PagedAttention & vLLM for Efficient LLM Inference

Woosuk Kwon



Contents

• PagedAttention

- vLLM: A real-world open-source LLM serving system
 - \circ Models
 - Parallelism
 - Inference optimizations

• Q&A

Contents

• PagedAttention

- vLLM: A real-world open-source LLM serving system
 - \circ Models
 - Parallelism
 - Inference optimizations

• Q&A

The era of Large Language Models (LLMs)



. . .

OPT-175B serving demo (Fall 2022)

tl I	Shubham Saboo here Retweeted Charly Wargnier @DataChaz · Aug 11 Introducing OPT-175B, a free version of GPT-3! 7000000000000000000000000000000000000
	No login, no credit card needed!
	👉 opt.alpa.ai
	via @Saboo_Shubham_ #ai #ml #nlp
	★ Pree, Unimited OPT-175B Text Generation Write: This model might enerates townething offensive. No safety measures are in place as a fire service. Write: Write: Arroro Code Themasterice Themast
	Shubham Saboo 🦒
	♀ 17 ℃↓ 250 ♥ 956 ①

Serving LLMs is (surprisingly) slow and expensive

- A single A100 GPU can only serve < 1 request per second
 - \circ $\,$ Moderate size of model (13B parameters) and input
- A ton of GPUs are required for production-scale LLM services



LLM inference process



LLM inference process



LLM inference process



•••

Why is LLM inference inefficient?



Sequential dependency \rightarrow Hard for GPU parallelization

Batching?



Attention KV cache



Intermediate **vector** repr. ("Attention key & value")

Attention KV cache size

of

ıre



Artific

KV Cache is huge:

- Each token: ~1 MB.
- One full request: ~several GBs.

KV cache management in previous systems



- **Pre-allocates contiguous** memory to the request's maximum length
 - Convention in previous deep learning workloads with **static** input/output shapes
- Results in memory fragmentation
 - **Internal fragmentation** due to the **unknown** output length.
 - **External fragmentation** due to **non-uniform** per-request max lengths.

Significant memory waste in KV cache space

Only **20–40%** of KV Cache space is utilized to store actual token states



Yu et al. "Orca: A Distributed Serving System for Transformer-Based Generative Models" (OSDI 22).

PagedAttention

Application-level memory **paging** and **virtualization** for attention KV Cache



Paging KV cache space into KV blocks

Contiguous KV Cache

Artificio	I Intelli- gence	is	the							
-----------	---------------------	----	-----	--	--	--	--	--	--	--



KV block: a fixed-size block of memory that stores KV cache from left to right

block 0 block 1 block 2 block 3 Intelliblock 4 Artificial is the gence block 5 Block size = 4

KV Blocks

Virtualizing KV Cache



Attention mechanism with virtualized KV cache

- 1. Fetch non-contiguous KV blocks using the block table
- 2. Apply attention operation on the fly



KV Cache











Memory efficiency of PagedAttention

- Minimal internal fragmentation
 - Only happens at the last block of a sequence
 - # wasted tokens/seq < block size
 - Sequence: ~1000 tokens
 - Block size: ~10 tokens
- No external fragmentation





Paging enables sharing



Multiple outputs

Sharing KV blocks



Sharing KV blocks



More complex sharing: beam search



- Similar to process tree (fork & kill)
- Efficiently supported by paged attention and copy-on-write mechanism

Memory saving via sharing



Percentage = (#blocks saved by sharing) / (#total blocks without sharing) OPT-13B on 1x A100-40G with ShareGPT dataset

How does PagedAttention benefit LLM Serving?

Reduce memory fragmentation with paging Further reduce memory usage with sharing

Out of KV block memory

Physical KV blocks

Request	The	future	of	cloud		
A	computer	scientist	and	mathe- matician		
	computing	is	likely	to	≻	Ful
Request	renowned					
В	Alan	Turing	is	a		

Out of KV block memory



33

Cannot allocate a



Notes on preemption & recovery

Swap/recompute the whole request, since all previous tokens are required every step.

Swapping: smaller block sizes \rightarrow higher overhead due to small data transfers.

Recomputation: surprisingly fast since all token's KV cache can be computed in parallel.



Figure: Swap/Recomputation latency of 256 tokens.

Our strategy: Use recomputation when possible with FCFS policy

Comparison with operating system virtual memory

Analogies

OS pages ↔ KV blocks

• Reduce memory fragmentation

Shared pages across processes ↔ Shared KV blocks across samples

• Reduce memory waste via sharing

Differences

Single-level block table

• Block table is tiny compared to the actual data.

Preemption & Recovery

- Request-level preemption
- Recomputation-based recovery
Throughput - greedy decoding



OPT-13B on 1xA100 40G with ShareGPT trace.

Throughput - beam search



OPT-13B on 1xA100 40G with Alpaca trace. Speedup: vLLM v.s. Orca(Pow2)

PagedAttention has become the industrial standard

🔍 🔍 🔍 🔍 🔍	DIA TensorRT-	LLM S	Superc	+				
\leftrightarrow \rightarrow C h	ttps:// dev		Gr	☆	*		۲	
Technical Blog	۹					Sub	oscrib	e >

requiring deep knowledge of C++ or NVIDIA CUDA.

TensorRT-LLM improves ease of use and extensibility through an open-source modular Python API for defining, optimizing, and executing new architectures and enhancements as LLMs evolve, and can be customized easily.

For example, MosaicML has added specific features that it needs on top of TensorRT-LLM seamlessly and integrated them into their inference serving. Naveen Rao, vice president of engineering at Databricks notes that "it has been an absolute breeze."

"TensorRT-LLM is easy to use, feature-packed with streaming of tokens, in-flight batching, <u>paged-attention</u>, quantization, and more, and is efficient," Rao said. "It delivers state-of-theart performance for LLM serving using NVIDIA GPUs and allows us to pass on the cost savings to our customers."

Performance comparison

Summarizing articles is just one of the many applications of

TensorRT-LLM

•••	huggingface/text-gen	eration-i ×	+ ~
$\leftrightarrow \rightarrow c$	🔒 h 🔍 😫 🔤	₫ ☆	• ≕ ⊡ @ :

E README.md

- Serve the most popular Large Language Models with a simple launcher
- Tensor Parallelism for faster inference on multiple GPUs
- Token streaming using Server-Sent Events (SSE)
- Continuous batching of incoming requests for increased total throughput
- Optimized transformers code for inference using flashattention and Paged Attention on the most popular architectures
- Quantization with bitsandbytes and GPT-Q
- Safetensors weight loading
- Watermarking with A Watermark for Large Language
 Models
- Logits warper (temperature scaling, top-p, top-k, repetition penalty, more details see transformers.LogitsProcessor)
- Stop sequences
- Log probabilities
- Production ready (distributed tracing with Open
 Televative Depresentation)

HuggingFace TGI



new and innovative products.

Fast, Affordable LLM Inference with Fireworks

Efficient Inference

Efficient inference of LLMs is an active area of research, but we are industry veterans from the PyTorch team specializing in performance optimization. We use model optimizations including <u>multi/group query attention</u> optimizations, sharding, quantization, kernel optimizations, CUDA graphs, and custom cross-GPU communication primitives. At the service level, we employ continuous batching, paged attention, prefill disaggregation, and pipelining to maximize throughput and reduce latency. We carefully tune deployment parameters including the level of parallelism and hardware choice for each model.

🖑 120

Fireworks AI

. . .

Contents

• PagedAttention

- vLLM: A real-world open-source LLM serving system
 - \circ Models
 - Parallelism
 - Inference optimizations

• Q&A

✓LLM's Goal

Build the fastest and

easiest-to-use open-source

LLM inference & serving engine

vLLM Today

https://github.com/vllm-project/vllm







vLLM Community



- 600+ PRs a month
- 890+ contributors
- 30+ major contributors from UCB, RedHat, Anyscale, Roblox, Meta, etc.

vLLM API (1): LLM class

A Python interface for offline batched inference

```
from vllm import LLM
# Example prompts.
prompts = ["Hello, my name is", "The capital of France is"]
# Create an LLM with HF model name.
llm = LLM(model="meta-llama/Meta-Llama-3.1-8B")
# Generate texts from the prompts.
outputs = llm.generate(prompts)
```

vLLM API (2): OpenAI-compatible server

A FastAPI-based server for online serving

Server

\$ <u>vllm serve</u> meta-llama/Meta-Llama-3.1-8B

Client

```
$ curl http://localhost:8000/v1/completions \
    -H "Content-Type: application/json" \
    -d '{
        "model": "meta-llama/Meta-Llama-3.1-8B",
        "prompt": "San Francisco is a",
        "max_tokens": 7,
        "temperature": 0
    }'
```

Key Focuses in vLLM

Models

Parallelism

Inference optimizations

Performance engineering

Key Focuses in vLLM

Models

Parallelism

Inference optimizations

Performance engineering

Broad Model Support

- Transformer-like LLMs (e.g., Llama)
- Mixture-of-Expert LLMs (e.g., DeepSeek)
- Multi-modal LLMs (e.g., Qwen-VL)
- State-Space Models (e.g., Mamba)
- Embedding Models (e.g. E5-Mistral)

Easy Model Integration

- vLLM is based on PyTorch and HuggingFace Transformers
- vLLM requires minimal changes to the model code



First-class Support for Multi-modal Models

vLLM provides a structured way to register models with various modalities (e.g., image, video, and audio)

```
@MULTIMODAL REGISTRY.register image input mapper(image input mapper)
@MULTIMODAL REGISTRY.register input mapper("video", video input mapper)
@MULTIMODAL REGISTRY.register max image tokens(get max image tokens)
@MULTIMODAL REGISTRY.register max multimodal tokens("video",
                                                    get max video tokens)
@INPUT REGISTRY.register_dummy_data(dummy_data_for_qwen2_vl)
@INPUT REGISTRY.register input processor(input processor for qwen2 vl)
class Qwen2VL(nn.Module, SupportsMultiModal):
   def init (self,
                config: Qwen2VLConfig,
                multimodal config: MultiModalConfig):
       super(). init ()
       . . .
```

Hybrid Attention Models

- Various attention variants to efficiently support long context:
 - Multi-head Latent Attention (DeepSeek)
 - Sliding window attention
 - Local chunked attention
 - State-space models (Mamba)
- Importantly, all of the above can be <u>mixed</u> in a single model
 - Gemma 3: Global attention + sliding window attention
 - Llama 4: Global attention + local chunked attention
 - Jamba: Global attention + Mamba

Hybrid memory allocator

• We are building a two-layer memory allocator/evictor to manage the heterogeneous memory:



Jenga: Effective Memory Management for Serving LLM
 with Heterogeneity

Manual CUDA kernels v.s. torch.compile

- vLLM once maintained a large set of CUDA kernels, with different kernel fusion and quantization.
- The modifications to the model code because of these kernels make the model code less readable and harder to maintain.



Key Focuses in vLLM



Parallelism

Inference optimizations

Performance engineering

Why do we need parallelism?

- Models are getting much bigger than a single GPU
 - Deepseek V3: 671B parameters \rightarrow ~700GB in FP8
 - NVIDIA H100: 80GB HBM
 - Memory is not only used for parameters, activations and KV cache also takes memory
- Use more GPUs to multiple the compute power
 - Reduce serving latency
 - Increase serving throughput

Things to consider for parallelism

- Theoretically, all tensors can be arbitrarily distributed and all dimensions of all tensors can be sharded.
- Goal 1: Minimize communication
 - Limit the number of bytes transferred
 - Hide the overhead via overlapping with computation
 - Network is also heterogeneous
 - Latest NVLink: 1.4TB/s
 - Infiniband: 50GB/s
- Goal 2: Minimize waiting from data dependency and stragglers' effect
 - A system is often bottlenecked by the slowest component

Data parallelism (mostly outside of vLLM)

- Replicated engines
- A load-balancer routes different requests to different engines based on
 - \circ Engine load
 - KV cache memory usage
 - (Prefix) Caching
- No communication between engines
- No parameter memory saving, which can limit the KV cache size, and thus serving throughput
- No reduction in latency



Tensor parallelism (in vLLM)

- Partition the weights in the linear layers
- Reduce the serving latency if the network is fast
- KV Cache can also be partitioned
- Heavy communication across TP shards. Often require NVLink to work effectively.



Expert parallelism (in vLLM)

- Mixture of Expert: A set of linear layers (experts) where each token only pick a subset to run on
- Distribute different experts to different nodes
- Friendly for GPU kernels
- Smaller communication than TP
- Only apply for MoE layers.
 Attention needs DP.
- Different experts have different loads, which needs to be balanced



Sequence and context parallelism (not in vLLM yet)

- Extension of DP to subsequences of a single long sequence
- Communicate Qs or KVs across different shards
- Parallelize the sequence dimension for very long sequences. Balance the KV Cache across shards
- No parameter memory saving
- For sampling one request, only one GPU will be running. Require balancing requests loads.



Pipeline parallelism (in vLLM)

- Distribute the layers to different GPUs and pipeline the execution across devices
- Lowest communication overhead
- Increased latency
- Throughput bottlenecked by the slowest stage



Communication kernel

• Communication kernels have different requirements in inference:





Prefill-Decode Disaggregation

- Different parallelization strategies have different characteristics in computation/memory/communication requirements
- Different stages of LLM inference have different characteristics as well:
 - Prefill: typically compute bound
 - Decode: typically memory bound and require large KV cache memory
- **Disaggregation:** Have different setups for prefill and decode. Transfer the KV cache between them.
- *Side benefit:* Separate concern on time-to-first-token/time-per-output-token

Key Focuses in vLLM



Parallelism

Inference optimizations

Performance engineering

Inference Optimizations

- Optimizing the the model:
 - Quantization

- Optimizing "prefill":
 - Prefix caching, CPU KV cache offloading, etc.

- Optimizing "decode":
 - Speculative decoding, Jump decoding, etc.

Quantization

- Use reduced precisions to store & compute the model
 - BFloat 16 & Float 16 are the standard for "unquantized" models
 - Quantization typically means <u>using 8 or lower bits</u> (e.g., FP8, INT8, FP4)

• Native hardware support makes quantization increasingly effective (e.g., FP8 in Hopper & FP4 in Blackwell)

LLM Quantization

- LLMs have 3 axes to quantize
 - Size: Weights >= KV cache >>> Activation
- 1. Weight quantization (e.g., FP8, INT8, GPTQ, AWQ)
 - Main benefits: Reduced storage & memory footprints
- 2. KV cache quantization (e.g., FP8)
 - Main benefits: Reduced KV cache storage & Faster attention
- 3. Activation quantization (e.g., FP8, INT8)
 - Main benefits: Faster GeMM & communication (in distributed inference)

Quantization support in vLLM

📀 LLM Compressor



• <u>LLM Compressor</u> library for quantizing the model weights

 vLLM provides highly-tuned GPU kernels for quantized ops

Automatic Prefix Caching

Example 1: Shared system prompt

Request A

A chat between a curious user and an artificial intelligence assistant. The assistant gives helpful, detailed, and polite answers to the user's questions. User: Hello!

Request B

A chat between a curious user and an artificial intelligence assistant. The assistant gives helpful, detailed, and polite answers to the user's questions. User: How are you? Shared

Automatic Prefix Caching

Example 2: Multi-round conversation

Prompt (round 1)

Human: What's AI?

LLM Result (round 1) LLM: AI is technology that simulates human intelligence, like Siri or Google Maps.

Shared Human: What's AI? LLM: AI is technology that simulates human intelligence, like Siri or Google Maps. Human: Cool, thanks!

> LLM Result (round 2) LLM: No problem!

Hash-Based Automatic Prefix Caching



Speculative Decoding

Small model writes a draft \rightarrow Large model verifies it

Much faster: verifying 5 tokens takes similar time as generating 1 token


Jump Decoding

"Jump" the token generation using the predefined JSON schema



Reference: Fast JSON Decoding for Local LLMs with Compressed Finite State Machine

Unified representation for scheduling

- The scheduling decision is simply represented as a dictionary of {request_id: num_tokens}
- "token budget" to control the per-step execution time



Unified representation for scheduling (cont'd)

- Unification of "prefill" and "decode"
 - There's no concept of prefill and decode
 - Schedule based on the difference between num_compute_tokens and len(all_token_ids)
- Ex1) "Prefill" & "Decode"



Unified representation for scheduling (cont'd)

- Unification of "prefill" and "decode"
 - There's no concept of prefill and decode
 - Schedule based on the difference between num_compute_tokens and len(all_token_ids)
- Ex2) Chunked prefills



Unified representation for scheduling (cont'd)

- Unification of "prefill" and "decode"
 - There's no concept of prefill and decode
 - Schedule based on the difference between num_compute_tokens and len(all_token_ids)
- Ex3) Prefix caching



Let's combine them all together!

Without optimizations

Optimizations Combined

Prompt	<system> You are a helpful, respectful and honest assistant Structure your response strictly in a given JSON format. <user> Generate a description for this item:</user></system>	Prompt	<system> You are a helpful, respectful and honest assistant Structure your response strictly in a given JSON format. <user> Generate a description for this item:</user></system>
Output		Output	

Key Focuses in vLLM

Models

Parallelism

Inference optimizations

Performance engineering

Performance Engineering

• The biggest lesson we've learned in vLLM

"To fully utilize the GPU, we need to pay close attention to everything happening on the CPU (i.e., <u>CPU overheads</u>)"

• Ex) If you build an inference engine in PyTorch without caring much about CPU overheads, you will likely get 10-20% GPU utilization

CPU Overheads

• CPU overheads in an old version of vLLM



- Why so much overheads?
 - Python is slow
 - We didn't utilize multiple threads/processes efficiently in Python
 - PyTorch has performance pitfalls
 - Continuous batching makes input preparation complicated

Optimized Engine Loop & API Server

Goal: Make sure GPU is not stalled

- By pre-processing
 - E.g., converting JPEG images into input tensors (resizing, cropping, ...)

- By post-processing
 - E.g., de-tokenizing output token IDs into output strings

- By HTTP request handling
 - E.g., streaming outputs to 100s of concurrent users

Optimized Engine Loop & API Server (cont'd)

- Two-process approach
 - Process 0 (Frontend): Pre-/post-processing & API Server
 - Importantly, de-tokenization happens in Process 0
 - **Process 1 (EngineCore):** Schedule & execute the model every step
 - A busy loop that is NOT blocked by Process 0



Incremental Input Preparation (Persistent Batch)

- In V0, input tensors are re-created from scratch at every step
- In V1, input tensors are cached, and we only apply the diff at each step
 - Typically, the diff is minimal because only a few requests join or finish at each step
 - The gain is larger for larger objects like block table



CUDA Graph for Low Latency



Piecewise CUDA Graphs



• V0: Single CUDA graph for the entire model

 Pros: Minimal CPU overheads in model execution

- Cons: Limited flexibility
 - Static shapes are required
 - \circ No CPU operations are allowed
 - \rightarrow Increased development burden

Piecewise CUDA Graphs (cont'd)



Piecewise CUDA Graphs (cont'd)



- V1: Splits the model into pieces
 - Runs the attention op in eager-mode PyTorch
 - Runs other ops with CUDA graphs
 - Easy to capture, since the ops are token-wise
- Pros: Maximum freedom in implementing the attention op
 - No restriction on shapes
 - Any CPU operations are allowed
- Cons: CPU overheads unhidden by CUDA graphs could slow down the model execution
 - Negligible for 8B+ models on H100

Q&A