# LLM Serving on Heterogeneous Hardware

Mingxing Zhang @ KVCache.AI

https://github.com/kvcache-ai
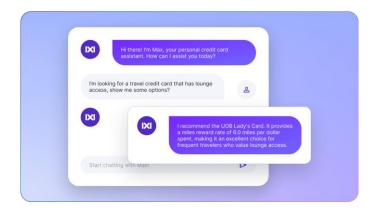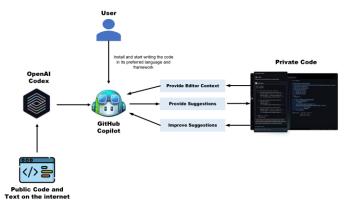
# Background: Large Language Models (LLMs)

Large Language Models (LLMs) are widely applied in industry and researched in academia.
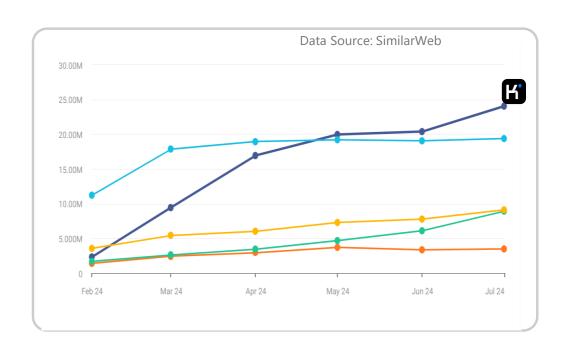


Knowledge Q&A



Content Creation

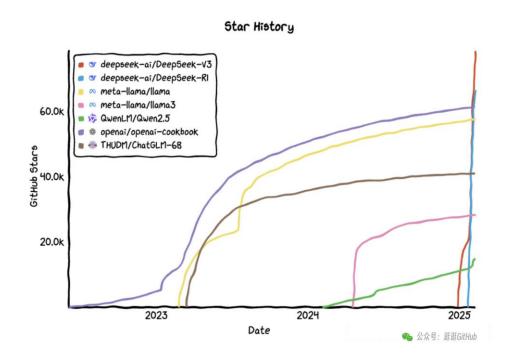

Code Generation



Office Assistant

# Challenge of Online Model as a Service System

More Data + Larger Model + **_Longer Context_** = 🙂 Higher Intelligence



Long input: Moonshot AI's Kimi Supports 2 Million Characters Input in March 2024, become a widely recognized app in China
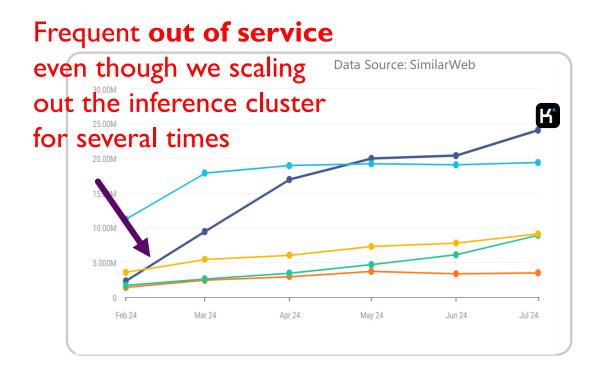


Long output: DeepSeek release V3/R1 at Dec 2024, Become a widely recognized app in global

More Data + Larger Model + ***Longer Context*** = ☹ Higher Service Loads

**Frequent out of service** even though we scaling out the inference cluster for several times

Data Source: SimilarWeb

服务器繁忙，请稍后再试。

**A lot of throttling**

已深度思考（用时 0 秒）

- deepseek-ai/DeepSeek-V3
- deepseek-ai/DeepSeek-R1
- meta-llama/llama
- meta-llama/llama3
- QwenLM/Qwen2.5
- openai/openai-cookbook
- THUDM/ChatGLM-6B

Long input: Moonshot AI's Kimi Supports 2 Million Characters Input in March 2024, become a widely recognized app in China

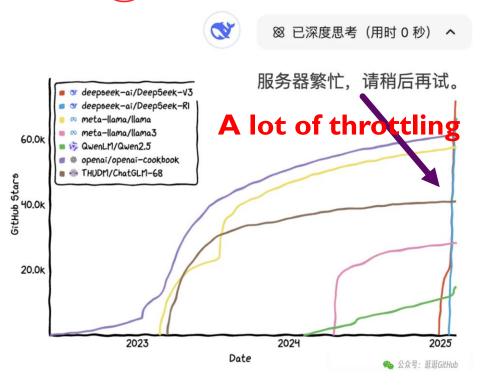Long output: DeepSeek release V3/R1 at Dec 2024, Become a widely recognized app in global

# Content

- Motivation for Heterogeneous LLM Serving

- Core Technologies of Mooncake

- Core Technologies of KTransformers

- Tutorial: Fine-Tune and Chat with Your Customized Model Locally

!!! The price numbers are not accurate, just a demonstration!

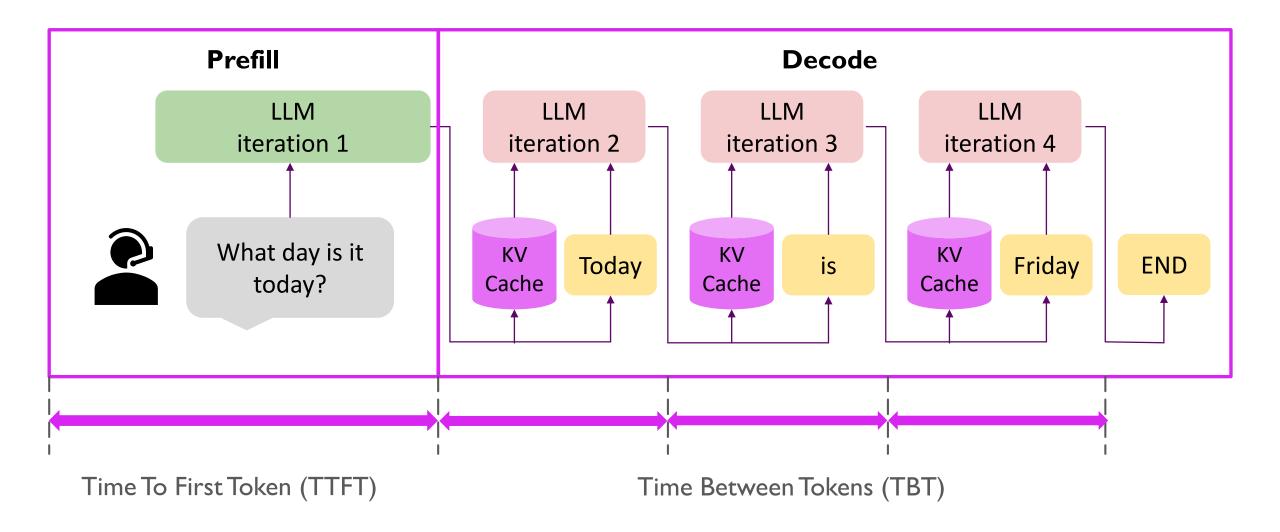|  | H800 | H20 | Xeon SPR + 8 * DDR5-4800 |
|---|---|---|---|
| Hardware Spec | 80GB VRAM, 3.3 TBps<br>~ 1 PFLOPS<br>> $ 10,000 | 96GB VRAM, 4 TBps<br>~ 200 TFLOPS<br>~ $50,000 | 8*64GB DRAM, 8*40GB/s<br>< 20 TFLOPS<br>~ ¥60,000 |
| Best for | Allround,<br>especially for TFLOPS/$ | Bandwidth/$ | Capacity/$ |

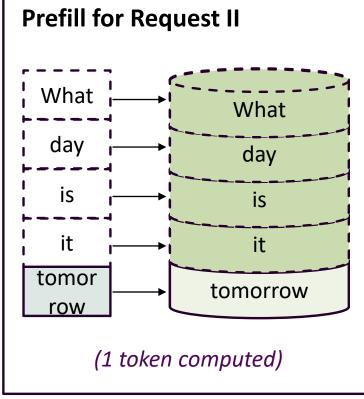**1  Motivation for Heterogeneous LLM Serving**

# LLM Inference

- KVCache can be shared across requests with the same prefix, reducing computation



"What day is it today"

**Prefill for Request I**

| What | → | What |
| day | → | day |
| is | → | is |
| it | → | it |
| today | → | today |

*(5 tokens computed)*

*KVCache Reuse*

"What day is it tomorrow"

**Prefill for Request II**

| What | → | What |
| day | → | day |
| is | → | is |
| it | → | it |
| tomor row | → | tomorrow |

*(1 token computed)*

# Different Hardware are Good at Different Dimension

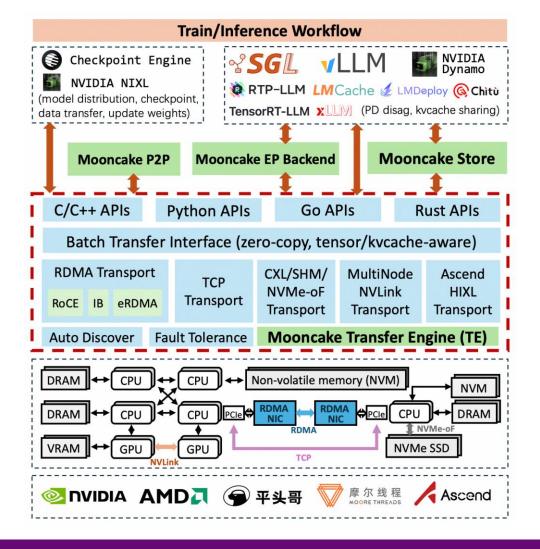| | H800 | H20 | Xeon SPR + 8 * DDR5-4800 |
|---|---|---|---|
| **Hardware Spec** | 80GB VRAM, 3.3 TBps<br>~ 1 PFLOPS<br>> $ 10,000<br><br>**For Prefill!** | 96GB VRAM, 4 TBps<br>~ 200 TFLOPS<br>~ $50,000<br><br>**For Decode!** | 8*64GB DRAM, 8*40GB/s<br>< 20 TFLOPS<br>~ ¥60,000<br><br>**For KVCache!** |
| **Best for** | Allround,<br>especially for TFLOPS/$ | Bandwidth/$ | Capacity/$ |

!!! The price numbers are not accurate, just a demonstration!

A KVCache-centric Disaggregated Architecture for LLM Serving

2  Core Technologies of Mooncake

- ⬛ The serving platform of Kimi

1. P/D disaggregation architecture centered around the distributed KVCache pool

2. Trading more storage of less compute! Increase the throughput of Kimi by 75%
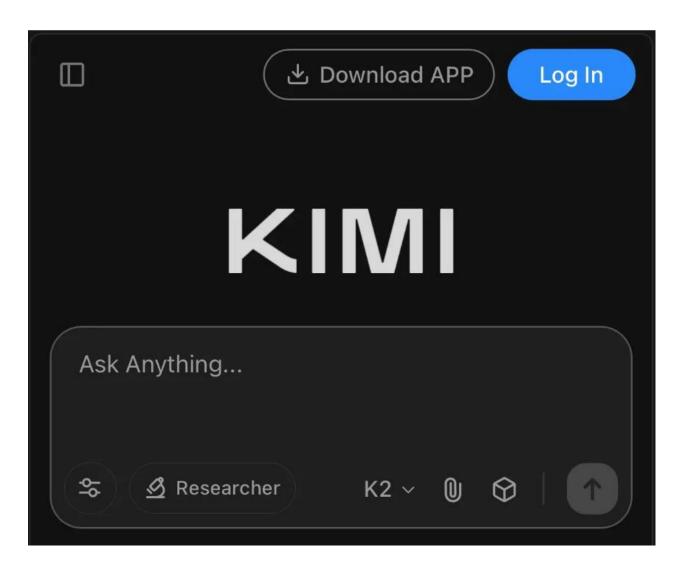
3. Meet SLO guarantee

Moonshot AI  +  KVCache.AI @ Tsinghua

**More：** https://github.com/kvcache-ai/Mooncake

Mooncake (1): 在月之暗面做月饼，Kimi 以 KVCache 为中心的分离式推理架构

# Kimi @ Moonshot AI

# P&D Disaggregated Inference

- Avoid interference between prefill and decoding in a mixed batch

- Decouple resources and parallelism to improve MFU (Model Flops Utilization)



[1] Agrawal et al. Taming Throughput-Latency Tradeoff in LLM Inference with Sarathi-Serve (OSDI'24)

(a) Conversation

(b) Tool Use

- Better SLO control

# P&D Disaggregation Becomes a Necessary

- Prefill and decode needs different parallelism strategy, e.g., DeepSeek V3/R1

# KVCache Cache introduces High Challenges to Storage System

- Each 1 token -> 2 * layers * hidden dimension = tens of KB KVCache

- Not only the size of KVCache is large, it also requires high transfer bandwidth to avoid stall of GPU

**Reusable KVCache**

**(Hundreds of TB ~ PB)**

**KVCache of TB Model (tens of TB)**

TB Model (TB)

100B+ Model (Hundreds of GB)

10B+ Model (GB)

- Cache hit ratio grows proportional to the size of the cache

- □ Different scenarios has different settings

- □ Overall, we need PB-level cache that exceeds to size of a single machine

- ## Key of KVCache Cache：Large size and bandwidth
  Utilize high performance connection like (GPUDirect) RDMA/Storage



Will be open sourced soon!

- Pooled memory as KVCache cache

- Independent to specific inference engine

- Optimized for multi-NIC scenario

**Inference Engine**

| vLLM, and other inference engine | Megatron, training |
|---|---|

↕ Zero-copy Object Put/Get ↕

**Mooncake Store**

| Managed Store (master-slave) | P2P Store (client-only) |
|---|---|

↕ Mooncake Store BatchTransfer API ↕

Mooncake Transfer Engine

↕ Memory Read/Write ↕

**Storage Resource**

Bare-metal Storage Resource

NUMA Node 1  NUMA Node 2

RDMA Remote Direct Memory Access

NVMe OVER FABRICS POWERED BY RoCE

NUMA Node 3  NUMA Node 4

CXL Compute Express Link

nvm EXPRESS

Local mem    Remote mem  (Remote) SSD

Third-party Memory Storage

- Effective request capacity: Number of requests that meet the latency requirements

- Achieve up to a **498%** increase in effective request capacity compared to vLLM, vLLM with prefix caching and with chunked prefill

# Evaluation: GPU Computation Cost

- Cache hit rate: global cache > local cache

- Save **29% - 61%** on GPU computation costs

# Mooncake – Open Sourced and Build with the Community

From flagship applications → To industry-wide adoption

2024.3 Kimi went viral for its long-context capabilities, using Mooncake to handle surging traffic
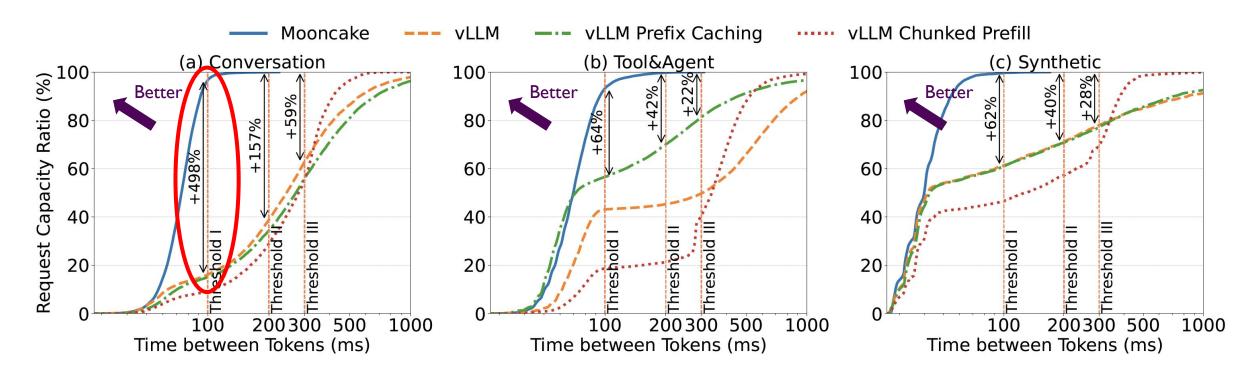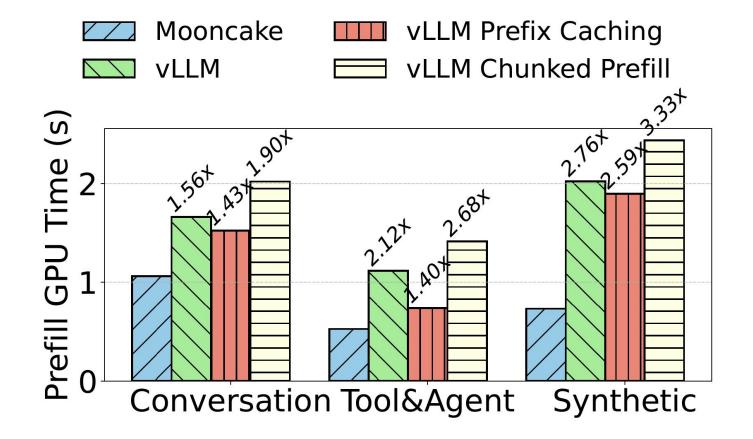
2024.11 Mooncake open-sourced; adopted by Alibaba and Ant Finance

Used in Dynamo, the distributed inference system highlighted at GTC 2025 Keynote

2024.6 Mooncake tech report sparked wide industry discussion

2025.2 USENIX FAST Best Paper Award



**Mooncake**

## USENIX
### THE ADVANCED COMPUTING SYSTEMS ASSOCIATION

**ERIK RIEDEL BEST PAPER AWARD**

presented to

Ruoyu Qin, Zheming Li, Weiran He, Jialei Cui, Feng Ren, Mingxing Zhang, Yongwei Wu, Weimin Zheng, and Xinran Xu

for

Mooncake: Trading More Storage for Less Computation — A KVCache-centric Architecture for Serving LLM Chatbot

Presented at the 23rd USENIX Conference on File and Storage Technologies

Amy Rich
President

2/25/2025
Date

**USENIX FAST2025 Best Paper**

kvcache-ai/Mooncake -- An open-source initiative co-launched by Moonshot AI and Tsinghua University, with collaboration from various large model and infrastructure providers

Moonshot AI    清华 MADSys 实验室    9#AI SOFT    ANT FINANCIAL

HUAWEI    Alibaba Cloud    趋境科技 APPROACHING.AI

VolcanoEngine    面壁智能 MODELBEST    金山云    NVIDIA    and more …

# Mooncake – Adopted/Collaborated with Other Famous Communities

## vLLM



⭐ **41.5K Stars**

- One of the most widely used inference engines, adopted by major cloud providers
- Its distributed inference is built on Mooncake

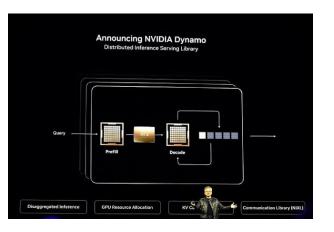**v0.8.3** (Latest)

github-actions released this 2 days ago · 57 commits to main since this release · v0.8.3 · 296c657

### Cluster Scale Serving

- Support XpYd disaggregated prefill with **MooncakeStore** (#12957)

## NVIDIA Dynamo



- Spotlighted by Jensen Huang at GTC 2025 Keynote
- Its architecture is inspired by Mooncake, with explicit acknowledgments

### Acknowledgement

We would like to acknowledge several open source software stacks for motivating us to create Dynamo.

- vLLM and vLLM-project
- SGLang
- DistServe
- **Mooncake**

- *Memory bottlenecks*: Large-scale inference workloads demand extensive capacity. KV cache offloading across memory hierarchies (HBM, DDR, N memory limits and speeds up latency. (**Mooncake**, AIBrix, LMCache)

## SGL

- Inference engine of xAI, widely used in DeepSeek inference
- Distributed architecture was co-developed with Mooncake



**LMSYS Org** ✔ @lmsysorg

The SGLang Team is honored to announce that the following well-known companies and teams, among others, have adopted SGLang for running DeepSeek V3 and R1. @AMD @nvidia @Azure @basetenco @novita_ai_labs @BytedanceTalk @DataCrunch_io @hyperbolic_labs @Vultr @runPod
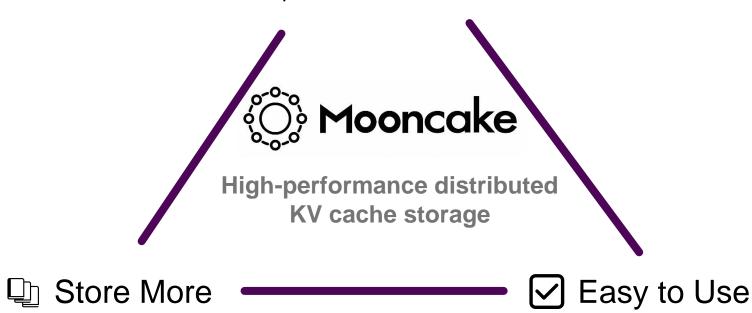
**LMSYS Org** ✔ @lmsysorg

SGLang has achieved a milestone with full support for PD Disaggregation, thanks to the **MoonCake team.** Teng Ma, Shangming Cai, Xuchun Shang and Yuan Luo were instrumental in this achievement. Special thanks to Atlas Cloud for their support with the H100s cluster. Let's go! 🚀

# Key to KVCache

- **Mooncake Transfer Engine**
- **End-to-end zero-copy**

⚡ Transfer Fast

**Mooncake**

**High-performance distributed
KV cache storage**

📄 Store More        ☑ Easy to Use

- **Elastic, Shared, and Multi-layer KV Cache**
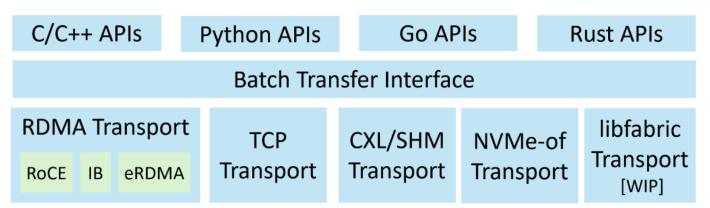- **Memory Allocator Optimized for LLM Inference**

- **Extensive and user-friendly APIs**

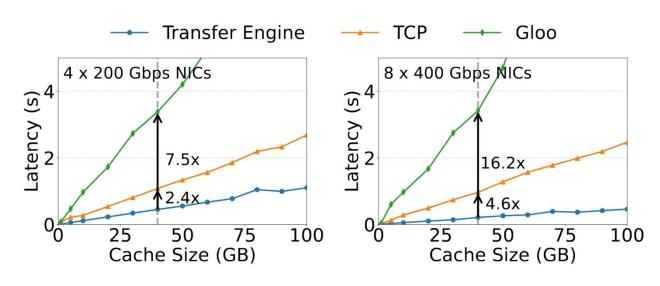# Transfer Fast: Mooncake Transfer Engine

- ## Key features
  - **Topology-aware path selection**
  - **Multi-NIC pooling**
  - **Supports multiple protocols and provides unified interfaces.**
  - **Multi-language APIs**



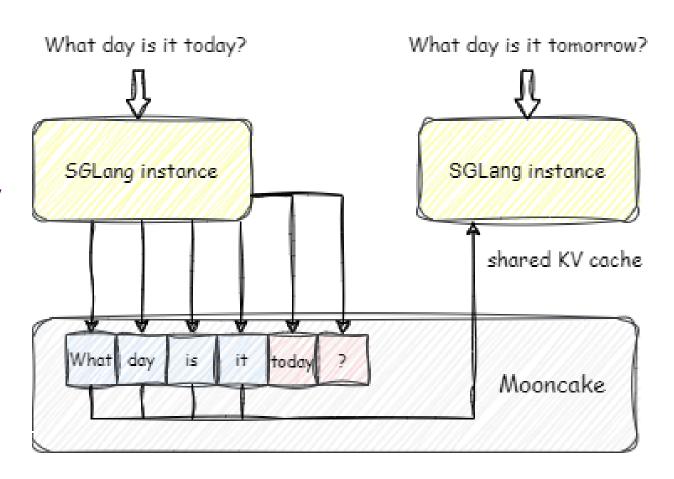Mooncake Transfer Engine



**Lightening fast over RDMA**

- 40 GB KVCache (128k tokens, LLaMA3-70B)
- **87 GB/s @ 4×200 Gbps**, RoCE
- **190 GB/s @ 8×400 Gbps**, RoCE

- **Key features**

  - **Distributed KV cache sharing:** storing one and usable by all

  - **Dynamic resource scaling:** dynamically adding and removing store nodes (startup in <80s for 500GB memory and 8 RDMA NICs)

  - **Multi-layer storage (WIP):** offloading cached data from RAM to SSD

# Extensive APIs, Easy to Use

## Put/Get APIs

- Put/Get single object
- Batch Put/Get
- **(Batch) Zero-copy Put/Get: recommended**
- (Batch and zero-copy) Put/Get from/into multi-parts
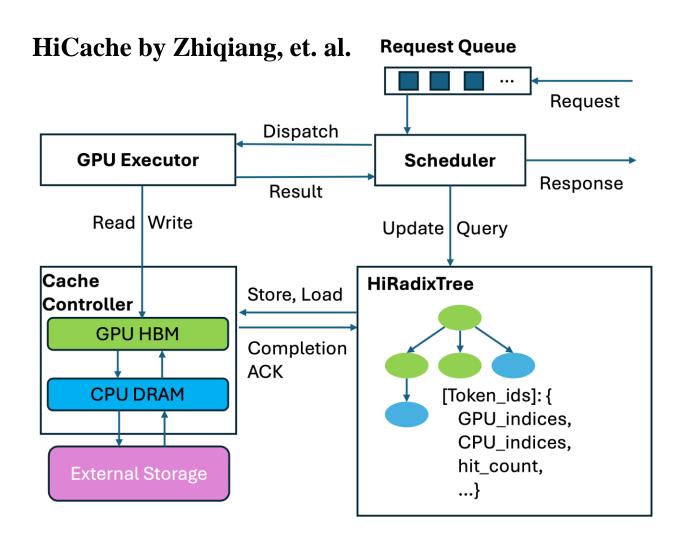
## Configurable KV cache placement

- Replica number
- With soft pin
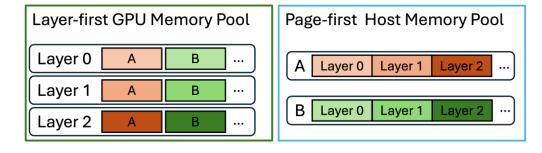- Preferred segment

## Hello world example

```python
from mooncake.store import MooncakeDistributedStore

# 1. Create store instance
store = MooncakeDistributedStore()

# 2. Setup with all required parameters
store.setup(
    "localhost",                # Your node's address
    "http://localhost:8080/metadata",    # HTTP metadata se
    512*1024*1024,              # 512MB segment size
    128*1024*1024,              # 128MB local buffer
    "tcp",                              # Use TCP (RDMA for
    "",                                 # Leave empty; Mooncake
    "localhost:50051"          # Master service
)

# 3. Store data
store.put("hello_key", b"Hello, Mooncake Store!")

# 4. Retrieve data
data = store.get("hello_key")
print(data.decode())  # Output: Hello, Mooncake Store!

# 5. Clean up
store.close()
```

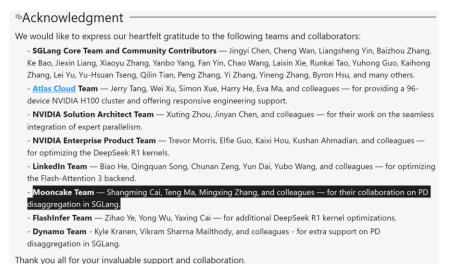**HiCache by Zhiqiang, et. al.**



Page First Layout
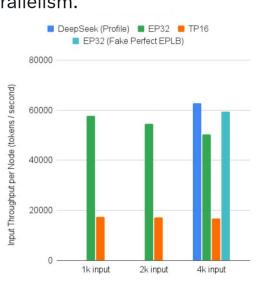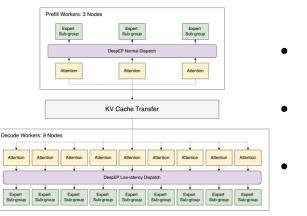
# SGLang + Mooncake

🚀 Breaking: SGLang provides the first open-source implementation to serve @deepseek_ai V3/R1 models with large-scale expert parallelism and prefill-decode disaggregation on 96 GPUs.

**It nearly matches the throughput reported by the official DeepSeek blog**, achieving 52.3K input tokens per second and 22.3K output tokens per second per node. This optimized strategy improves output throughput by up to 5x compared to vanilla tensor parallelism.

Prefill Workers: 3 Nodes / KV Cache Transfer / Decode Workers: 9 Nodes

Decode Performance
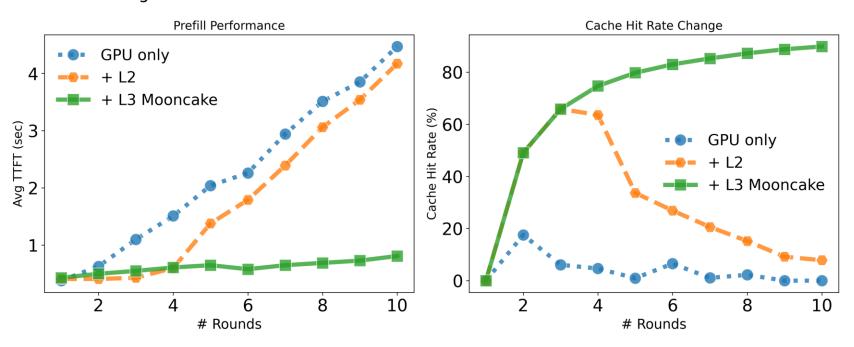
- 3 Prefill + 9 Decode
- DeepEP + EPLB
- Double Batch Overlap

- 52.3k input tps
- 22.3k output tps
- per node

# SGLang + HiCache + Mooncake

## SGLang HiCache with Mooncake Backend on Multi-turn Conversation Benchmark



Prefill Performance

- GPU only
- + L2
- + L3 Mooncake

Cache Hit Rate Change

- GPU only
- + L2
- + L3 Mooncake

Effective KV caching significantly reduces TTFT by eliminating redundant and costly re-computation. **Integrating SGLang HiCache with the Mooncake service** enables scalable KV cache retention and high-performance access. In our evaluation, we tested the DeepSeek-R1-671B model under PD-disaggregated deployment using in-house online requests sampled from a **general QA scenario**. On average, **cache hits achieved an 84% reduction in TTFT compared to full re-computation**.
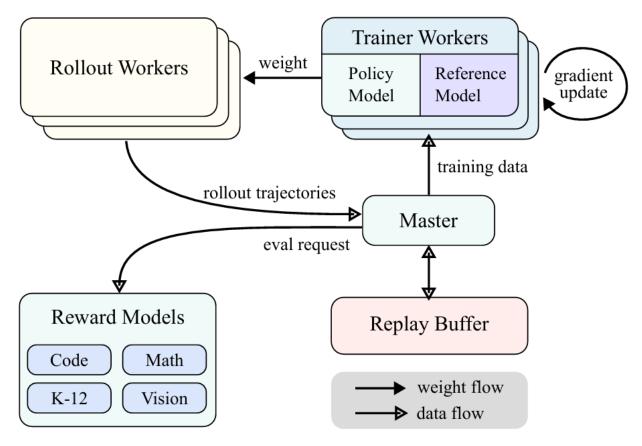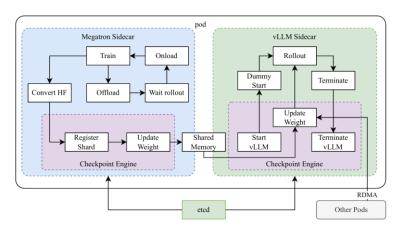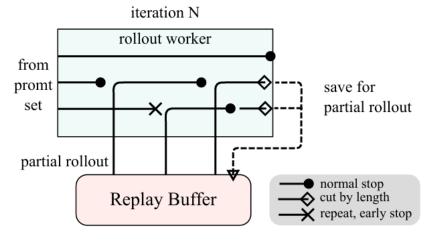
– Ant Group

Thanks :

https://github.com/MoonshotAI/checkpoint-engine/



(a) System overview

Fast Checkpoint Transfer

(b) Partial Rollout

**Train/Inference Workflow**

Checkpoint Engine

NVIDIA NIXL

SGL vLLM

NVIDIA Dynamo

**Mooncake Store**

C/C++ APIs | Python APIs | Go APIs | Rust APIs

Batch Transfer Interface

| RDMA Transport | TCP Transport | CXL/SHM Transport | MultiNode NVLink Transport | Ascend HIXL Transport |
| RoCE IB eRDMA | | | | |

**Mooncake Transfer Engine (TE)**

DRAM ↔ CPU ↔ CPU ↔ Non-volatile memory (NVM) → NVM

DRAM ↔ CPU ↔ CPU — PCIe ↔ RDMA NIC ↔ RDMA NIC — PCIe ↔ CPU → DRAM

RDMA

VRAM ↔ GPU ↔ GPU

NVLink

NVMe SSD

mooncake

- Transfer Engine as the core
- Disaggregated LLM Serving
  - Reinforcement Learning

Moonshot AI | 清华 MADSys 实验室

9#AI SOFT | HUAWEI | 阿里云

火山引擎 | 蚂蚁集团 ANT GROUP | 趋境科技 APPROACHING.AI

金山云 | 面壁智能 MODELBEST | NVIDIA

and more ...

- Multiple paths coexist within the same cluster
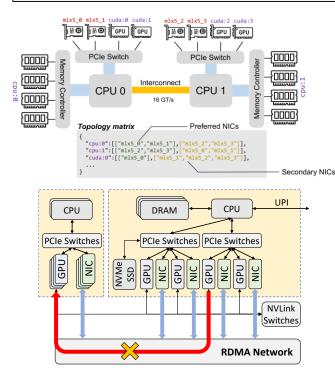
```
auto engine = new TransferEngine();
engine.installTransport("rdma", args);

auto id = engine.allocateBatch();
engine.submitTransfer(id, reqs);
while (true) engine.getStatus(id, st);
engine.freeBatch(id);
```
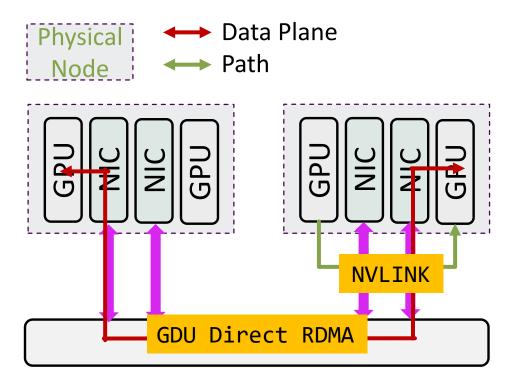


- **The Imperative Path Selection Paradigm**

◎ made static binding decisions once **at startup**

◎ executed a **fixed, state-blind** path scheduling policy

◎ **executed fragilely**, and lacks mechanisms to detect & bypass unavailable paths

- ## Static Binding

  - ### Creates communication silos



Physical Node

Data Plane Path

NVLINK

GDU Direct RDMA

**Different workloads, different transports**

With GPU-Direct

GDU Direct RDMA

Without GPU-Direct (e.g., GTX series)

???

**Different hardware, different transports**

# State-Blind Scheduling

- Increases latency and wastes bandwidth

☐ Slices from this request
☐ Slices from concurrent requests

**A transfer request with 2560 KB**
40 slices in total (each 64KB)

**Topology matrix snippets:**
"cpu:0": {[NIC0, NIC1, NIC2, NIC3], […]}



Transfer Latency

- ## Fragile Execution
  - ### Requires manual intervention and heavy troubleshooting

Process 1

Process 2

NIC  NIC  NIC  NIC

NIC  NIC  NIC  NIC

RDMA Switches

➢ What if a single RDMA fabric failed?

➢ What if a single process crashed?

➢ What if the RDMA switch failed?

■ Goal: Make all transports first-class citizens

| vLLM | SGLang | Dynamo | LMCache | Checkpoint Engine | NIXL | Mooncake Store | ... |

**Transfer Engine APIs**

**Orchestration**

Discover  Routing  Staging

Plugin Loaders  Error Handling

**Unified Segment Abstraction**

Buffers  Devices

Topology

**Transport Plugins** (load all if possible)

RDMA (with GDR/eRDMA)  NVLink  MNNVL  TCP  RoCM

SHM/CXL  GDS  io_uring  HIXL  ...

**Control Plane Plugins**

P2P  etcd

Redis  HTTP

**Heterogeneous GPU Interconnects**

## Dynamic Orchestration

◎ Unified Segment Abstraction

◎ Application-Oblivious Topology Discovery

◎ Dynamic Per-Request Orchestration

## Adaptive Slice Spraying

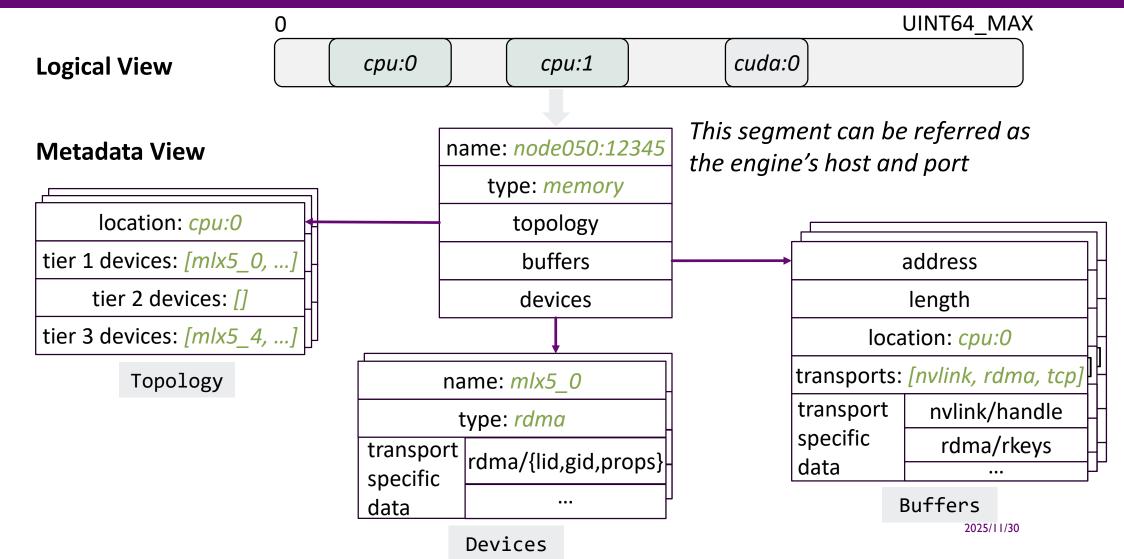■ Latency Prediction based NIC Selection

◎ Cross-Process Fairness

## Resilient Self-Healing

◎ Link-Level Resilience

◎ Transport-Level Resilience

**Logical View**

0                                                                    UINT64_MAX

*cpu:0*     *cpu:1*     *cuda:0*

**Metadata View**

*This segment can be referred as the engine's host and port*

name: *node050:12345*

type: *memory*

topology

buffers

devices

location: *cpu:0*

tier 1 devices: *[mlx5_0, ...]*

tier 2 devices: *[]*

tier 3 devices: *[mlx5_4, ...]*

Topology

name: *mlx5_0*

type: *rdma*

| transport specific data | rdma/{lid,gid,props} |
| | ... |

Devices

address

length

location: *cpu:0*

transports: *[nvlink, rdma, tcp]*

| transport specific data | nvlink/handle |
| | rdma/rkeys |
| | ... |

Buffers

- Step 1: Probe hardware information
  - List of memory/NIC devices
  - Their NUMA affinity, PCIe Bus ID
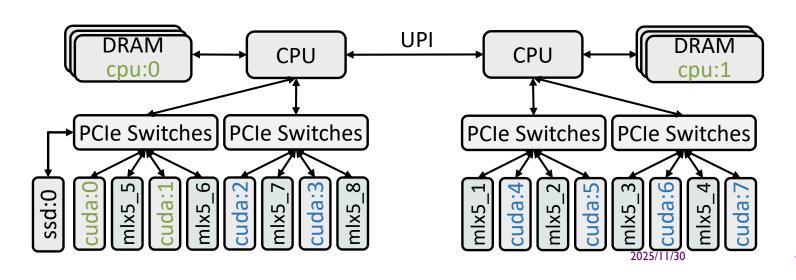  - Capabilities: bandwidth, direct-access, etc.
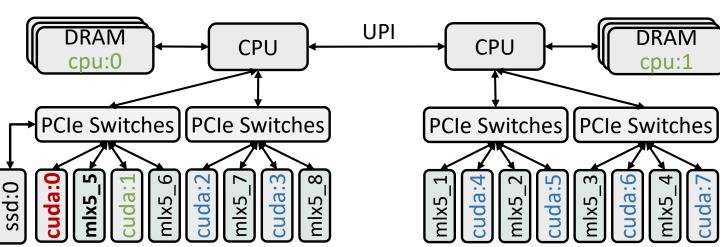
```
MEM:
cpu:[0-1],
cuda:[0-7]
NIC:
mlx5_[1-8]
```



2025/11/30                                                                41

- # Step 1: Probe hardware information

- # Step 2: Maps NICs for each MEM

  - ## Tier 1: NIC(s) with the shortest PCIe hop
  - ## Tier 2: same NUMA but not in tier 1
  - ## Tier 3: cross NUMA

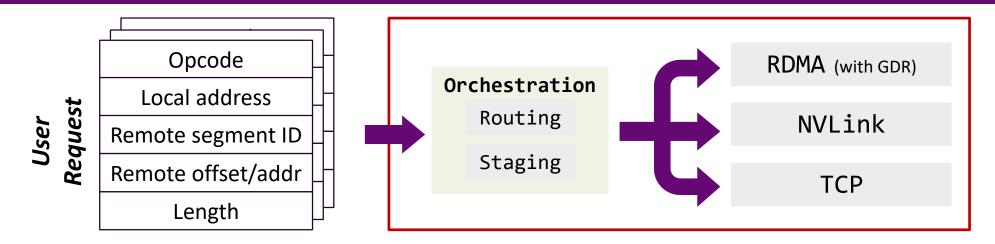- Step 1: Probe hardware information

- Step 2: Maps NICs for each MEM

- Step 3: Load transports
    - Runtime support and transports can be dynamic libraries
    - They can be loaded on runtime (e.g., if CUDA is enabled)

# Dynamic Per-Request Orchestration



**User Request**

| Opcode |
|---|
| Local address |
| Remote segment ID |
| Remote offset/addr |
| Length |

**Orchestration**
- Routing
- Staging

- RDMA (with GDR)
- NVLink
- TCP

- Decide transports for each request using <u>the Unified Segment</u>

  - **Local address**

    find local MEM type

    find local installed transports

  ◎ Remote segment ID & offset/address

    find remote MEM type

    find remote installed transports

Prefer to use a transport with the <u>highest</u>
speed, and supported by <u>both</u> sides

**Dynamic Orchestration**

◎ Unified Segment Abstraction

◎ Application-Oblivious Topology Discovery

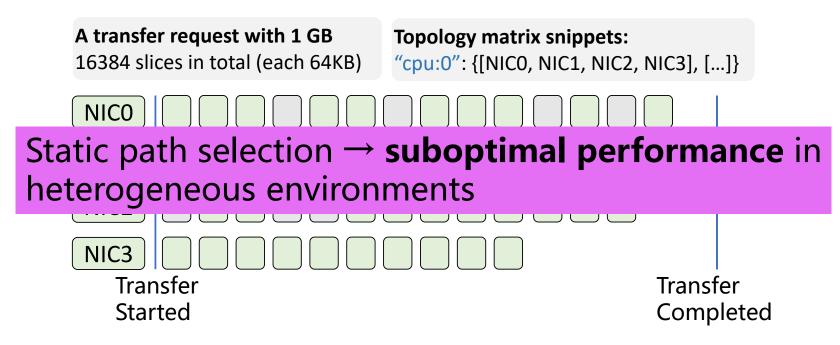◎ Dynamic Per-Request Orchestration

**Adaptive Slice Spraying**

▪ Latency Prediction based NIC Selection

◎ Cross-Process Fairness

**Resilient Self-Healing**

◎ Link-Level Resilience

◎ Transport-Level Resilience

A transfer request with 1 GB
16384 slices in total (each 64KB)

Topology matrix snippets:
"cpu:0": {[NIC0, NIC1, NIC2, NIC3], [...]}

NIC0

Static path selection → **suboptimal performance** in heterogeneous environments

NIC3

Transfer
Started

Transfer
Completed

- How to reduce transfer latency and avoid bandwidth waste

  - **Predict**: select local-remote NIC pair based on historical telemetry

  - **Feedback**: use measured latency to update prediction parameters

# Local NIC Selection

- Latency estimation

- $L_{pred}(NIC_k) =$
  $$\beta_{1,k} \frac{IF_k + P_k * S}{BW_k} + \beta_{0,k}$$

  - $IF_k$ : NIC inflight bytes

  - $S$ : Slice length

  - $P_k$ : Penalty factor (e.g., higher for cross-NUMA)

  - $BW_k$ : NIC bandwidth

  - $\beta_{0,k}, \beta_{1,k}$ : Prediction parameters

## Pre-transfer

- ◎ Calculate $L_{pred}$
- ◎ Find the best NIC
- ◎ Reserve inflight bytes $IF_k$

## Post-transfer

- ◎ Return inflight bytes $IF_k$
- ◎ Measure latency
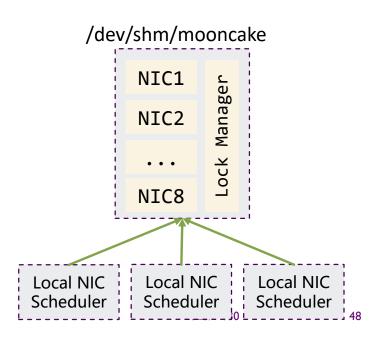- ◎ Update $\beta_{0,k}, \beta_{1,k}$ using EWMA

(make estimation more accurate)

# Cross-Process Fairness

- How to avoid any single process from saturating NICs?

  - Coarse-grained quota allocation

- Decentralized, shared-memory implementation

  - Not every scheduling task enters the global level
    Interval: ~tens of milliseconds

  - Tolerant shutdown/crashes of any process

/dev/shm/mooncake

## Dynamic Orchestration

◎ Unified Segment Abstraction

◎ Application-Oblivious Topology Discovery

◎ Dynamic Per-Request Orchestration

## Adaptive Slice Spraying

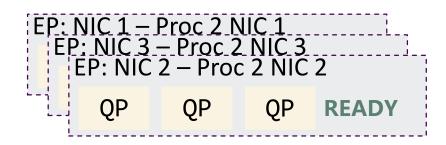■ Latency Prediction based NIC Selection

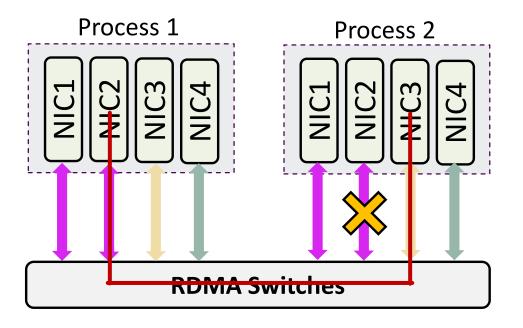◎ Cross-Process Fairness

## Resilient Self-Healing

◎ Link-Level Resilience

◎ Transport-Level Resilience

- ## RDMA Resource Lifecycle
  - ### Endpoint == NIC-to-NIC connection

EP: NIC 1 – Proc 2 NIC 1
EP: NIC 3 – Proc 2 NIC 3
EP: NIC 2 – Proc 2 NIC 2

| QP | QP | QP | **READY** |



**Link-Level Resilience**

◎ Detect failed/unstable link: PAUSE, CLOSE or TERMINATE
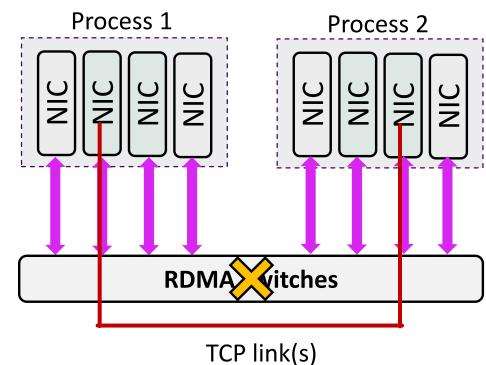
◎ Allow suboptimal path

# Proactive Dual-Layer Resilience

- ## Transport-Level Resilience
  - Transparent fallback to other transports
    (e.g., RDMA→TCP)
  - Driven by Dynamic Per-Request Orchestration

- ## Recovery
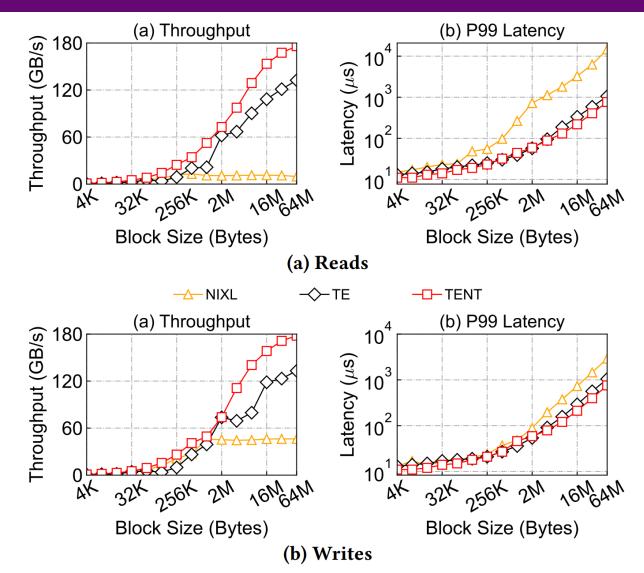  - Transport/path will be reused after link recovery

- Test Cluster: NVIDIA H800 Platform
  - Each node is equipped with two Intel Xeon Platinum 8468V CPUs and eight NVIDIA H800 GPUs
  - NVLink & 200 Gbps $\times$ 8 RoCE interconnect

- **Competitors**
  - Mooncake TENT
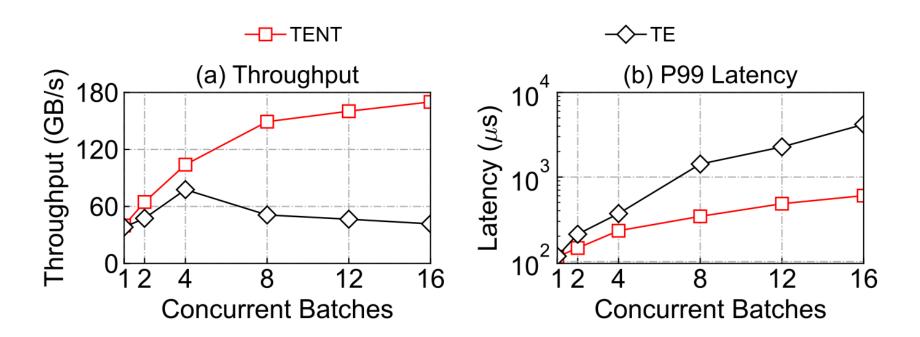  - Mooncake TE (mainstream version)
  - NVIDIA NIXL (UCX backend)

- ## Synthetic Workload

  - ### Block size range from 4KB to 64MB

  - ### Two concurrent threads, 8 NICs are fully utilized



(a) Throughput

(b) P99 Latency

**(a) Reads**

NIXL — TE — TENT

(a) Throughput

(b) P99 Latency

**(b) Writes**

- KVCache Transfer Benchmark

  - **Workload:** DeepSeek-R1-W8A8 model with a 4K input

  - Comprises 61 layers, each containing 32 blocks of 144 KB, consisting of a 128 KB NoPE block and a 16 KB RoPE block
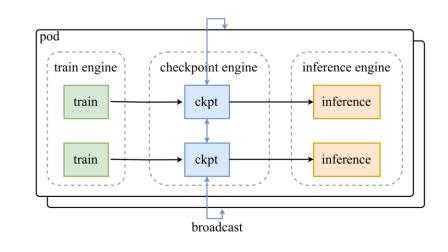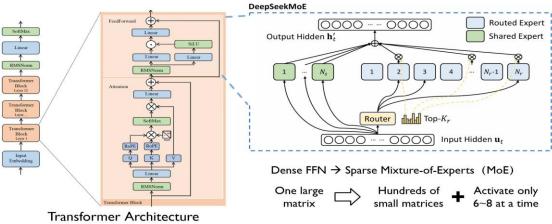
**Moonshot AI**



- **Case Study: Moonshot AI Checkpoint Engine**

  - Open source with Kimi K2

  - Update model weights in LLM inference engines

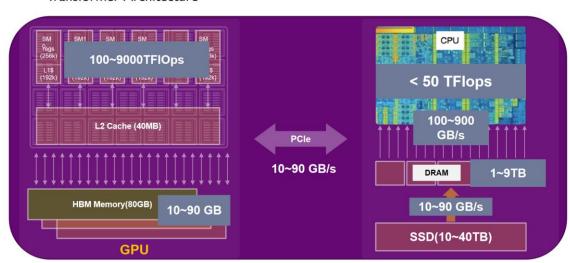  - Update time $\propto$ transfer latency

| Model | TE | TENT |
|---|---|---|
| Qwen3-235B-A22B-Instruct-2507 | 28.56 | 18.97 |
| GLM-4.5-Air | 14.75 | 11.26 |

Attention + MoE

GPU + CPU

Background and Observation of LLM and Sparse Mixture-of-Experts (MoE)

# Background: Sparsification Trends in LLMs



2 out of top 10 open-source models are MoE

All top 10 open-source models are MoE

Transformer Architecture

**DeepSeekMoE**

Output Hidden $\mathbf{h}'_t$

Input Hidden $\mathbf{u}_t$

Router    Top-$K_r$

Routed Expert

Shared Expert

Dense FFN → Sparse Mixture-of-Experts（MoE）

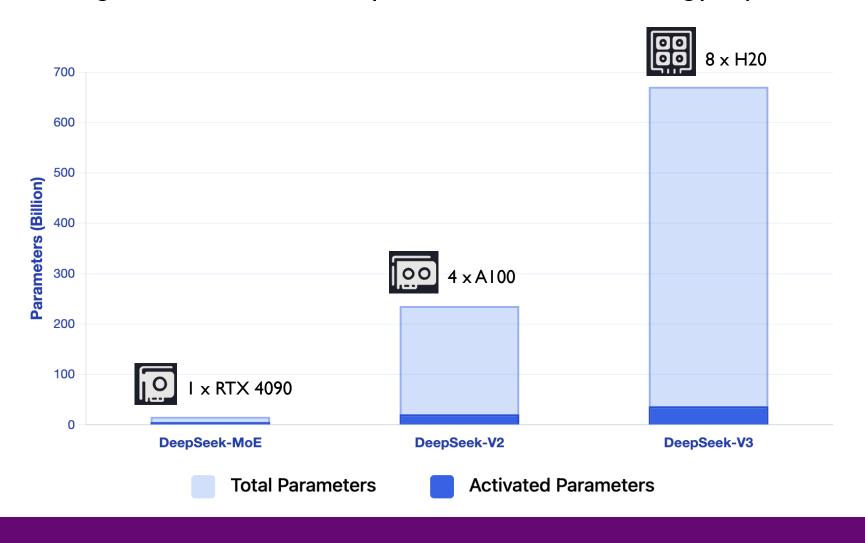One large matrix ⟹ Hundreds of small matrices **+** Activate only 6~8 at a time
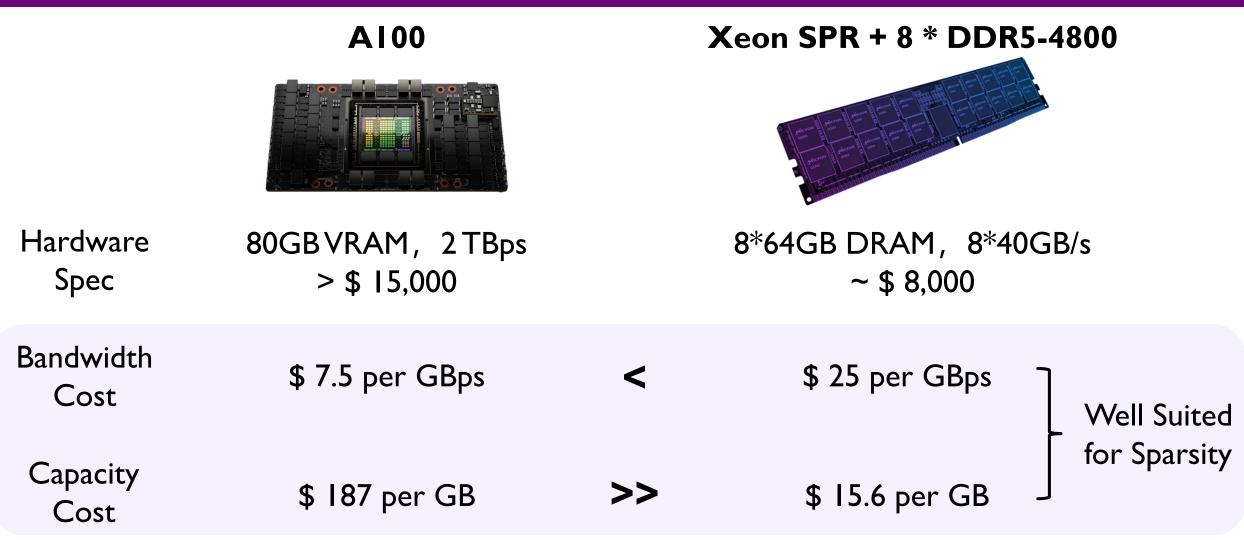
# Background: New Challenge in local deployment

As model sizes grow, traditional GPU-only solutions demand increasingly expensive hardware.

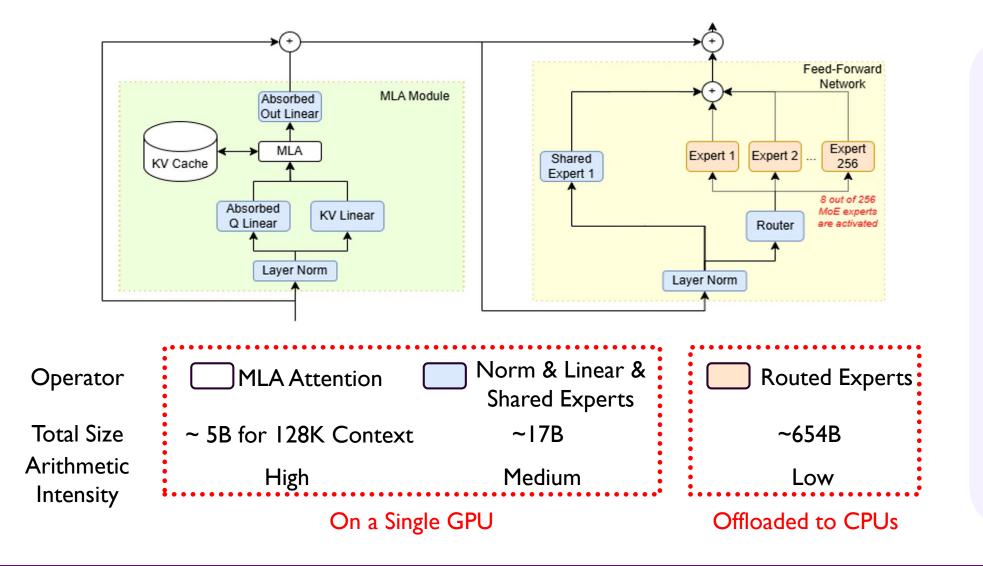# Observation: CPU DRAM is More Suitable for Sparse Models

**A100**                    **Xeon SPR + 8 * DDR5-4800**

| | **A100** | | **Xeon SPR + 8 * DDR5-4800** | |
|---|---|---|---|---|
| Hardware Spec | 80GB VRAM,  2 TBps<br>> $ 15,000 | | 8*64GB DRAM,  8*40GB/s<br>~ $ 8,000 | |
| Bandwidth Cost | $ 7.5 per GBps | < | $ 25 per GBps | Well Suited for Sparsity |
| Capacity Cost | $ 187 per GB | >> | $ 15.6 per GB | |

The price numbers are not accurate, just a demonstration!

Overall KT-System and Optimize in Prefill & Decode

4  Core Technologies of KTransformers

# KTransformers: Challenges and Key Solutions

**Prefill**

**Decode**

**Challenges**

CPU is the Bottleneck for

Intense Computation

Latency of CPU/GPU Coordination

Poor CPU/GPU Overlap

**Solutions**

Advanced CPU Instructions:

Intel AMX

CUDA Graph

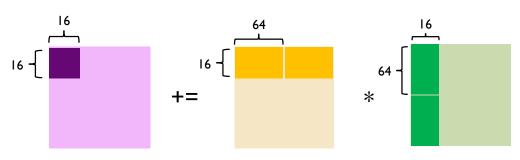Numa-aware Tensor Parallel

Expert Deferral

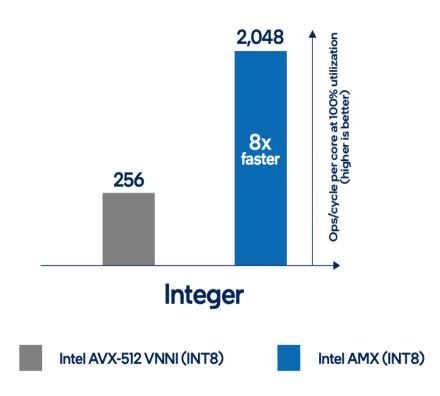**How AVX-512** solves INT8 matrix multiplication problems



128OPS/cycle/FMA. 256OPS/cycle/core

**How AMX** solves INT8 matrix multiplication problems


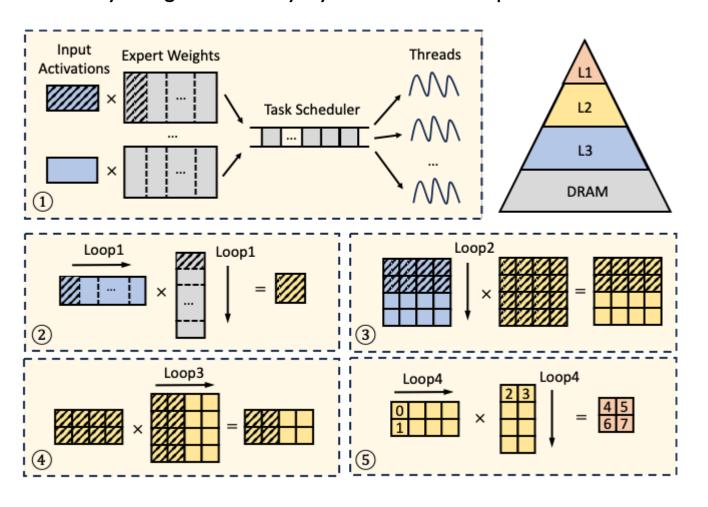
32768OPS/16cycle/core. 2048OPS/cycle/core

**AMX** is 8x faster than AVX-512



2,048

8x faster

256

Integer

Ops/cycle per core at 100% utilization (higher is better)
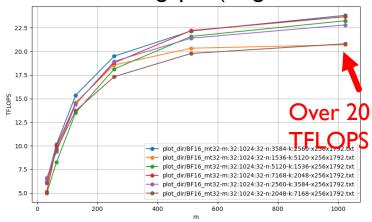
Intel AVX-512 VNNI (INT8)    Intel AMX (INT8)

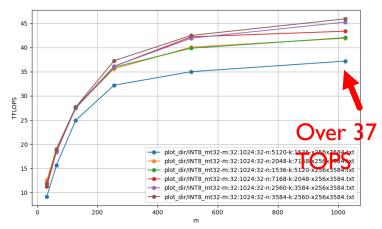# Prefill: AMX Tiling-aware GEMM Kernel

Carefully designed memory layouts and cache-optimized kernels.



BF16 GEMM Throughput (Single Xeon4 CPU).



Over 20 TFLOPS
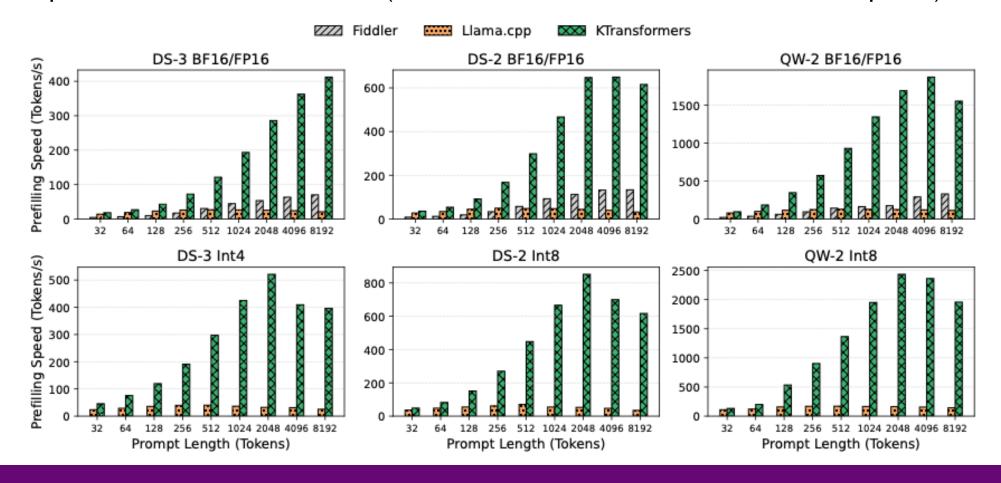
INT8/INT4 GEMM Throughput (Single Xeon4 CPU).



Over 37 TOPS

Up to 19.74✕ faster than Llama.cpp (which does not use AMX kernel)

Up to 5.88✕ faster than Fiddler (which uses Torch's native AMX kernel, sub-optimal)

# Decode: CUDA Graph

**Challenge: Inefficient CPU-GPU coordination**

Fiddler/Llama.cpp forward (a single token) requires **~7000/3000** CUDA kernels, with launch time taking **73%/21%** of total.



Fiddler / Llama.cpp

- Kernel Execution Time
- Average Kernel Launch Time

Only 8% of kernels exceed average launch time

75% of kernels exceed average launch time



**Solution: CUDA Graph**

Capture the **full forward** in a CUDA Graph to remove launch overhead, while carefully avoiding CPU-based operations that introduce breakpoints.

# Decode: Numa-aware Tensor Parallel



CPU architecture

**Challenge: Inefficient CPU-CPU coordination**

Modern systems span multiple NUMA nodes, **cross-NUMA** memory access has worse **latency/bandwidth**.



Local Weights/States
Remote Weights/States

**Solution: Numa-aware Tensor Parallel**

Place expert weight slices in the **local memory** of each NUMA node so that memory access is mostly local, avoiding expensive cross-NUMA memory traffic.

# Decode: Expert Deferral Mechanism



CPU busy

CPU idle

CPU busy

CPU and GPU work alternately

CPU process experts a, b

CPU continue with experts c, d

CPU process experts a, b

CPU and GPU work concurrently

# Decode: Determining the Number of Deferred Experts

## Concern 1: Decoding Speedup



## Concern 2: Model Accuracy Drop



**Balanced Configuration:**

defer as few experts as needed to **saturate the CPU**, while keep at least 2 non-deferred experts per layer to **protect model accuracy**.

Full-accuracy implementation is up to 1.92× faster than Llama.cpp and up to 4.09× faster than Fiddler.

Expert Deferral provides up to 1.45× additional speedups.

# Open Source: KTransformers High-performance Heterogeneous Inference System

Exploratory Open-Source Framework → Widely Used

Jul. 2024. First open release. DeepSeek-V2 with Single GPU + 136GB DRAM

Feb. 2025. DeepSeek-V3/R1 with Single GPU + 382GB DRAM
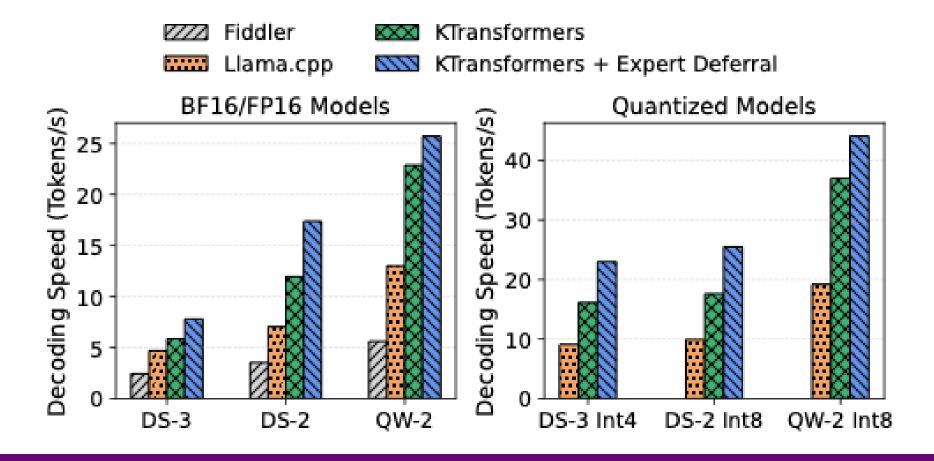
May. 2025. Release AMX-based CPU kernel.

Future. Integrating more features. Supporting more hardware and models.

Aug. 2024. Support 1M-level long context.

Apr. 2025. Support multiple batch size.

Oct. 2025. Integrating into SGL

(a) Flexible Framework

(b) Top 0.01% on Github

(c) Various models and hardware supported

GitHub

https://qwenlm.github.io › blog › qwen3   · 翻译此页   ⋮

## Qwen3: Think Deeper, Act Faster | Qwen

2025年4月29日 — For local usage, tools such as Ollama, LMStudio, MLX, llama.cpp, and KTransformers are highly recommended. These options ensure that users ...

⬤ moonshotai / **Kimi-K2-Thinking** ⧉

**KTransformers Deployment**

**KTransformers+SGLang Inference Deployment**

Launch with KTransformers + SGLang for CPU+GPU heterogeneous inference:

Transformer Architecture

Dense Attention → Sparse Attention

Full Attention ⟹ Hundreds of small chunks **+** Scan only a few at a time

**K** Kimi Mixture of Block Attention（MoBA）

**MoBA: Mixture of Block Attention for Long-Context LLMs**

🐋 deepseek

**Native Sparse Attention: Hardware-Aligned and Natively Trainable Sparse Attention**

# More Open Source Integration

**FineTuning** – Integrated into LLaMA-Factory for local fine-tuning

[Roadmap] Integration of KTransformers as a LoRA Fine-Tuning Backend for LLaMA-Factory #9266  https://github.com/hiyouga/LLaMA-Factory/issues/9266

**Inference** – Integrated into SGLang for wider model support and multi-GPU acceleration

[Feature] KTransformers Integration to Support CPU/GPU Hybrid Inference for MoE Models #11425  https://github.com/sgl-project/sglang/issues/11425

You will be able to fine-tuning and inference 671B DeepSeek and 1TB Kimi K2 locally with consumer GPUs + server CPUs!

# Thanks!

## kvcache.ai

KVCache.AI is a joint research project between MADSys and top industry collaborators, focusing on efficient LLM serving.

👥 **903** followers 🔗 https://madsys.cs.tsinghua.edu.cn/ ✉ zhang_mingxing@mail.tsinghua.edu.cn

---

## Pinned

Customize pins

### 📖 Mooncake  [Public]

Mooncake is the serving platform for Kimi, a leading LLM service provided by Moonshot AI.

🔴 C++  ⭐ 4.1k  🔱 388

### 📖 ktransformers  [Public]

A Flexible Framework for Experiencing Cutting-edge LLM Inference Optimizations

🔵 Python  ⭐ 15.2k  🔱 1.1k

### 📖 TrEnv-X  [Public]

🔵 Go  ⭐ 58

https://github.com/kvcache-ai