

---

# Week 9 Reading Summaries: PagedAttention and FlashAttention

---

**Amogh Patankar**  
University of California, San Diego  
3869 Miramar Street Box #3037, La Jolla, CA- 92092  
apatankar@ucsd.edu

## Abstract

1 In this set of readings, we go over PagedAttention and FlashAttention, which  
2 are two methods of attention algorithms that use various strategies to reduce the  
3 number of memory read/writes that occur between high bandwidth memory (HBM)  
4 and GPU on-chip SRAM. Paged Attention uses classic virtual memory and paging  
5 techniques for operating systems, and the authors build vLLM, which is a serving  
6 system for minimal KV cache waste and sharing of KV cache. Flash Attention is  
7 extended to block-sparse attention, and they train transformers in a more efficient  
8 manner

## 9 1 PagedAttention

### 10 1.1 Introduction

11 In recent times, with LLMs on the rise, especially GPT, PaLM, and many more, there have been a  
12 wider variety of applications. Namely, chatbots and assistants are rising, with many cloud companies  
13 providing these services as features and hosted services. But, these applications require a lot of  
14 compute, i.e. GPUs. Given these costs, LLM serving systems come to the forefront of our attention,  
15 in order to improve throughput.  
16 Improving throughput can be achieved by multibatching, but memory needs to be managed properly.  
17 The authors note that existing systems fall short because the KV cache memory isn't managed well.  
18 Moreover, these systems can't use memory sharing properly. The authors, in this paper, build vLLM,  
19 which is a LLM servicing engine using PagedAttention, achieving close to zero waste in the KV  
20 cache memory

### 21 1.2 Background

#### 22 1.2.1 Transformer-Based LLMs

23 For transformer based LLMs, the whole idea is modeling probabilities as a list of tokens; this is done  
24 using a technique called autoregressive decomposition. A self-attention layer within a transformer  
25 uses linear transformations, attaining query, key, and value vectors, finally computing the attention  
26 score.

#### 27 1.2.2 LLM Service and Autoregressive Generation

28 LLMs are essentially used as a generational API service; the KV cache stores key and value vectors  
29 in the generation process. The prompt phase takes the prompt and computes probabilities of new  
30 tokens, while the autoregressive generation phase generates the new tokens sequentially.

### 31 1.2.3 Batching Techniques for LLMs

32 Batching multiple requests in serving LLMs can help to decrease compute utilization; however, it's  
33 non-trivial. As requests arrive at different times, and the requests' input and output lengths may be  
34 different, standard batching doesn't work. Hence, batching techniques like cellular or iteration-level  
35 scheduling are proposed.

## 36 1.3 Memory Challenges in LLM Serving

37 Surpassing memory requirements has a few challenges, namely, large KV Cache, complex decoding  
38 algos, and scheduling for unknown input/output lengths.

### 39 1.3.1 Memory Management in Existing Systems

40 LLM serving systems require a dynamic allocation of memory for newer systems, storing KV caches  
41 efficiently. Typical static allocation of memory ends with huge pure memory waste.

## 42 1.4 Method

43 The authors develop PagedAttention to solve this exact problem.

### 44 1.4.1 PagedAttention

45 PagedAttention allows for the storage of keys and values in non-contiguous memory space, by  
46 creating KV blocks. Each block stores an attention score, at the  $ij^{th}$  index. In doing so, the KV  
47 blocks are stored in non-contiguous memory, meaning vLLM can handle memory more efficiently.

### 48 1.4.2 KV Cache Manager

49 vLLM's memory manager operates similarly to virtual memory in OS'. vLLM organizes KV cache  
50 as fixed-size KV blocks (like pages), with a request's KV cache represented as a series of KV blocks.  
51 The block manager maintains block tables.

### 52 1.4.3 Decoding with PagedAttention and vLLM

53 When doing decoding, vLLM doesn't reserve memory; it reserves KV blocks. In the first autore-  
54 gressive decoding step, vLLM generates new tokens with PagedAttention on various blocks, and  
55 stores the newest KV cache in a new block for the second decoding step. Finally, vLLM dynamically  
56 assigns new physical blocks to logical blocks as tokens get generated.

### 57 1.4.4 Application to Other Decoding Scenarios

- 58 • Parallel Sampling: For parallel sampling scenarios, only one copy of the prompt's state will  
59 be reserved, and we keep a reference count. When generating, we use the copy-on-write  
60 mechanism.
- 61 • Beam Search: For beam search, the beam width determines the top candidates for each step.  
62 We have beam candidates vLLM allocates blocks, and vLLM reduces overhead by sharing  
63 physical blocks.
- 64 • Shared prefix: vLLM reserves a set of physical blocks for shared prefixes.
- 65 • Mixed Decoding Methods: vLLM allows for different decoding preferences.

### 66 1.4.5 Scheduling and Preemption

67 vLLM adopts a first-come-first-serve scheduler, ensuring fairness and preventing starvation.

- 68 • Swapping: vLLM selects a set of sequence to remove, then swaps their KV cache to the  
69 CPU.
- 70 • Recomputation: vLLM recomputes the KV cache for each preempted sequences that are  
71 rescheduled.

## 72 1.5 Distributed Execution

73 vLLM efficiently manages memory for large language models using Megatron-LM’s tensor model  
74 parallelism and an SPMD execution strategy. Linear layers are partitioned for block-wise matrix  
75 multiplication, with GPUs synchronizing intermediate results via all-reduce. A centralized KV cache  
76 manager handles memory mapping, which all GPU workers share. Each worker stores only part  
77 of the KV cache for its assigned attention heads. The scheduler coordinates memory management  
78 upfront, allowing GPU workers to execute independently and sync only intermediate results.

## 79 1.6 Implementation

80 vLLM is an end-to-end serving system with a FastAPI frontend and a GPU-based inference engine,  
81 featuring Python and CUDA code for control and custom kernels, supporting popular LLMs like  
82 GPT, OPT, and LLaMA with NCCL for distributed communication.

### 83 1.6.1 Kernel-level Optimization

84 To optimize PagedAttention’s memory access patterns, vLLM introduces custom GPU kernels that  
85 fuse key operations: (1) reshaping, block writing, and saving KV cache to minimize kernel launch  
86 overheads; (2) block reading and attention to improve efficiency with coalesced memory access; and  
87 (3) a fused block copy kernel to reduce overhead from fragmented data movement.

## 88 1.7 Supporting Various Decoding Algorithms

89 vLLM supports decoding algorithms using three key methods: ‘fork’ (to create new sequences),  
90 ‘append’ (to add tokens), and ‘free’ (to delete sequences), enabling techniques like parallel sampling,  
91 beam search, and prefix sharing.

## 92 1.8 Evaluation

93 The authors evaluate vLLM under a bunch of different workloads; they find that vLLM sustains much  
94 higher request rates than existing solutions.

## 95 1.9 Ablation Studies

96 The authors evaluate design choices through ablation experiments. vLLM’s PagedAttention introduces  
97 some overhead (20–26% higher latency) in the attention kernel due to dynamic block mapping,  
98 but this impact is limited to the attention operator. Optimal performance is achieved with a block  
99 size of 16, balancing GPU utilization and minimizing fragmentation. For recovery, recomputation  
100 outperforms swapping with small block sizes, while both methods perform similarly for medium  
101 block sizes (16–64).

### 102 1.10 Discussion

103 The authors discuss the idea of applying virtual memory and paging to other GPU workloads.

## 104 2 FlashAttention

### 105 2.1 Introduction

106 The authors develop FlashAttention, a novel attention algorithm used to improve the efficiency of  
107 transformers by reducing memory access overhead. They use tiling and recomputation methods to  
108 restructure attention computation, avoiding repeatedly read/writes of the attention matrix to GPU  
109 memory. The authors’ approach enables FlashAttention to run close to 8x faster, and consumes less  
110 memory than vanilla attention. FlashAttention is able to accelerate model training, getting much  
111 faster results on BERT-large, GPT-2, and long-range tasks. FlashAttention also improves model  
112 quality; block-sparse FlashAttention enhances speed and scalability even further, achieving faster  
113 performance than all known approximate attention methods.

## 114 2.2 Background

115 The authors provide some background about performance of DL operations on GPUs.

### 116 2.2.1 Hardware Performance

117 The authors focus on GPU performance, specifically its memory hierarchy and execution model.  
118 GPUs have larger, slower HBM, with small, fast on-chip SRAM, making it crucial to use SRAM  
119 efficiently. For kernel execution, data is loaded from HBM to SRAM, computed, and written back to  
120 HBM. Depending on the balance of computation and memory access, operations can be compute-  
121 bound or memory-bound. To optimize memory-bound operations, the authors used kernel fusion,  
122 reducing redundant HBM access by combining multiple ops. However, in model training, saving  
123 intermediate values for backpropagation limits the full benefits of naive kernel fusion.

### 124 2.2.2 Standard Attention Implementation

125 In regular attention,  $S = QK^T$ ,  $P = \text{softmax}(S) \in R^{N \times N}$ , and  $O = PV \in R^{N \times d}$ , with row-  
126 wise softmax. The problem with this, is the memory-bound operations, slowing down the entire  
127 computation.

## 128 2.3 FlashAttention: Algorithm, Analysis, and Extensions

### 129 2.3.1 An Efficient Attention Algorithm With Tiling and Recomputation

130 The authors compute the attention and write it to HBM, reducing the r/w to HBM. Then, they apply  
131 tiling and recomputation to overcome the exact computations for attention. With tiling, the attention  
132 is computed by blocks, and compute softmax block by block, then combine. For recomputation, the  
133 authors recompute S and P easily, and often, with gradient checkpoint.

### 134 2.3.2 Analysis: IO Complexity of FlashAttention

135 FlashAttention significantly reduces high-bandwidth memory (HBM) accesses compared to standard  
136 attention, having faster execution and lower memory usage, especially for typical head dimensions.  
137 The authors show that exact attention algorithms can't improve HBM accesses, and show that HBM  
138 accesses influence runtime. Increasing block sizes in FlashAttention reduces HBM accesses and  
139 runtime, though performance flattens due to arithmetic ops.

### 140 2.3.3 Extension: Block-Sparse FlashAttention

141 Block-sparse FlashAttention is an extension of FlashAttention, where the authors compute only  
142 the non-zero blocks of the attention matrix, i.e. dense blocks. This optimization has a  $O(N\sqrt{N})$   
143 complexity, decreasing from FlashAttention by a factor proportional to the sparsity.

## 144 2.4 Experiments

145 In their experiments, the authors show that FlashAttention outperforms standard attention, on models  
146 such as GPT-2 small, medium, and BERT by up to 3.5x. Moreover, FlashAttention and Block-Sparse  
147 FlashAttention outperform every model in long range, long sequences, and long-context.

## 148 2.5 Limitations and Future Directions

149 The authors note that future work can be compiling to CUDA, IO-Aware Deep Learning methods,  
150 and Multi-GPU IO-Aware methods, for data transfer between GPUs.