v2.5.8) [9]. As shown in Figure 9, at a 25% sparsity level (achieved with a page size and sparse page
 287 budget of 64), the final error after 256 training steps is approximately 0.03.

288 6 Conclusion and Limitations

289 **Conclusion.** In this work, we introduced SeCOv2, a highly memory-efficient training framework that
 290 significantly lowers the resource barrier for training LLMs on million-token contexts. By combining
 291 a block-recurrent training strategy with on-the-fly activation recomputation, we achieve a constant
 292 memory footprint for activations, irrespective of sequence length. We address the subsequent KV
 293 cache bottleneck through a suite of co-designed optimizations: a paged memory manager for the
 294 cache and its gradients, autograd-decoupled cache operations, and an asynchronous CPU offloading
 295 mechanism. The synergy of these techniques enables training a 4-million-token context Qwen2.5-7B
 296 model on a single 96GB GPU, a task that previously required large-scale clusters. Our framework
 297 is designed to be a plug-and-play wrapper, compatible with existing training systems, to help
 298 democratize research and development in ultra-long-context LLMs.

299 **Limitations.** The primary limitation of SeCOv2 is the trade-off between memory efficiency and
 300 computational time. The serial, block-wise processing and activation recomputation introduce
 301 latency overhead compared to standard parallel training. While we demonstrate that this overhead is
 302 manageable and can be partially mitigated by larger block sizes and integration with tensor parallelism
 303 or sparse attention, it remains an inherent characteristic of the design. Furthermore, our asynchronous

304 CPU offloading mechanism relies on the computational workload of a block being sufficient to hide
305 data transfer latency, which may not hold for smaller models or block sizes. Finally, this work focuses
306 on the systems-level efficiency of the training framework itself; we do not propose new long-context
307 learning recipes or data curation strategies, which are orthogonal but essential for achieving optimal
308 model performance.

309 **References**

- 310 [1] L. AI. Litgpt. <https://github.com/Lightning-AI/litgpt>, 2023.
- 311 [2] J. Ainslie, J. Lee-Thorp, M. de Jong, Y. Zemlyanskiy, F. Lebron, et al. GQA: Training
312 generalized multi-query transformer models from multi-head checkpoints. In *EMNLP*, 2023.
- 313 [3] C. An, F. Huang, J. Zhang, S. Gong, X. Qiu, et al. Training-free long-context scaling of large
314 language models. In *ICML*, 2024.
- 315 [4] Y. Bai, X. Lv, J. Zhang, Y. He, J. Qi, L. Hou, et al. LongAlign: A recipe for long context
316 alignment of large language models. In *EMNLP*, 2024.
- 317 [5] W. Bencheikh, J. Finkbeiner, and E. Neftci. Optimal gradient checkpointing for sparse and
318 recurrent architectures using off-chip memory. *arXiv preprint arXiv:2412.11810*, 2024.
- 319 [6] W. Brandon, A. Nrusimha, K. Qian, Z. Ankner, T. Jin, et al. Striped attention: Faster ring
320 attention for causal transformers. 2023.
- 321 [7] S. Chen, S. Wong, L. Chen, and Y. Tian. Extending context window of large language models
322 via positional interpolation. *arXiv preprint arXiv:2306.15595*, 2023.
- 323 [8] Y. Chen, S. Qian, H. Tang, X. Lai, Z. Liu, et al. Longlora: Efficient fine-tuning of long-context
324 large language models. In *ICLR*, 2024.
- 325 [9] T. Dao. FlashAttention-2: Faster attention with better parallelism and work partitioning. In
326 *ICLR*, 2024.
- 327 [10] Y. Ding, L. L. Zhang, C. Zhang, Y. Xu, N. Shang, et al. Longrope: extending llm context
328 window beyond 2 million tokens. In *ICML*, 2024.
- 329 [11] D. Y. Fu, T. Dao, K. K. Saab, A. W. Thomas, A. Rudra, et al. Hungry hungry hippos: Towards
330 language modeling with state space models. In *ICLR*, 2023.
- 331 [12] D. Y. Fu, E. L. Epstein, E. Nguyen, A. W. Thomas, M. Zhang, T. Dao, A. Rudra, and C. Re.
332 Simple hardware-efficient long convolutions for sequence modeling. In *ICLR*, 2023.
- 333 [13] T. Gao, A. Wettig, H. Yen, and D. Chen. How to train long-context language models (effectively).
334 *arXiv preprint arXiv:2410.02660*, 2024.
- 335 [14] T. Gao, A. Wettig, H. Yen, and D. Chen. How to train long-context language models (effectively).
336 In *ICLR*, 2025.
- 337 [15] A. Grattafiori, A. Dubey, A. Jauhri, A. Pandey, A. Kadian, et al. The llama 3 herd of models.
338 *arXiv preprint arXiv:2407.21783*, 2024.
- 339 [16] A. Gu, K. Goel, and C. Ré. Efficiently modeling long sequences with structured state spaces. In
340 *ICLR*, 2022.
- 341 [17] S. Gugger, L. Debut, T. Wolf, P. Schmid, Z. Mueller, et al. Accelerate: Training and infer-
342 ence at scale made simple, efficient and adaptable. [https://github.com/huggingface/](https://github.com/huggingface/accelerate)
343 *accelerate*, 2022.
- 344 [18] C. Han, Q. Wang, H. Peng, W. Xiong, Y. Chen, et al. Lm-infinite: Zero-shot extreme length
345 generalization for large language models. In *NAACL*, 2024.
- 346 [19] S. A. Jacobs, M. Tanaka, C. Zhang, M. Zhang, R. Y. Aminadabi, et al. Deepspeed ulysses:
347 System optimizations for enabling training of extreme long sequence transformer models. In
348 *PODC*, 2024.
- 349 [20] H. Jin, X. Han, J. Yang, Z. Jiang, Z. Liu, et al. Llm maybe longlm: Selfextend llm context
350 window without tuning. In *ICML*, 2024.
- 351 [21] W. Kwon, Z. Li, S. Zhuang, Y. Sheng, L. Zheng, C. H. Yu, J. Gonzalez, H. Zhang, and I. Stoica.
352 Efficient memory management for large language model serving with pagedattention. In *SOSP*,
353 2023.

- 354 [22] D. Li, R. Shao, A. Xie, E. P. Xing, X. Ma, et al. DISTFLASHATTN: Distributed memory-
355 efficient attention for long-context LLMs training. In *CoLM*, 2024.
- 356 [23] W. Li, Y. Zhang, G. Luo, D. Yu, and R. Ji. Training long-context llms efficiently via chunk-wise
357 optimization. In *ACL*, 2025.
- 358 [24] W. Liang, T. Liu, L. Wright, W. Constable, A. Gu, et al. Torch titan: One-stop pytorch native
359 solution for production ready llm pre-training. *arXiv preprint arXiv:2410.06511*, 2025.
- 360 [25] Z. Lin. Ring flash attention: Ring attention implementation with flash attention. <https://github.com/zhuolin/ring-flash-attention>, 2025.
- 362 [26] A. Liu, B. Feng, B. Xue, B. Wang, B. Wu, et al. Deepseek-v3 technical report. *arXiv preprint
363 arXiv:2412.19437*, 2025.
- 364 [27] H. Liu, M. Zaharia, and P. Abbeel. Ringattention with blockwise transformers for near-infinite
365 context. In *ICLR*, 2024.
- 366 [28] J. Liu, D. Zhu, Z. Bai, Y. He, H. Liao, H. Que, Z. Wang, C. Zhang, G. Zhang, J. Zhang, et al. A
367 comprehensive survey on long context language modeling. *arXiv preprint arXiv:2503.17407*,
368 2025.
- 369 [29] E. Lu, Z. Jiang, J. Liu, Y. Du, T. Jiang, et al. Moba: Mixture of block attention for long-context
370 llms. *arXiv preprint arXiv:2502.13189*, 2025.
- 371 [30] B. Peng, J. Quesnelle, H. Fan, and E. Shippole. YaRN: Efficient context window extension of
372 large language models. In *ICLR*, 2024.
- 373 [31] M. Poli, S. Massaroli, E. Nguyen, D. Y. Fu, T. Dao, et al. Hyena hierarchy: towards larger
374 convolutional language models. In *ICML*, 2023.
- 375 [32] J. W. Rae, A. Potapenko, S. M. Jayakumar, C. Hillier, and T. P. Lillicrap. Compressive
376 transformers for long-range sequence modelling. In *ICLR*, 2020.
- 377 [33] S. Rajbhandari, J. Rasley, O. Ruwase, and Y. He. Zero: memory optimizations toward training
378 trillion parameter models. In *SC*, 2020.
- 379 [34] M. Shoeybi, M. Patwary, R. Puri, P. LeGresley, J. Casper, et al. Megatron-lm: Training multi-
380 billion parameter language models using model parallelism. *arXiv preprint arXiv:1909.08053*,
381 2020.
- 382 [35] N. S. Sohoni, C. R. Aberger, M. Leszczynski, J. Zhang, and C. Ré. Low-memory neural network
383 training: A technical report. *arXiv preprint arXiv:1904.10631*, 2022.
- 384 [36] R. C. Staudemeyer and E. R. Morris. Understanding lstm – a tutorial into long short-term
385 memory recurrent neural networks. *arXiv preprint arXiv:1909.09586*, 2019.
- 386 [37] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, et al. Attention is all you need. In
387 *NeurIPS*, 2017.
- 388 [38] X. Wang, M. Salmani, P. Omidi, X. Ren, M. Rezagholizadeh, et al. Beyond the limits: a survey
389 of techniques to extend the context length in large language models. In *IJCAI*, 2024.
- 390 [39] C. Xu, W. Ping, P. Xu, Z. Liu, B. Wang, et al. From 128k to 4m: Efficient training of ultra-long
391 context large language models. *arXiv preprint arXiv:2504.06214*, 2025.
- 392 [40] A. Yang, B. Yang, B. Zhang, B. Hui, B. Zheng, B. Yu, et al. Qwen2.5 technical report. *arXiv
393 preprint arXiv:2412.15115*, 2025.
- 394 [41] Z. Ye, L. Chen, R. Lai, W. Lin, Y. Zhang, et al. Flashinfer: Efficient and customizable attention
395 engine for llm inference serving. *arXiv preprint arXiv:2501.01005*, 2025.
- 396 [42] J. Yuan, H. Gao, D. Dai, J. Luo, L. Zhao, et al. Native sparse attention: Hardware-aligned and
397 natively trainable sparse attention. 2025.