Tsinghua University

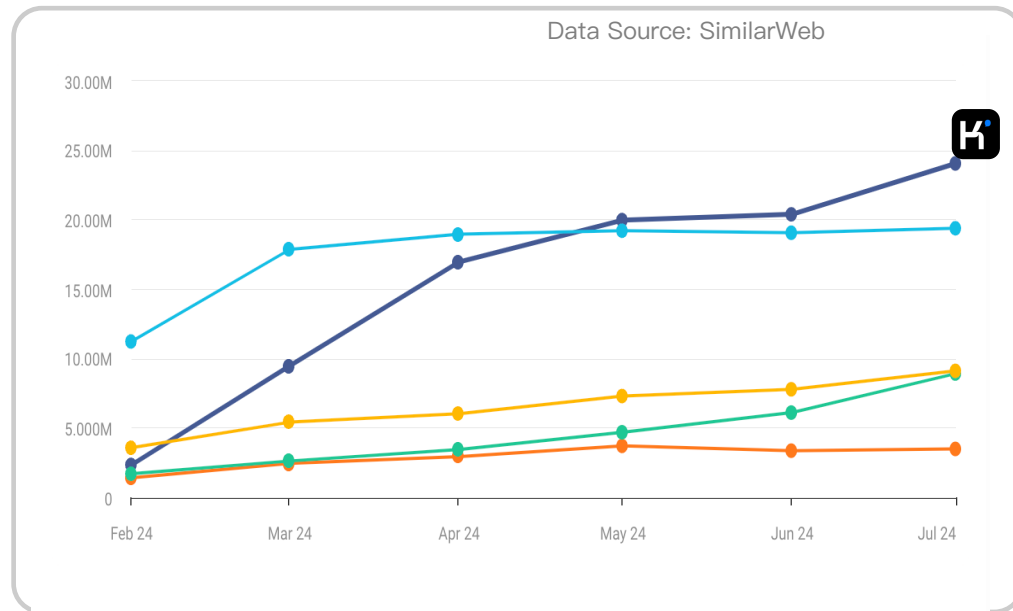# KTransformers: Unleashing the Full Potential of CPU/GPU Hybrid Inference for MoE Models

Mingxing Zhang @ KVCache.AI

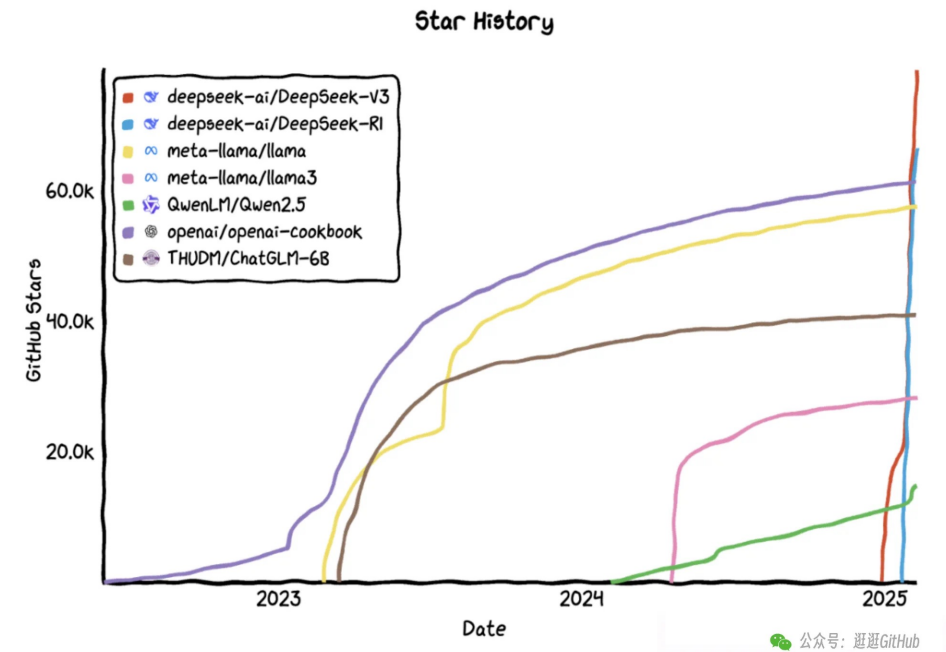https://github.com/kvcache-ai

More Data + Larger Model + ***Longer Context*** = ☺ Higher Intelligence



Data Source: SimilarWeb



Star History

Long input: Moonshot AI's Kimi Supports 2 Million Characters Input in March 2024, become a widely recognized app in China
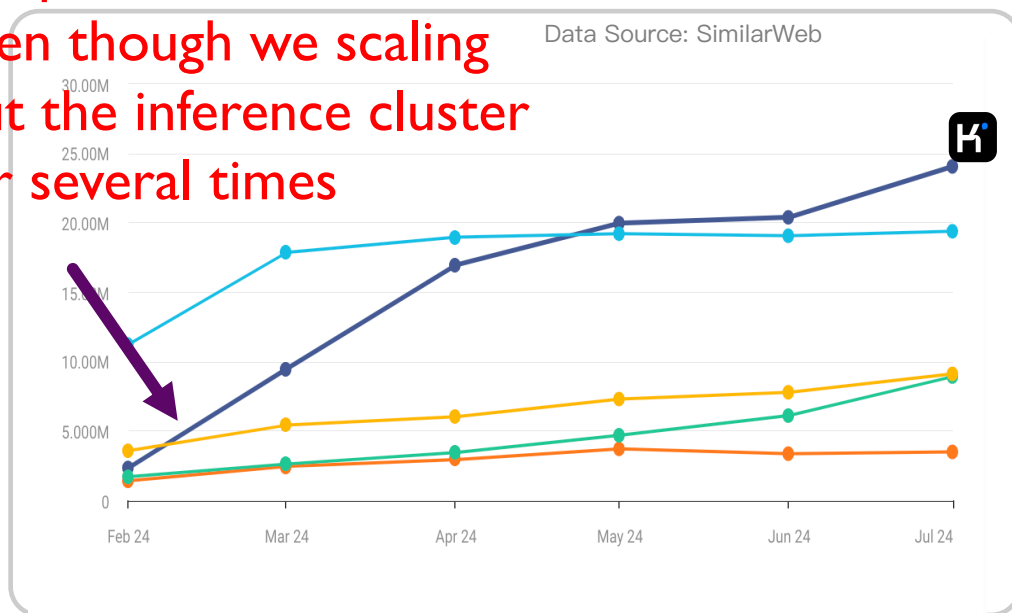
Long output: DeepSeek release V3/R1 at Dec 2024, Become a widely recognized app in global
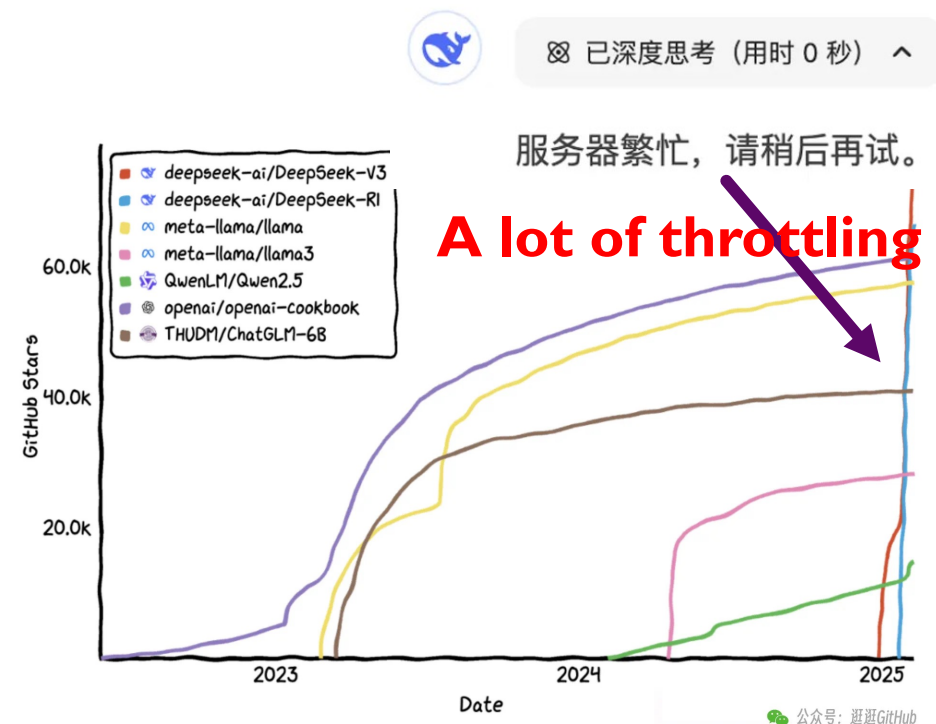
# Challenge of Online Model as a Service System

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

**A lot of throttling**

已深度思考（用时 0 秒）

服务器繁忙，请稍后再试。

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

GitHub Stars axis: 20.0k, 40.0k, 60.0k
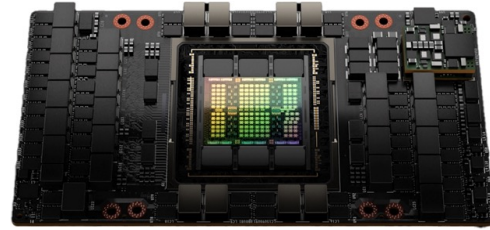Date axis: 2023, 2024, 2025

公众号: 逛逛GitHub

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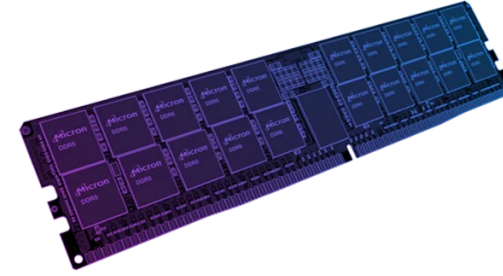Input in March 2024, become a widely recognized app in China

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

# Different Hardware are Good at Different Dimension



H800

Xeon SPR + 8 * DDR5-4800

**Hardware Spec**

80GB VRAM， 3.3 TBps
~ 1 PFLOPS
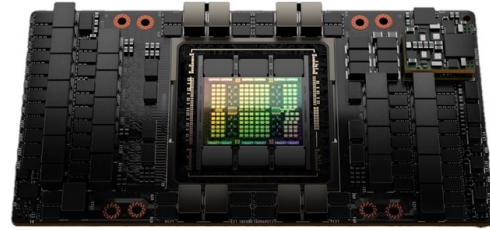> $ 10,000

8*64GB DRAM， 8*40GB/s
< 20 TFLOPS
~ ¥60,000

**Good for Capacity, bad for Bandwidth and Compute Which part is more suitable?**

**Best for**

Allround,
especially for TFLOPS/$

Capacity/$

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

# Different Hardware are Good at Different Dimension

**H800**

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

Hardware
Spec

80GB VRAM，3.3 TBps
~ 1 PFLOPS
> $ 10,000

8*64GB DRAM，8*40GB/s
< 20 TFLOPS
~ ¥60,000

Best
for

Allround,
especially for TFLOPS/$

Capacity/$
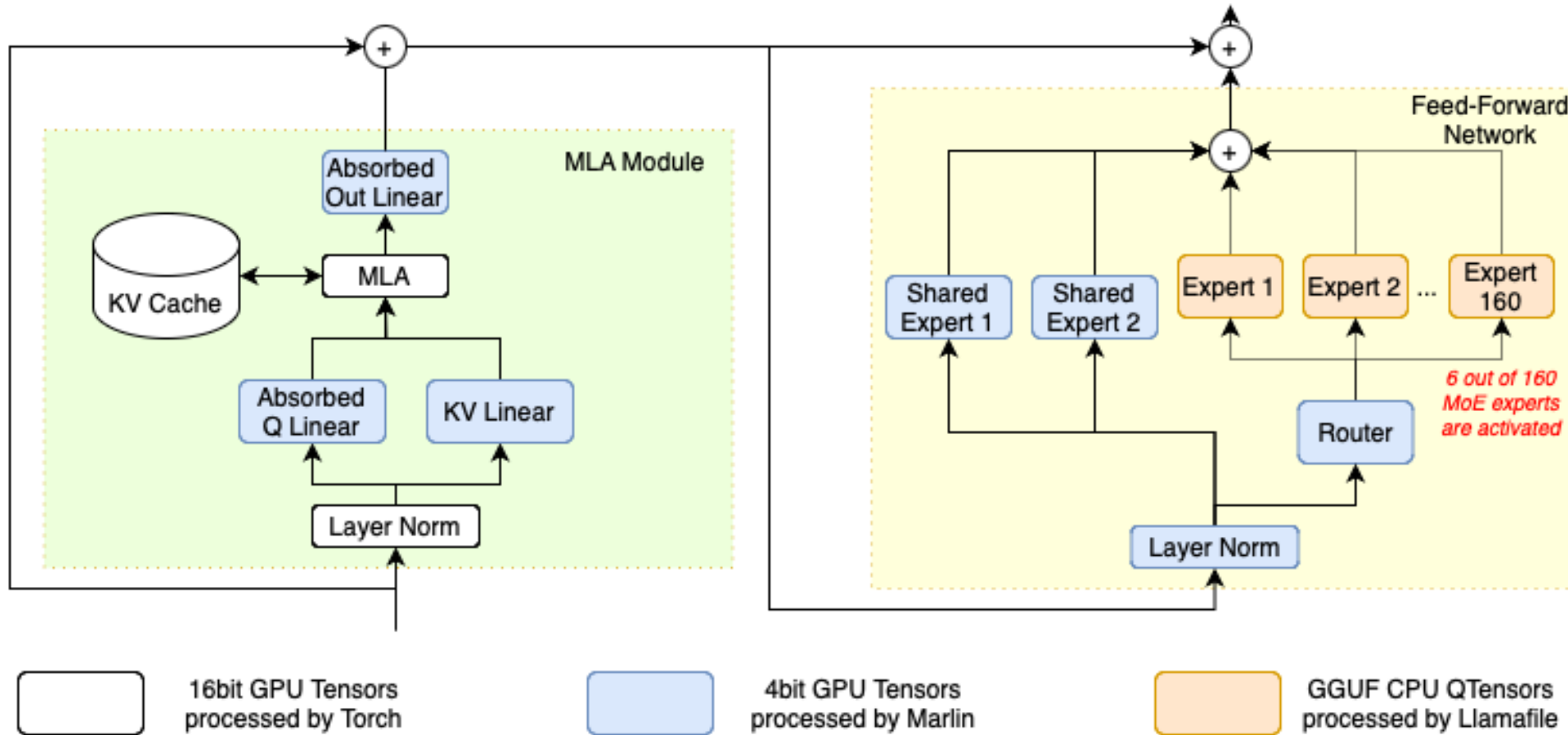
**Good for Capacity,
bad for Bandwidth
and Compute
Which part is
more suitable?**

**Sparsity!**

!!! The price numbers are not accurate, just a demonstration!
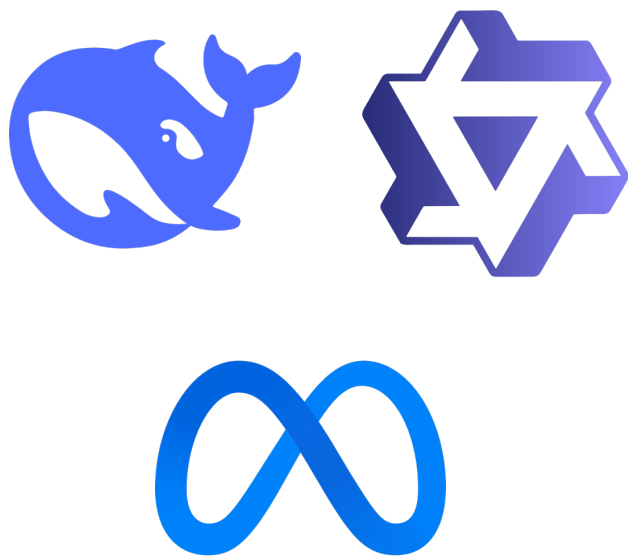
- DeepSeek Archtecture



Offload
Priority

Routed
Experts

>

Shared
Experts
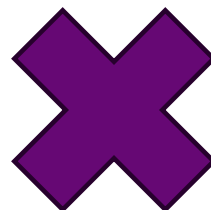
>

MLA
Attention

Different Models ✕ Different Hardware

## Decode

Latency, again,

the latencies!

CUDAGraph

is the key

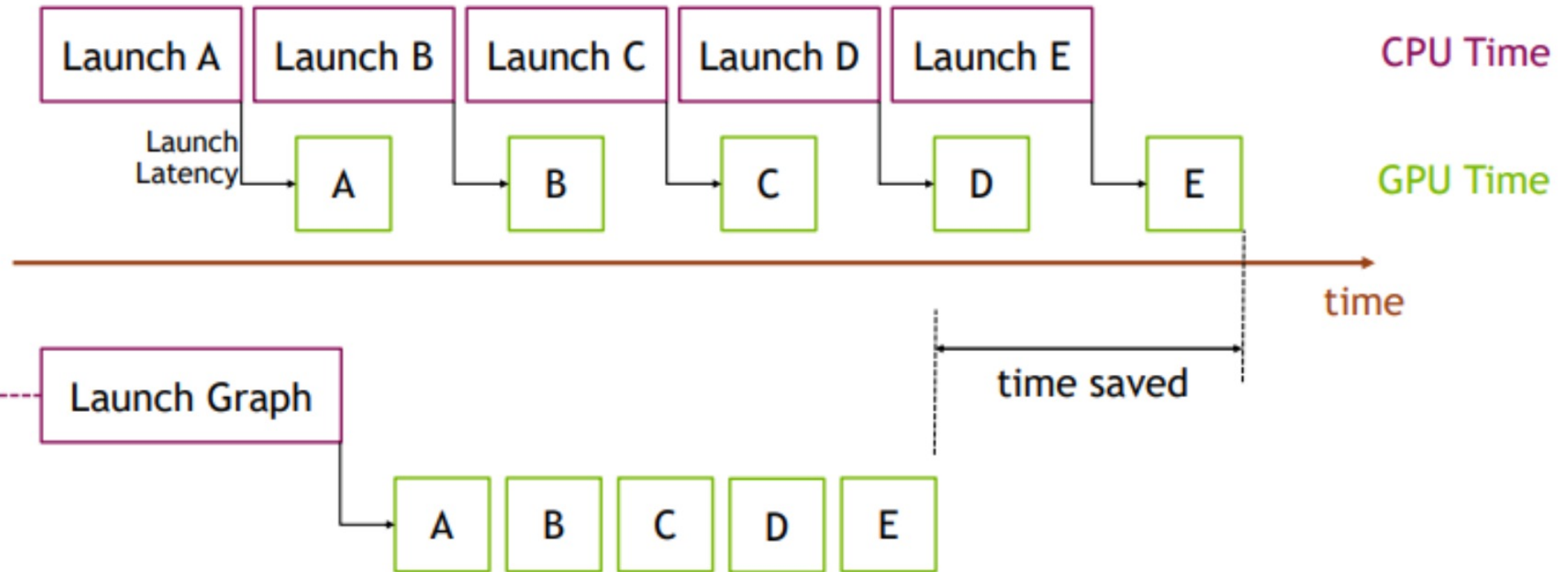(only one launch each forward)

## Prefill

CPU is too weak, even GQA 8

becomes compute bound
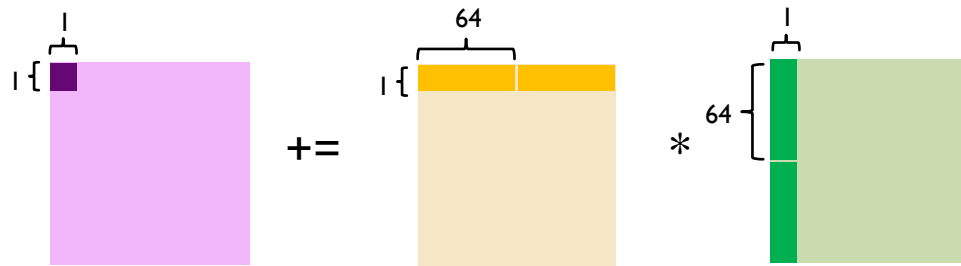
New hardware:

Intel AMX

# CUDA Graph



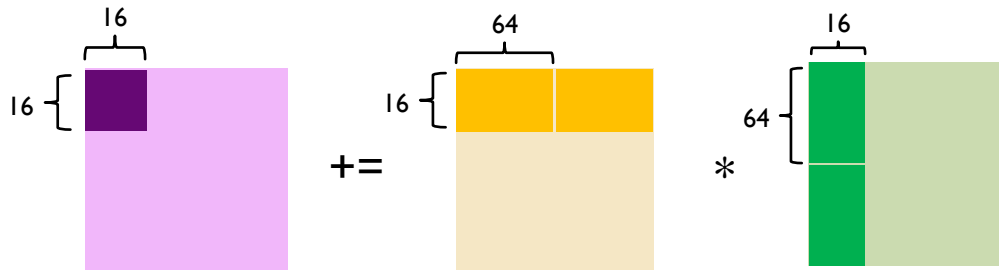**How to handle dynamic shape in continues-batched decoding?**

# Intel Advanced Matrix Extensions (Intel AMX)
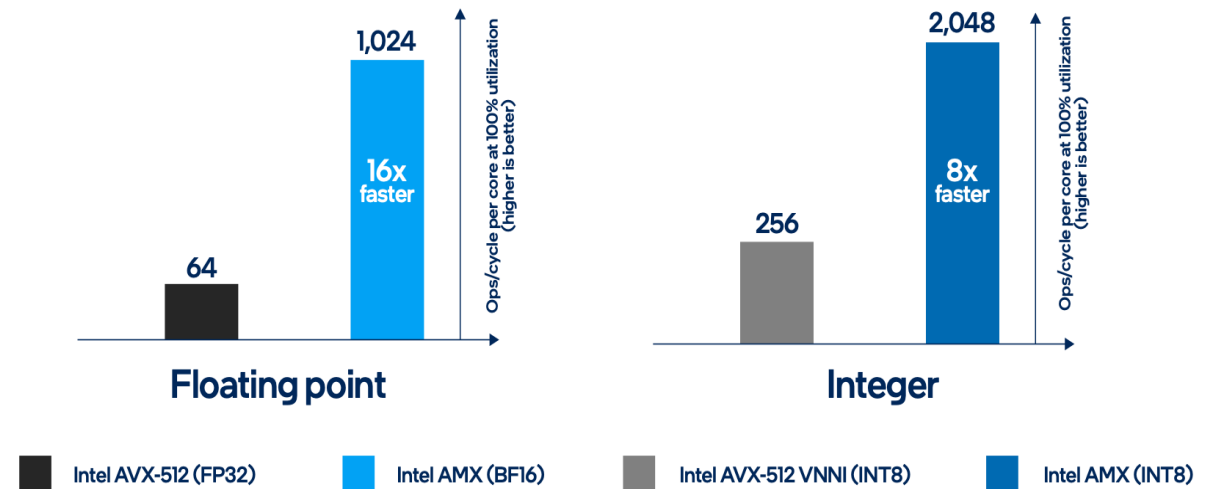
How AVX-512 solves INT8 matrix multiplication problems

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

How AMX solves INT8 matrix multiplication problems

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

AMX offers better performance than AVX-512  for INT8 and BF16 data types.



**Floating point**
- Intel AVX-512 (FP32): 64
- Intel AMX (BF16): 1,024 — 16x faster

**Integer**
- Intel AVX-512 VNNI (INT8): 256
- Intel AMX (INT8): 2,048 — 8x faster

Ops/cycle per core at 100% utilization (higher is better)

# Key optimization of matrix multiplication with Intel AMX

Due to the high computational capability of AMX instructions, memory bandwidth becomes a bottleneck, and the key to optimization lies in <span style="color:red">improving cache hit rates.</span>
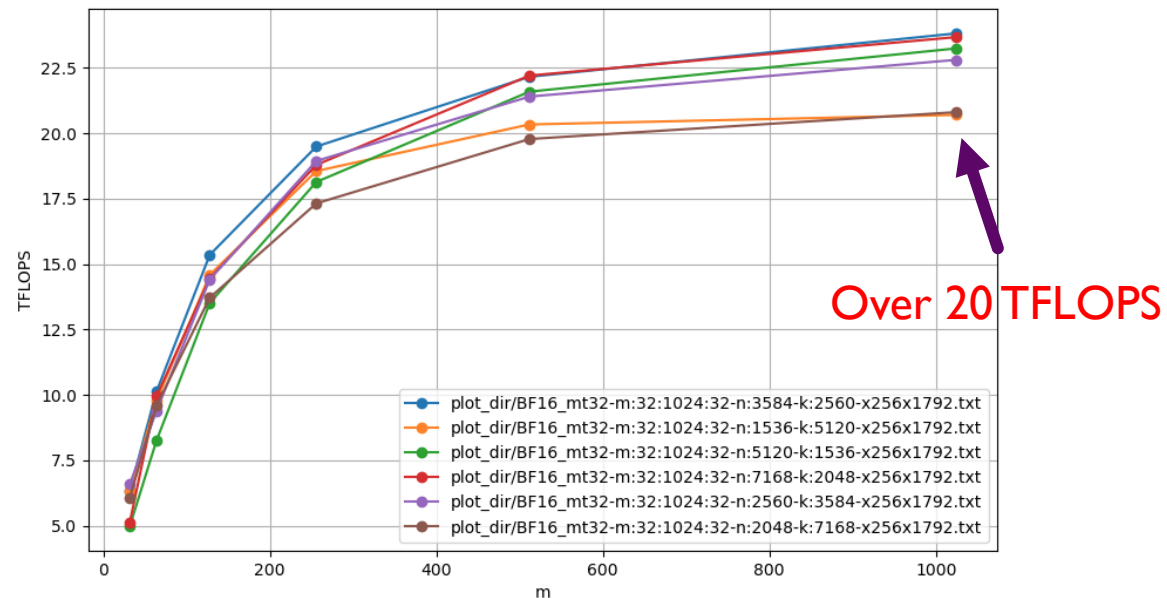
1. Design memory layout based on access patterns to enhance data locality, fully leverage hardware prefetch capabilities, and improve **L1 cache** hit rates.

2. Split large matrices into smaller tiles based on L2 cache size, ensuring that only the current tile is accessed at a time, improving **L2 cache** hit rates.

3. Dynamic work-scheduling，increasing data sharing between threads and optimizing **L3 cache** hit rates.

4. Quantize matrices by rows/columns, maintaining precision while reducing the number of scaling operations (and memory access) for each output element.
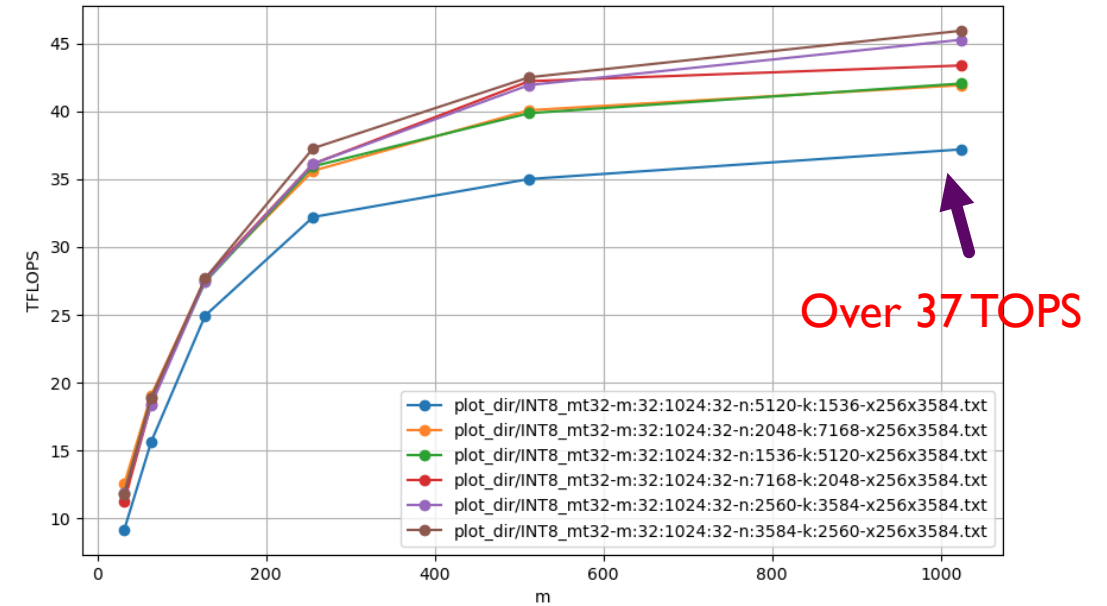
# Applications of Intel AMX

Applications： AMX kernel for sparse layers of MoE models (Deepseek R1/V3/V2, Mixtral, etc.)



BF16 TFLOPS vs. avg. # of selections

INT8 TOPS vs. avg. # of selections

Over 20 TFLOPS

Over 37 TOPS

In the matrix multiplication micro-benchmark, achieve over 20 TFLOPS and 37 TOPS of computational performance.
In the sparse MOE layer, achieve over 18TFLOPS and 30 TOPS of end-to-end computational performance.

# The Results: KTransformers



**KTransformers**

**A flexible heterogeneous inference framework**

- 24GB VRAM +382GB DRAM for 761B q4
- Currently several times faster than llama.cpp

https://github.com/kvcache-ai/ktransformers

# Impact of quantization

| DataSet | CPU Weight Format | CPU Kernel | GPU Weight Format | GEMM Kernel | MLA Kernel | Siliconflow | Ktrans Point |
|---------|-------------------|------------|-------------------|-------------|------------|-------------|--------------|
| MMLU (shuffle 1k) | | | | | | | |
| 1 | bf16 | cpuinfer | bf16 | torch | torch | 81.6 | 81.9 |
| 2 | q8_0 | cpuinfer | bf16 | torch | torch | 81.6 | 83.1 |
| 3 | q4km | cpuinfer | bf16 | torch | triton | 81.6 | 81.4 |
| 4 | q4km | cpuinfer | q4km->marlin 8 | marlin | triton | 81.6 | 81.1 |
| 5 | q4km | cpuinfer | q4km->marlin 4 | marlin | triton | 81.6 | 81 |
| 6 | q4km | cpuinfer | fp8 | fp8gemm | triton | 81.6 | 81.5 |
| MMLU-pro | | | | | | | |
| 1 | q4km | cpuinfer | fp8 | fp8gemm | triton | 57.7 | 57.6 |
| 2 | q4km | cpuinfer | q4km->marlin 4 | marlin | triton | 57.7 | 57.5 |
| HumanEval | tbd | tbd | tbd | tbd | tbd | tbd | tbd |
| GSM8K | tbd | tbd | tbd | tbd | tbd | tbd | tbd |

We only sample 1k from 10k MMLU dataset, test once, and do not use few shot, thus the score is lower than the number reported in paper. More will come, updated on Github repo.

- Support Continues Batch and Chunked Prefill via an asynchronous architecture learnt from SGLang

- Better total output tokens due to the share parts in GPU

# KTransformers v0.3：Qwen3 Support, the dawn of real AI PC

The post-trained models, such as **Qwen3-30B-A3B**, along with their pre-trained counterparts (e.g., **Qwen3-30B-A3B-Base**), are now available on platforms like **Hugging Face**, **ModelScope**, and **Kaggle**. For deployment, we recommend using frameworks like **SGLang** and **vLLM**. For local usage, tools such as **Ollama**, **LMStudio**, **MLX**, **llama.cpp**, and **KTransformers** are highly recommended. These options ensure that users can easily integrate Qwen3 into their workflows, whether in research, development, or production environments.



Intel 14900KF + 4090

Still much room for optimization!

| Qwen3 Model | Prefill | Single Decode | 4 Decode | (tokens/s) |
|---|---|---|---|---|
| 30B-A3B | 204 | 12 | 27 | smooth |
| 235B-A22B | 45 | 2.5 | 6 | workable |

# KTransformers v0.3.2：L3 Cache and Kimi K2

## Enabling Prefix Cache Mode in KTransformers

Balance serve now supports prefix cache reuse! To enable **Prefix Cache Mode** in KTransformers, you need to modify the configuration file and recompile the project.

### Step 1: Modify the Configuration File

Edit the `./ktransformers/configs/config.yaml` file with the following content (you can adjust the values according to your needs):

```yaml
attn:
  page_size: 16 # Size of a page in KV Cache.
  chunk_size: 256
kvc2:
  gpu_only: false # Set to false to enable prefix cache mode (Disk + CPU + GPU KV storage)
  utilization_percentage: 1.0
  cpu_memory_size_GB: 500 # Amount of CPU memory allocated for KV Cache
  disk_path: /mnt/data/kvc # Path to store KV Cache on disk
```

### Step 2: Update Submodules and Recompile

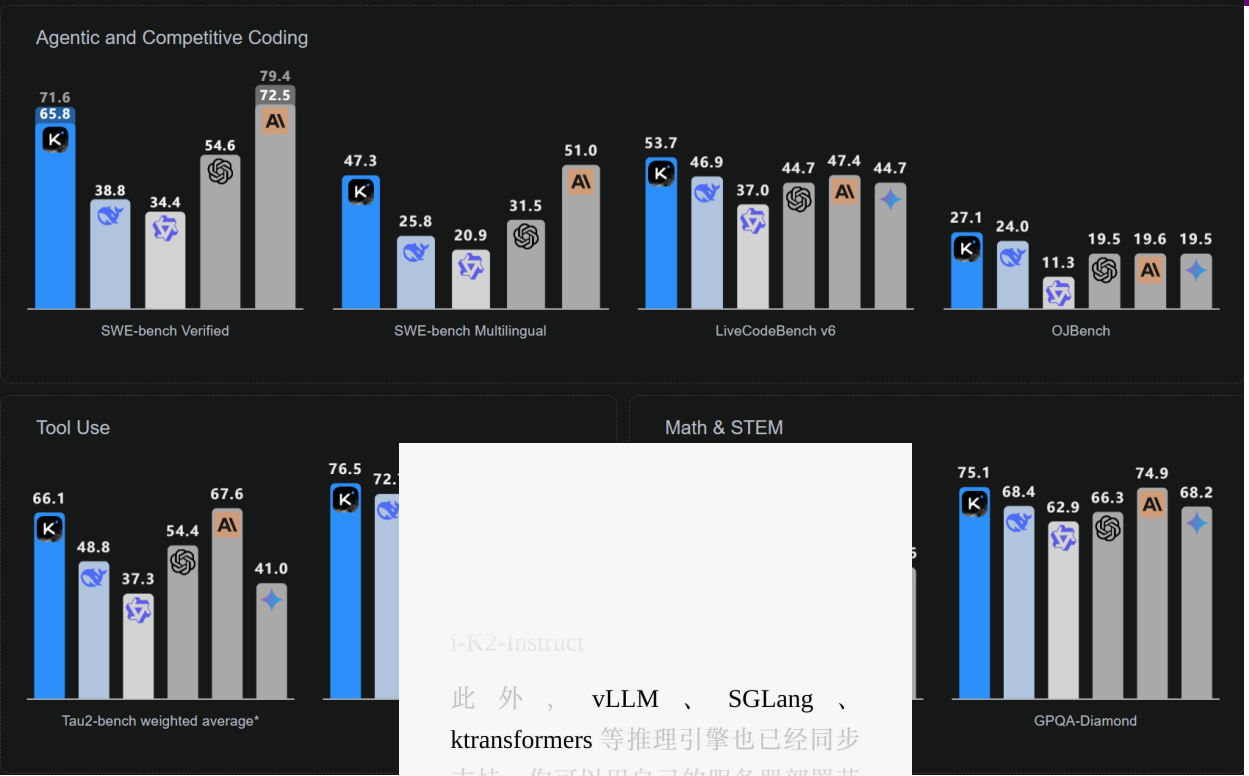If this is your first time using prefix cache mode, please update the submodules first:

```
git submodule update --init --recursive # Update PhotonLibOS submodule
```

Then recompile the project:

```
# Install single NUMA dependencies
USE_BALANCE_SERVE=1  bash ./install.sh
# For those who have two cpu and 1T RAM（Dual NUMA）：
USE_BALANCE_SERVE=1 USE_NUMA=1 bash ./install.sh
```

## Note

Balance serve utilizes a 3-layer (GPU-CPU-Disk) scheme to store and reuse KVCache. Deleting KVCache is not supported now. If you have too much KVCache, you can simply delete them by remove kvcache files.

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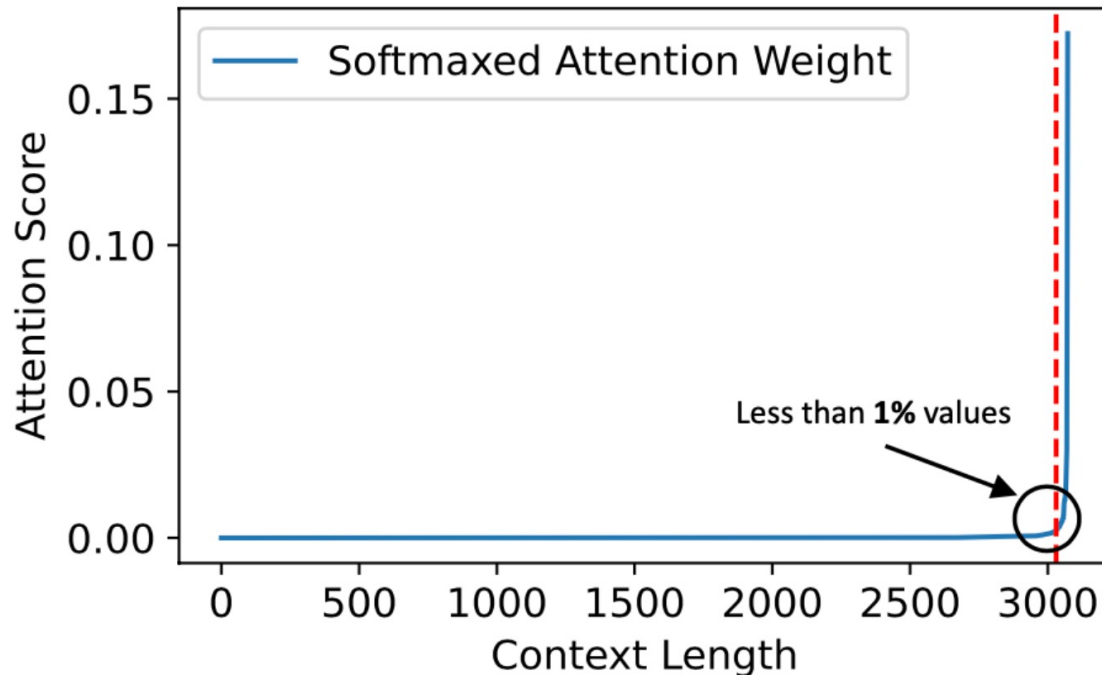
**MoBA: Mixture of Block Attention for Long-Context LLMs**

**deepseek**

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

Natural sparsity because of Softmax

Learnable Block-based Sparsity

# A flexible CPU offload framework for sparse attention



- Initial and local tokens, split into blocks and dynamically selected

- Integrate with Quest, SnapKV, InfLLM, …

- 100% in NIHS <- but this is a weak bmk

- A combination of single-operator optimizations with each other



```
1  - match:
2      name: "^model\\.layers\\..*\\.mlp\\.experts$"
3      replace:
4        class: ktransformers.operators.experts.KTransformersExperts  # custom MoE Kernel
5        kwargs:
6          generate_device: "cpu"
7          generate_op:  "KExpertsCPU"
```

```
1  - match:
2      name: "^model\\.layers\\..*\\.mlp\\.shared_experts"  # regular expression
3      class: torch.nn.Linear  # only match modules matching name and class simultaneously
4      replace:
5        class: ktransformers.operators.linear.KTransformersLinear  # optimized Kernel on
   quantized data types
6        kwargs:
7          generate_device: "cuda"
8          generate_op: "KLinearMarlin"
```

# Thanks!

## kvcache.ai

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

👥 **758 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++   ⭐ 3.3k   🍴 264

**ktransformers** (Public)

A Flexible Framework for Experiencing Cutting-edge LLM Inference Optimizations

🔵 Python   ⭐ 14.2k   🍴 1k

https://github.com/kvcache-ai