Infinite-LLM: Efficient LLM Service for Long Context with DistAttention and Distributed KVCache

Bin Lin[†], Tao Peng[†], Hanyu Zhao[†], Wencong Xiao[†], Minmin Sun[†], Zhipeng Zhang[†], Lanbo Li[†], Xiafei Qiu[†], Shen Li[†], Yong Li[†], Wei Lin[†], Chen Zhang[‡], Zhigang Ji[‡], Anmin Liu[◊], Tao Xie[◊]

主讲: 肖同欢 任鑫

† Alibaba Group‡ Shanghai Jiao Tong University♦ Peking University

Trend of model token growth

Trend of model token growth



Model

Trend of model token growth

□ Use BF16, Layer=50, Heads=64, Dimension=128

3500 Infinite-LLM, 196K 3000 2500 Cost 2000 Memoery RingAttention, 1000K 1500 1000 Claude 3, 200K GPT-4 Turbo, 128K 500 Claude 2,10K LLaMA-7B, 4K LLaMA-13B, 8K GPT-3 13B. 2K 0 GPT-3 13B LLaMA-7B LLaMA-13B Claude 2 GPT-4 Turbo Claude 3 RingAttention Infinite-LLM

Trend of model token growth

Trend of GPU memory growth

160

140

□ Maximum VRAM of a single GPU card is 141 GB

Trend of GPU memory growth



141

How to cut LLM cloud service costs?

Optimizations of Attention

D Parallelism Method

Optimizations of Attention

MHA vs GQA vs MQA vs MLA



Flash Attention & Paged Attention

□ Flash Attention: Block-by-Block Sequential Computation

♦ Calculation in HBM: $\mathbf{S} = \mathbf{Q}\mathbf{K}^{\Box} \in \Box^{N \times N}$, $\mathbf{P} = \operatorname{softmax}(\mathbf{S}) \in \Box^{N \times N}$, $\mathbf{O} = \mathbf{P}\mathbf{V} \in \Box^{N \times d}$



Flash Attention & Paged Attention

Paged Attention: Paged Dynamic Loading



How to cut LLM cloud service costs?

Optimizations of Attention

Parallelism Method

- Data Parallelism
 - Data Parallelism
 - Distributed Data Parallelism
- Model Parallelism

Data Parallelism



Distributed Data Parallelism



How to cut LLM cloud service costs?

Optimizations of Attention

Parallelism Method

- Data Parallelism
- Model Parallelism
 - Pipeline Parallelism
 - Tensor Parallelism

D Pipeline Parallelism

Pipeline Parallelism in 4 GPUS, use model with 4T layers



D Pipeline Parallelism

Gpipe: use smaller, equally-sized microbatches to improve efficiency



Pipeline Parallelism

- Gpipe: use smaller, equally-sized microbatches to improve efficiency
- Bubble: idle time slots caused by dependencies
- PipeDream: BWD begins when a microbatch completes its FWD



□ Tensor Parallelism

Parameters (tensors) within a layer are split across different devices



Tensor Parallelism

 $\text{ Row Parallelism } : XA = [X1, X2] \begin{bmatrix} A1 \\ A2 \end{bmatrix} = X1A1 + X2A2 = Y1 + Y2 = Y$



□ Tensor Parallelism

♦ Column Parallelism: XA = [X][A1, A2] = [XA1, XA2] = [Y1, Y2] = Y



Challenges of deploying LLM services

More GPUs support long-text but perform poorly

Dynamic KV Cache Complicates Resource Management

More GPUs support long-text but perform poorly

Traditional parallel methods do not differentiate between attention and nonattention layers, causing performance degradation due to over-segmentation.



More GPUs support long-text but perform poorly

Traditional parallel methods do not differentiate between attention and nonattention layers, causing performance degradation due to over-segmentation.



More GPUs support long-text but perform poorly

Traditional parallel methods do not differentiate between attention and nonattention layers, causing performance degradation due to over-segmentation.



Challenges of deploying LLM services

Over the other set of the other set o

Dynamic KV Cache Complicates Resource Management

Final context length of inference tasks is unknown and highly variable



(a) Many instances are in a state of low computational/memory utilization.

- Final context length of inference tasks is unknown and highly variable
 - Large memory causes waste, while small memory leads to a reduced batch size



(b) Long request competes with short requests, resulting in low batch size and low GPU util.

- Final context length of inference tasks is unknown and highly variable
 - Large memory causes waste, while small memory leads to a reduced batch size



- Final context length of inference tasks is unknown and highly variable
 - Small memory leads to a reduced batch size, lowering GPU computational utilization



(b) Long request competes with short requests, resulting in low batch size and low GPU util.

- Final context length of inference tasks is unknown and highly variable
 - Small memory leads to a reduced batch size, lowering GPU computational utilization

batch size =4



batch size =2





(b) Long request competes with short requests, resulting in low batch size and low GPU util.

- Dynamically allocating appropriate video memory for each task becomes crucial !
 - Large memory causes waste, while small memory leads to a reduced batch size



DistAttention

Cluster-scale Throughput Optimization



DistAttention





DistAttention

Distributed Attention computation



- **Debtors and Creditors**
- **Greedy Algorithm based on Performance model**

Debtors and Creditors

- Local First
- Limited Offload

Debtors and Creditors

- Local First
- Limited Offload



The scale of offloaded KVCache not exceed the local's, for overlapping the local computation with data transfers and remote computation

□ Greedy Algorithm based on Performance model

- Performance model
- Greedy Algorithm

□ Greedy Algorithm based on Performance model

Performance model

$$T^{lyr}(\beta, S) = T^{natn}(\beta) + T^{atn}(S) = -\frac{W(\beta)}{f(\beta)} + \sum_{r=1}^{\beta} \frac{S^r}{g(S)}$$

One Transformer layer's Throughput β is batch size S is length of Sequence

□ Greedy Algorithm based on Performance model

Performance model

$$T^{lyr}(\beta, S) = T^{natn}(\beta) + T^{atn}(S) = -\frac{W(\beta)}{f(\beta)} + \sum_{r=1}^{\beta} \frac{S^r}{g(S)}$$

Non Attention Layer is dictated by batch size Attention Layer is dictated by the requests' length *S*.

□ Greedy Algorithm based on Performance model

Performance model

$$T^{lyr}(\beta, S) = T^{natn}(\beta) + T^{atn}(S) = \left(\frac{W(\beta)}{f(\beta)} + \sum_{r=1}^{\beta} \frac{S^r}{g(S)}\right)$$

 $W(\beta)$: Workload of non-Attention layer S^r : rth Sequence of batch $f(\beta), g(S)$: The GPU's real performance

□ Greedy Algorithm based on Performance model

Performance model

$$\begin{split} T^{lyr}_{dbt}(\beta',K^d) = T^{lyr} &- \frac{K^d}{g(S)} \\ T^{lyr}_{cdt}(\beta,K^c) = T^{lyr} &+ \frac{K^c}{g(S)} \end{split}$$

Affect of KVCache's computation performance

□ Greedy Algorithm based on Performance model

Performance model

$$T_{dbt}^{lyr}(\beta', K^d) = T^{lyr} - \frac{K^d}{g(S)}$$
$$T_{cdt}^{lyr}(\beta, K^c) = T^{lyr} + \frac{K^c}{g(S)}$$

 $TPS_{cluster} = \sum_{i=1}^{M} TPS_i$

Affect of KVCache's computation performance

$$TPS = \frac{\beta}{n \cdot T^{layer}} \quad \text{A LLM instance's Throughput with n layer, batch } \beta$$

A cluster deployed with *M* instances'TPS

□ Greedy Algorithm based on Performance model

- Greedy Algorithm
 - maximize the cluster throughput with DistAttention approximately
 - > pairing overloaded debtor instances with the free creditor

□ Greedy Algorithm based on Performance model



□ Greedy Algorithm based on Performance model



□ Greedy Algorithm based on Performance model

Greedy Algorithm



Each debtor handles multiple requests



□ Greedy Algorithm based on Performance model

Greedy Algorithm



Select debtor's longest request

□ Greedy Algorithm based on Performance model



□ Greedy Algorithm based on Performance model



□ Greedy Algorithm based on Performance model



□ Greedy Algorithm based on Performance model













Given Service Schuld Schulder General Schuld Sch



□ workflow

gManager









Req ID	Length	Instance0		Instance1	
		Block num	Is local	Block num	Is local
0	100k	1562	False	4688	False



Req ID	Length	Instance0		Instance1		
		Block num	Is local	Block num	Is local	k
0	100k	1562	False	4688	False	

Experimental Setup

- ✤ 4 nodes and 32 GPUs. (each node 8 A100, 80G)
- ✤ GPUs are connected via NVLink (600GB/s) within each node
- Ethernet (125MB/s) across nodes
- ✤ Model: LLaMA2(7B, 13B, 70B)

□ Traces

- ✤ 9 traces with different context length ranges and length distributions
- Trace 0
 - open-source dataset ShareGPT4
 - ➤ contains GPT-4 conversation records.
- ✤ Trace 1-2
 - subsets of ShareGPT4
 - > assess performance on different length distribution
- Trace 3-8
 - > opensource dataset L-Eval and these online service
 - assess performance on longer contexts

Comparison

focuses on comparing Infinite-LLM with static model parallelism and resource planning (vLLM signal worker & multi workers).

Infinite-LLM & vLLM multi wokers support multiple instances

Comparison(context length)



Comparison(context length)



Comparison(context length)



Comparison(context length)



For short sequences Infinite-LLM achieves **1.4x to 5.3x throughput improvement.**

For long sequences both frameworks support **similar maximum context lengths**, but Infinite-LLM reduces **OOM risks**.

□ Comparison(end-to-end performance)



□ Comparison(end-to-end performance)



□ Comparison(end-to-end performance)



Summary

□ Background

- Different Context length
- waste of memory and computation resources

Infinite-LLM

- DistAttention
- Greedy Algorithm

□ thinking