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

**Le Zhihao**

# Outline

❑**Background**

❑Motivation

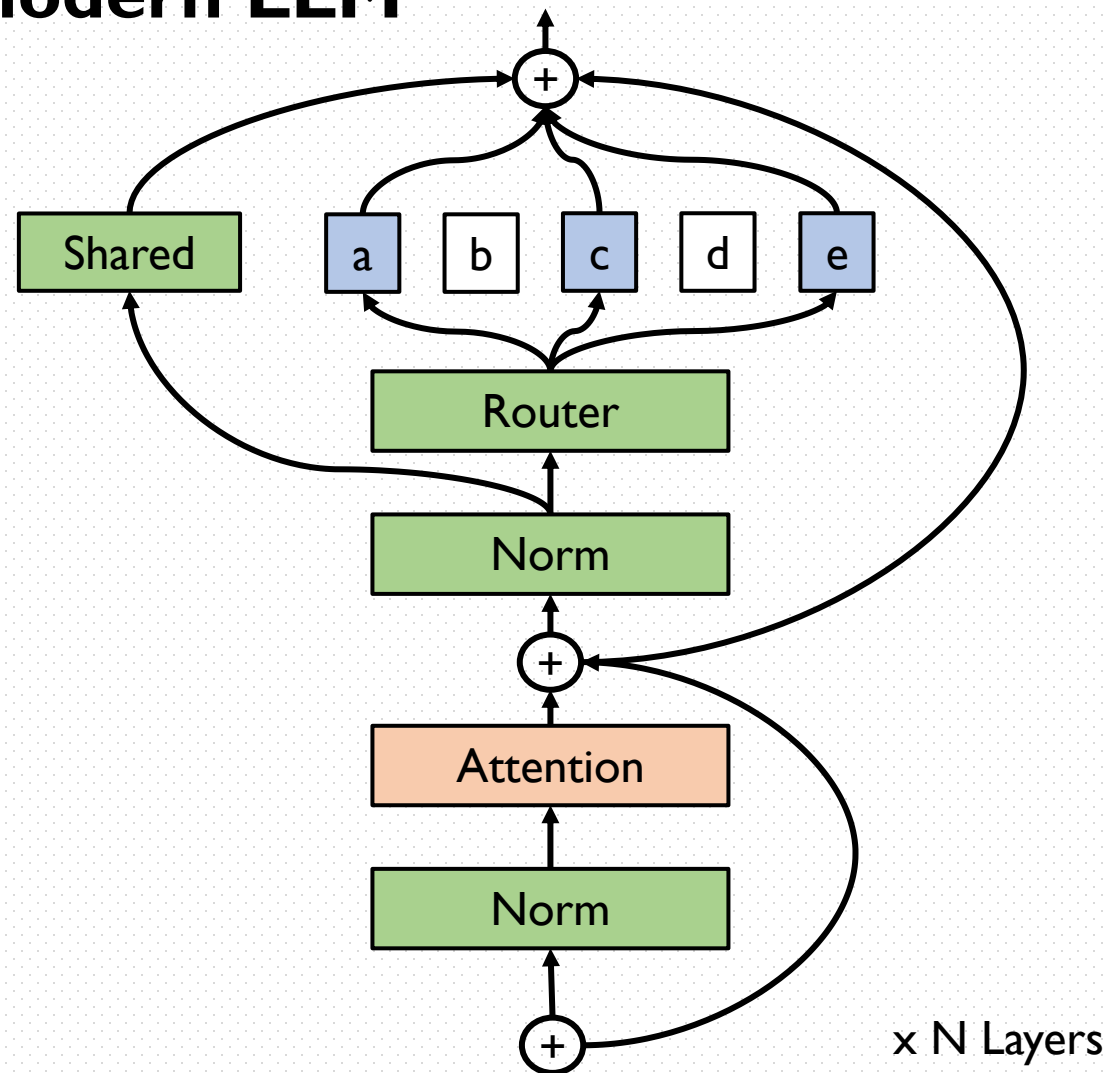❑Challenge

❑Design

❑Evaluation

❑Conclusion

# Background

❑ **MoE model is everywhere in modern LLM**
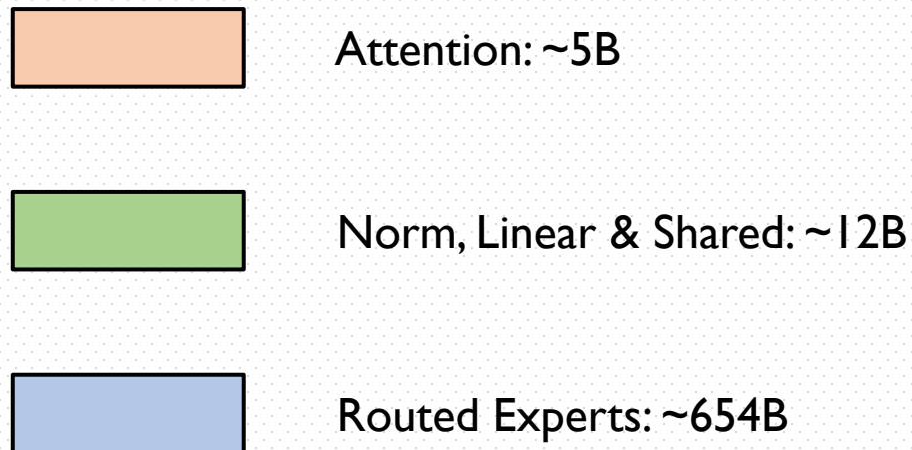
❖ **Qwen3, DeepSeekV3/R1**



x N Layers

# Background

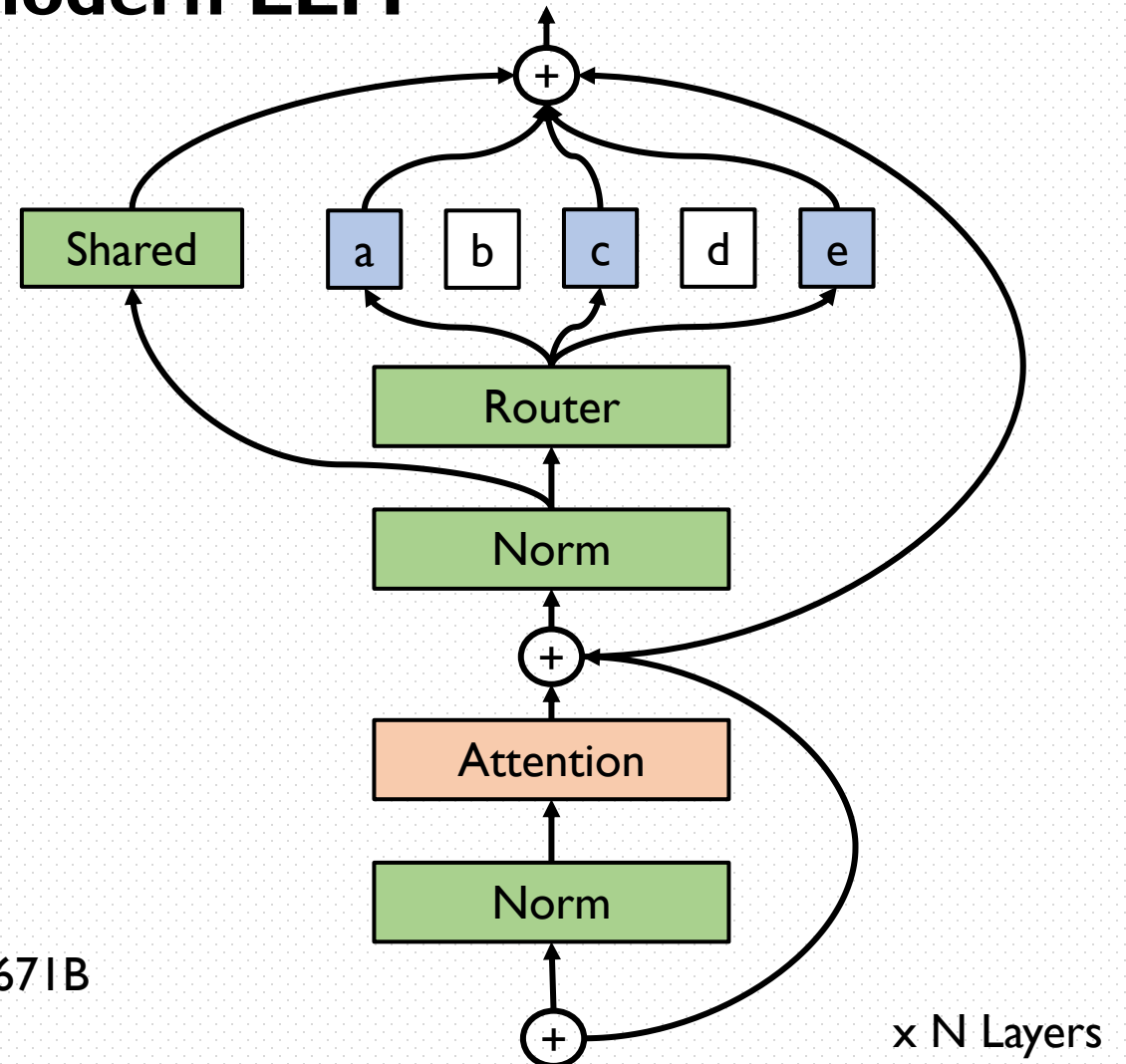❑ **MoE model is everywhere in modern LLM**

❑ **Memory becomes bottleneck**

How to deal with memory bottleneck with constrained GPU memory?

Attention: ~5B

Norm, Linear & Shared: ~12B

Routed Experts: ~654B

For DeepSeekV3 MoE 671B



x N Layers

# Background

❏ **MoE model is everywhere in modern LLM**

❏ **Memory becomes bottleneck**

❏ **Hybrid CPU/GPU inference**

Attention: ~5B

Norm, Linear & Shared: ~17B

Routed Experts: ~654B

GPU

CPU

For DeepSeekV3 MoE 671B

Shared | a | b | c | d | e

Router

Norm

Attention

Norm

+

x N Layers

# Background

❑ **MoE model is everywhere in modern LLM**

❑ **Memory becomes bottleneck**

❑ **Hybrid CPU/GPU inference**

Compute in GPU

Attention: ~5B

Norm, Linear & Shared: ~17B

GPU

PCIe: 32 GB/s

Routed Experts: ~654B

CPU  Mem: 440 GB/s

Compute in CPU      For DeepSeekV3 MoE 671B

Shared    a  b  c  d  e

Router

Norm

Attention

Norm

x N Layers

# Background

☐ **GPU only**

Attention    Router    Shared    Experts

GPU | Attn | | Shared | a | b | c | Attn | | Shared | a | b | c | Attn

☐ **Hybrid CPU/GPU inference**

GPU | Attn | | Shared | | Attn | | Shared | | Attn

CPU | a | b | c | | a | b | c

# Recent Work

❑ **Llama.cpp[1]**

   ❖ **C++ based LLM inference enabling heterogeneous execution.**

❑ **Fiddler[2]**

   ❖ **Support expert offloading and selectively reload experts.**

One A100 and two Intel Xeon CPUs:

- Prefill: 70.02 tokens per second
- Decode: 4.68 tokens per second
- Low GPU utilization (below 30%)

[1] Georgi Gerganov 2023. ggerganov/llama.cpp. Retrieved Feb 8, 2025 from https://github.com/ggerganov/llama.cpp
[2] Keisuke Kamahori, Yile Gu, Kan Zhu, and Baris Kasikci. 2024. Fiddler: CPU-GPU Orchestration for Fast Inference of Mixture-of-Experts Models. arXiv:2402.07033 [cs.LG]

# Outline

❑ **Background**
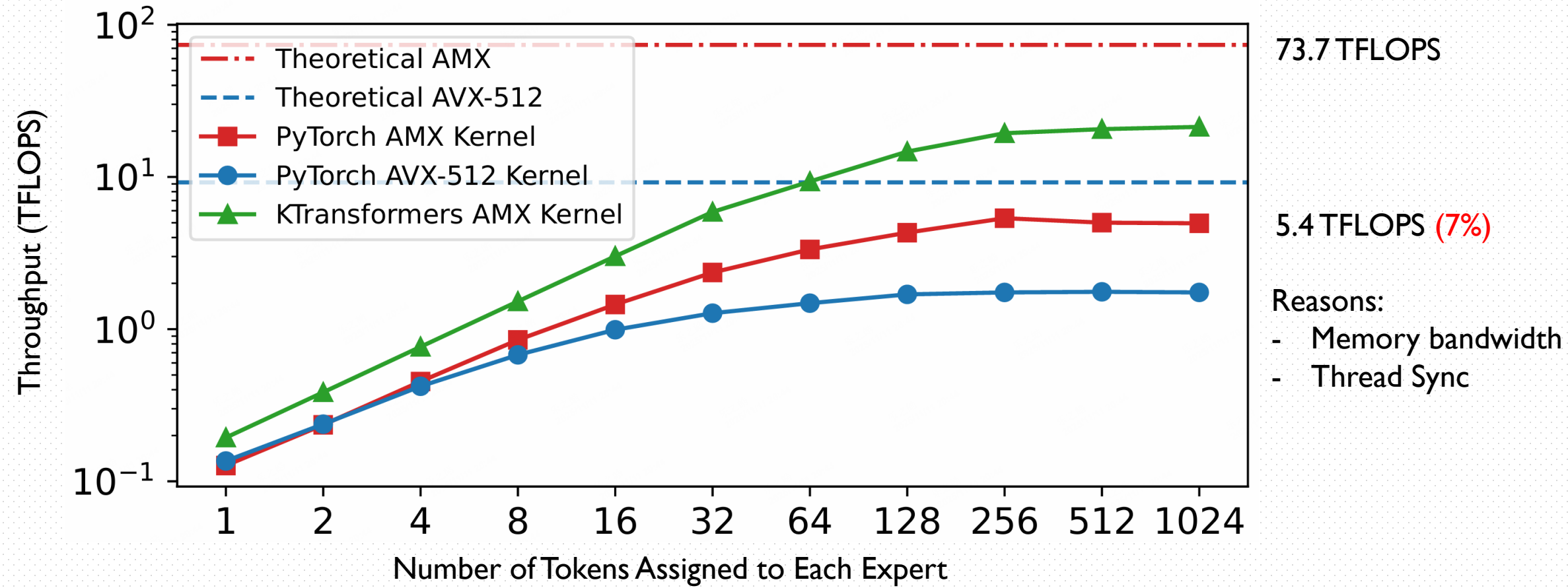
❑ **Motivation**

❑ Challenge

❑ Design

❑ Evaluation

❑ Conclusion

# Motivation

## ❑Underutilized CPU compute resources



Throughput of the MoE Layers on DeepSeek-V3 using PyTorch's AMX and AVX-512 kernels
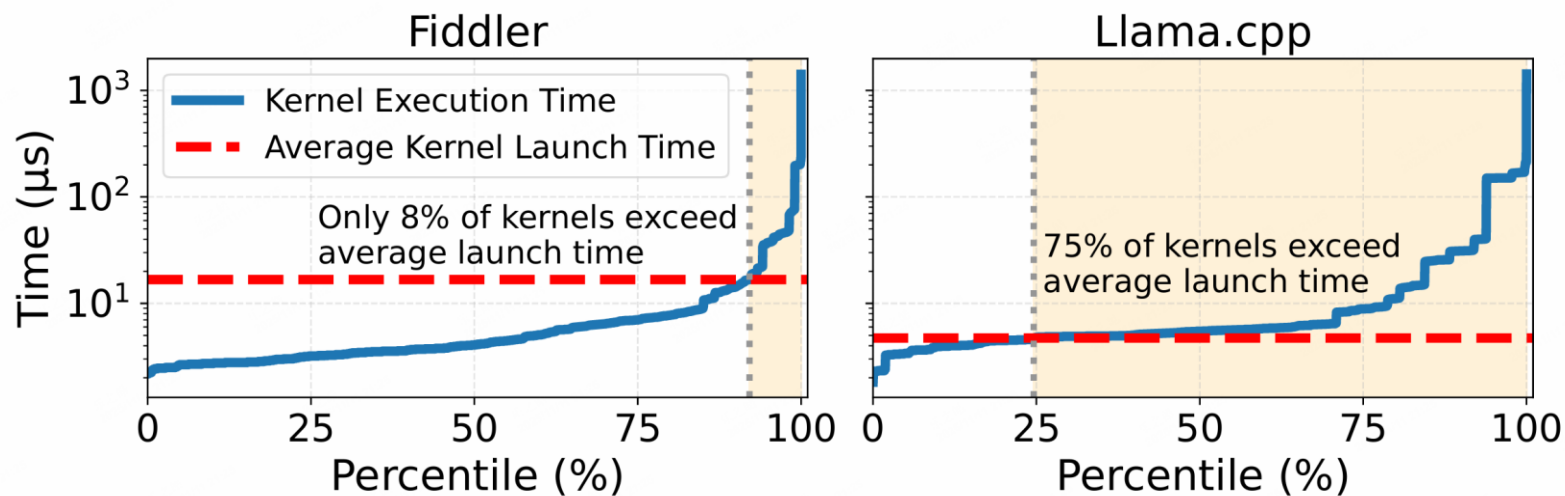
# Motivation

❑ **Underutilized CPU compute resources**

❑ **CPU-GPU/CPU coordination**

    ❖ **CPU-GPU coordination**

       ➤ **High kernel launch latency**



GPU kernel launch and execution time analysis of DeepSeek-V3 in A100

# Motivation

❑ **Underutilized CPU compute resources**

❑ **CPU-GPU/CPU coordination**

   ❖ **CPU-GPU coordination**

      ➢ **High kernel launch latency**

      ➢ **CUDA graph fails to support CPU and GPU overlapping computation**

# Motivation

❑ **Underutilized CPU compute resources**

❑ **CPU-GPU/CPU coordination**

❖ **CPU-GPU coordination**

❖ **CPU-CPU coordination**

➢ **Inefficient memory access NUMA nodes**

- **DeepSeek-V3 using Fiddler on a single socket: 6.9ms**
- **DeepSeek-V3 using Fiddler on two sockets: 5.8ms (-16%)**

# Outline

❑**Background**

❑**Motivation**

❑**Challenge**

❑**Design**

❑**Evaluation**

❑**Conclusion**

# Challenge

❑ **Underutilized CPU compute resources**

    ❖ **Memory bandwidth constraints**

    ❖ **Thread synchronization overhead**

❑ **CPU-GPU/CPU coordination**

    ❖ **CPU-GPU: high overhead of kernel invocation and synchronization**

    ❖ **CPU-CPU: inefficient cross-socket memory access**

# Outline
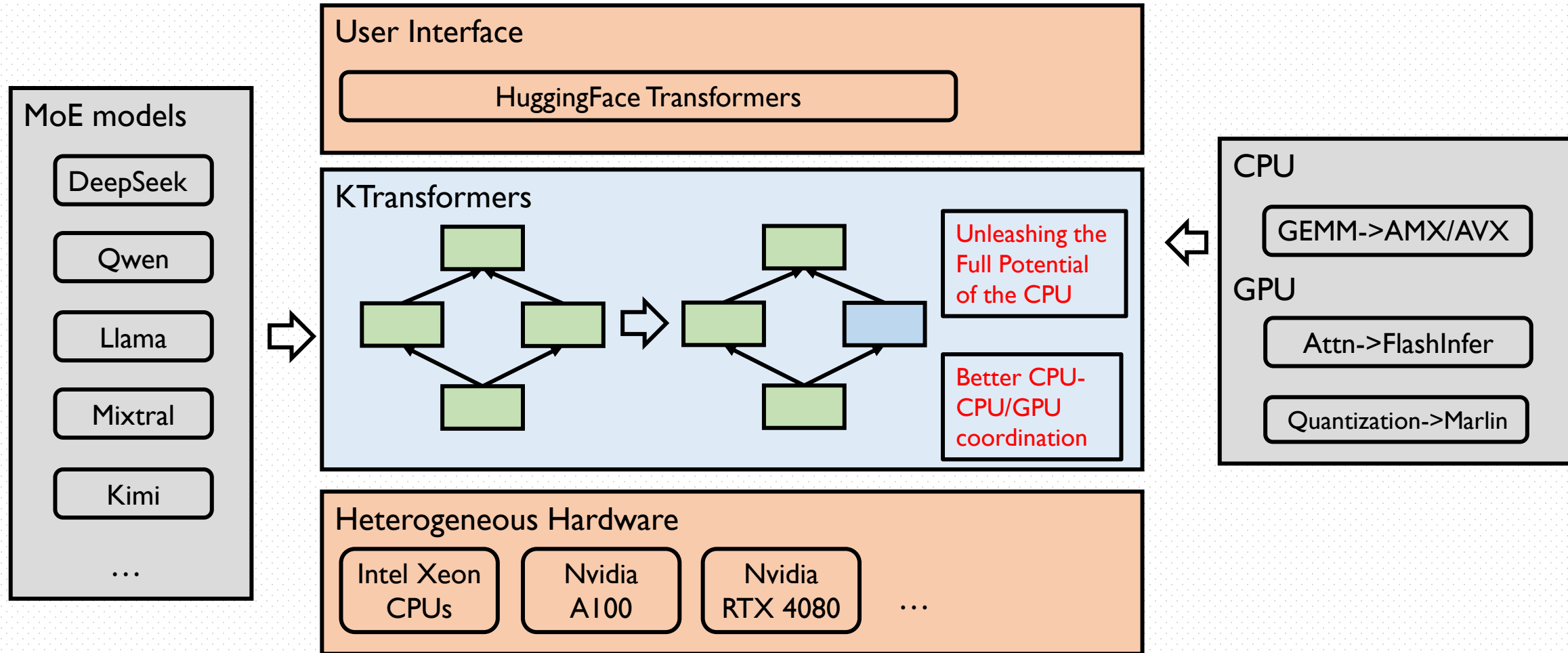
❑**Background**

❑**Motivation**

❑**Challenge**

❑**Design**

❑**Evaluation**

❑**Conclusion**

# Design - Overview

**MoE models**
- DeepSeek
- Qwen
- Llama
- Mixtral
- Kimi
- …

**User Interface**
- HuggingFace Transformers

**KTransformers**

Unleashing the Full Potential of the CPU

Better CPU-CPU/GPU coordination

**CPU**
- GEMM->AMX/AVX

**GPU**
- Attn->FlashInfer
- Quantization->Marlin

**Heterogeneous Hardware**
- Intel Xeon CPUs
- Nvidia A100
- Nvidia RTX 4080
- …

# Design - Unleashing the Full Potential of the CPU

## ❑ AMX Tiling-aware Memory Layout

| Expert Weight | → preprocessed → | 64 bytes submatrix |
|---|---|---|

64 bytes submatrix

64 bytes submatrix

…

64 bytes submatrix

→ AMX instructions

# Design - Unleashing the Full Potential of the CPU

❑ **AMX Tiling-aware Memory Layout**

❑ **Cache-Friendly AMX Kernels**

# Design - Unleashing the Full Potential of the CPU

❑ **AMX Tiling-aware Memory Layout**

❑ **Cache-Friendly AMX Kernels**

❑ **Adaptive AVX-512 Kernel for Low ARI Scenarios**



Lower is better

# Design - Unleashing the Full Potential of the CPU

❑ **AMX Tiling-aware Memory Layout**

❑ **Cache-Friendly AMX Kernels**

❑ **Adaptive AVX-512 Kernel for Low ARI Scenarios**
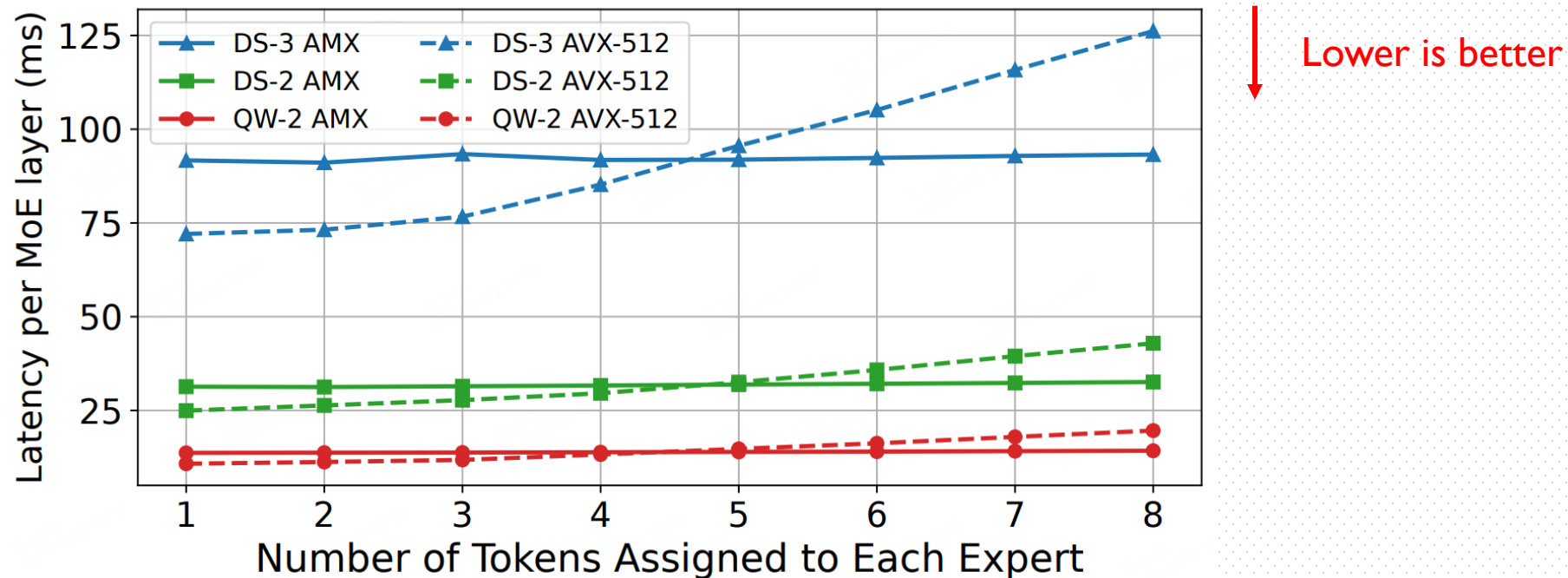
# Design - Unleashing the Full Potential of the CPU

❑ **AMX Tiling-aware Memory Layout**

❑ **Cache-Friendly AMX Kernels**

❑ **Adaptive AVX-512 Kernel for Low ARI Scenarios**

❑ **Fuse MoE Ops**



2. Combine gate and up projections where no dependencies.

1. Merge gate projections, up and down across experts.

# Design - Unleashing the Full Potential of the CPU

❑ **AMX Tiling-aware Memory Layout**

❑ **Cache-Friendly AMX Kernels**

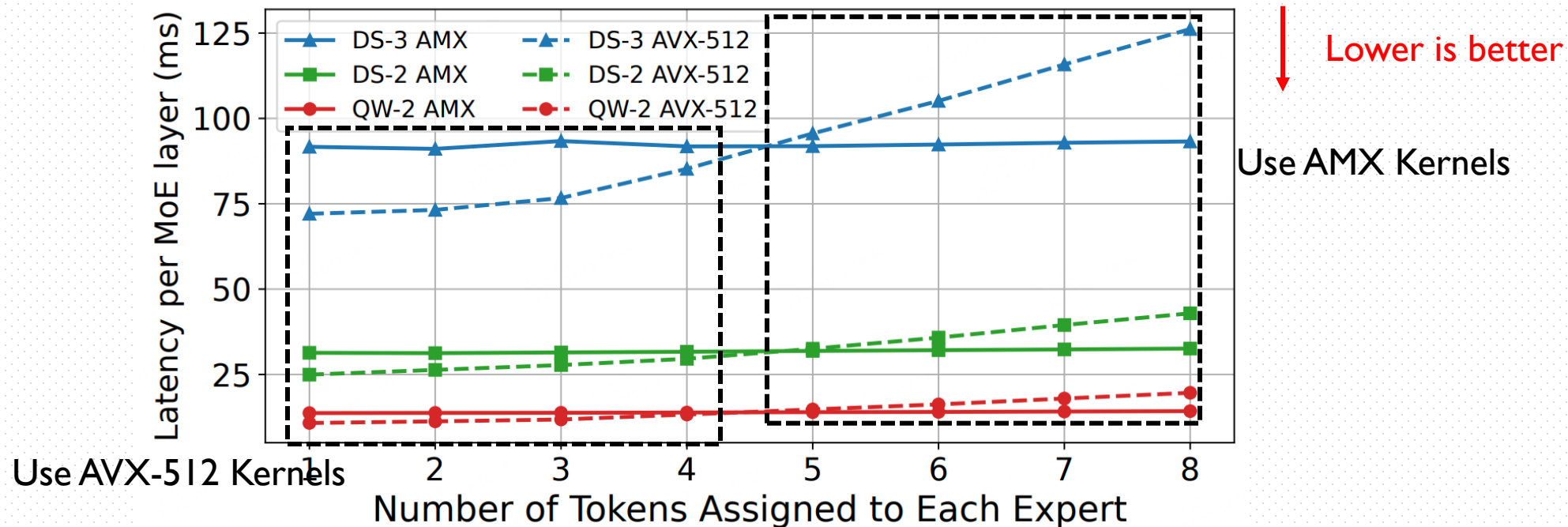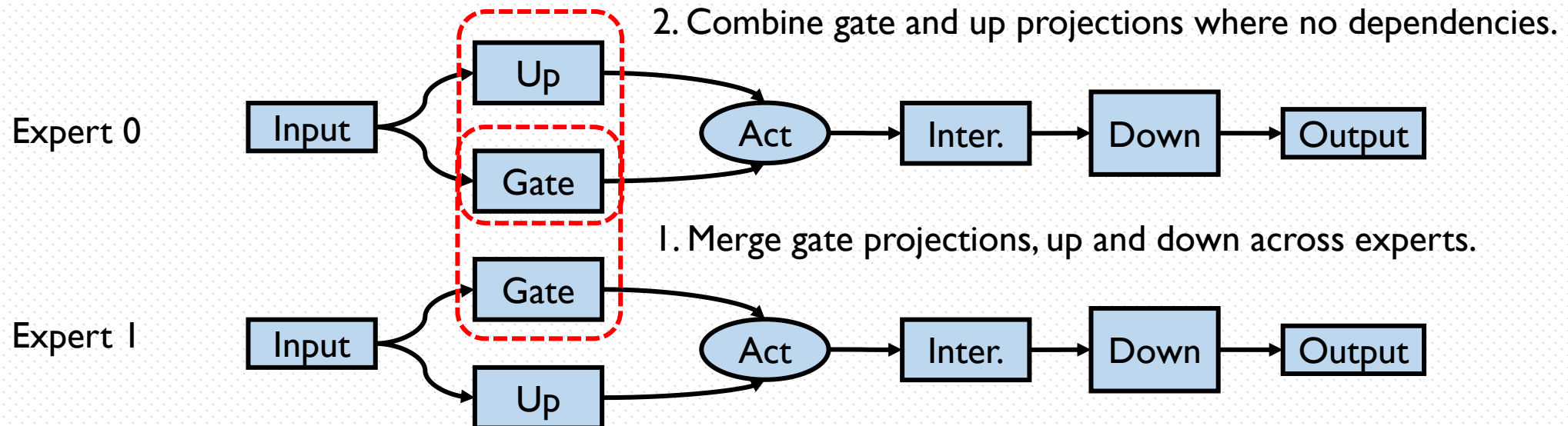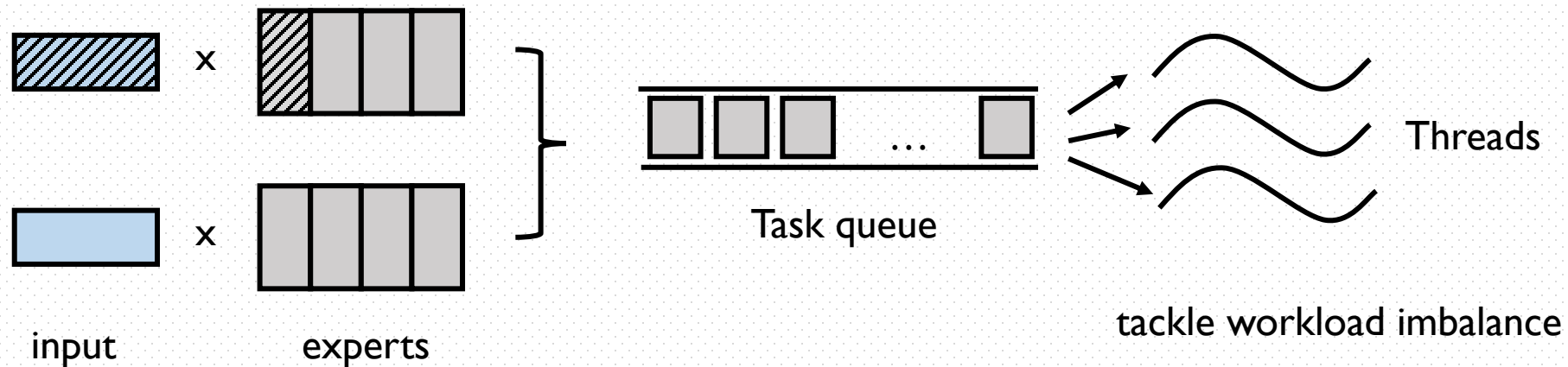❑ **Adaptive AVX-512 Kernel for Low ARI Scenarios**

❑ **Fuse MoE Ops**

❑ **Dynamic Task Scheduling**
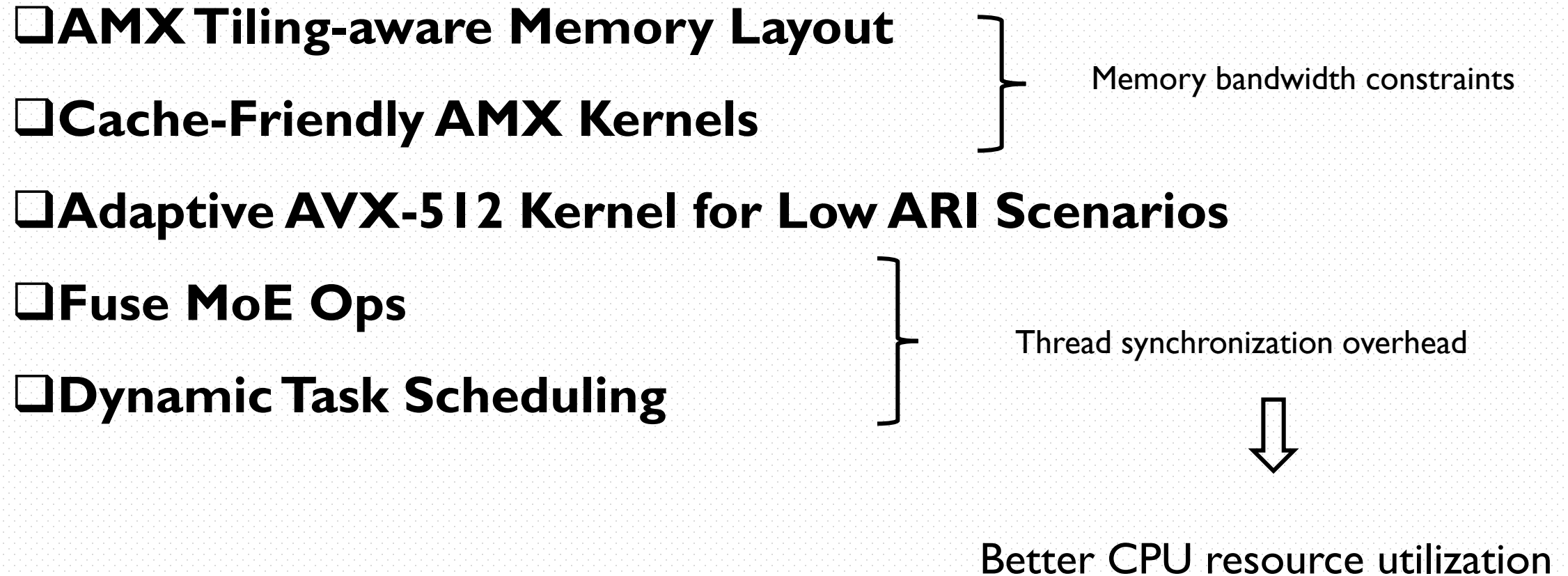
input    experts    Task queue    Threads

tackle workload imbalance

# Design - Unleashing the Full Potential of the CPU

❑ **AMX Tiling-aware Memory Layout**

❑ **Cache-Friendly AMX Kernels**

Memory bandwidth constraints

❑ **Adaptive AVX-512 Kernel for Low ARI Scenarios**

❑ **Fuse MoE Ops**

Thread synchronization overhead

❑ **Dynamic Task Scheduling**

Better CPU resource utilization

# Design - Better CPU-CPU/GPU Coordination

❑ **Asynchronous CPU-GPU Task Scheduling Mechanism**



Can't be captured by CUDA graph

# Design - Better CPU-CPU/GPU Coordination

❑**Asynchronous CPU-GPU Task Scheduling Mechanism**

# Design - Better CPU-CPU/GPU Coordination

❑ **Asynchronous CPU-GPU Task Scheduling Mechanism**

❑ **NUMA-aware Tensor Parallelism**



Local Weights/States

Remote Weights/States

Expert parallel is inefficient when imbalance.

# Design - Better CPU-CPU/GPU Coordination

❑ **Asynchronous CPU-GPU Task Scheduling Mechanism**

❑ **NUMA-aware Tensor Parallelism**
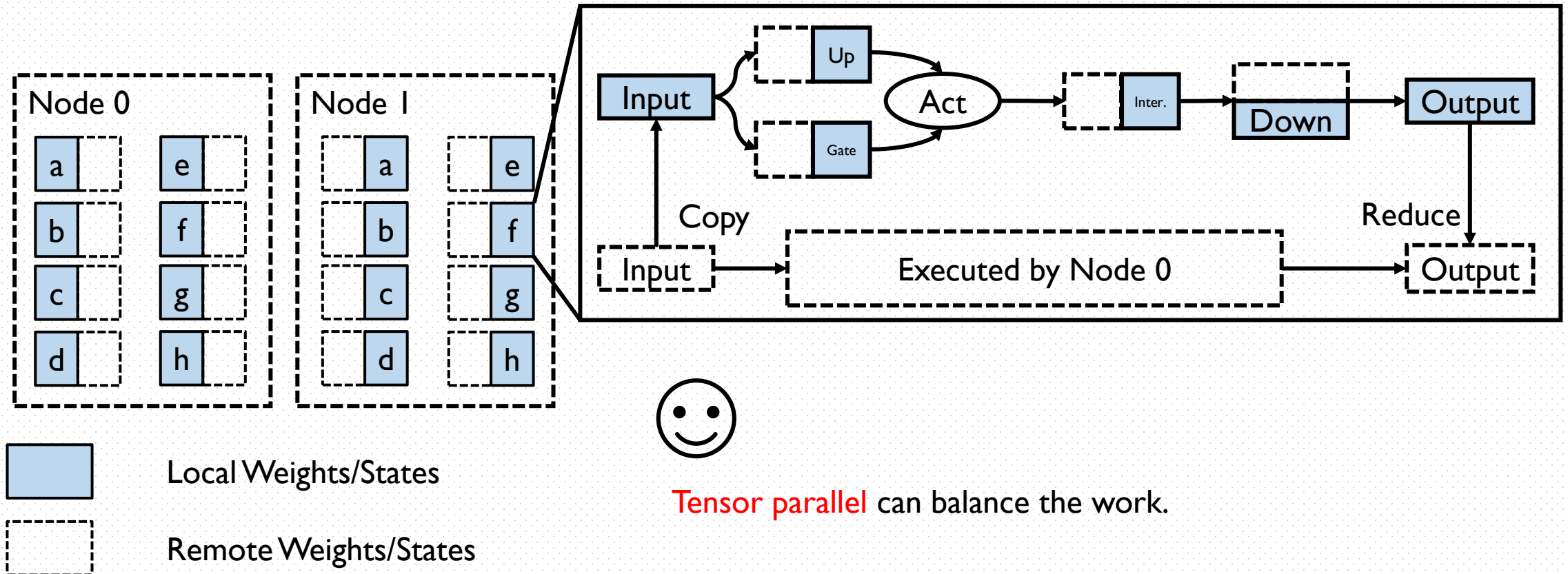


Tensor parallel can balance the work.

# Design - Better CPU-CPU/GPU Coordination

❑ **Asynchronous CPU-GPU Task Scheduling Mechanism**

❑ **NUMA-aware Tensor Parallelism**

High overhead of kernel invocation and synchronization

Inefficient cross-socket memory access

⇩

Better CPU-CPU/GPU coordination

# Design – Expert Deferral

❑**Hybrid CPU/GPU inference**



| GPU | Attn | | Shared | idle | Attn | | Shared | idle | Attn |

| CPU | a | b | c | idle | a | b | c |

❑**Hybrid CPU/GPU inference + Expert Deferral**



| GPU | 0.Attn | | 0.Shared | 1.Attn | | 1.Shared | 2.Attn |

| CPU | 0.a | 0.b | 0.c | 1.a | 1.b | 1.c |

# Design – Expert Deferral



Layer k+1 MoE

Shared

a b c d e

Router

Norm

Layer k+1 Attention

Attention

Norm

Layer k MoE

Shared

a b c d e

Router

Norm

Workflow

# Design – Expert Deferral



Layer k+1 MoE

Shared

a  b  c  d  e

Router

Norm

Layer k+1 Attention

Attention

Norm

Layer k MoE

Shared

a  b  c  d  e

Router

Norm

a  Immediate experts

c  deferred experts

Workflow

Layer k+1 MoE

Shared

a  b  c  d  e

Router

Norm

Layer k+1 Attention

Attention

Norm

Layer k MoE

Shared

a  b  c  d  e

Router

Norm

# Design – Expert Deferral

❑ **How to decide the expert deferral configuration?**

# Design – Expert Deferral

❑ **How to decide the expert deferral configuration?**

GPU utilization: 28%
CPU utilization: 74%

**Legend:**
Wait · Send + Submit · Immediate Experts
Attention · Shared Experts · Deferred Experts
Gate · Sync + Receive

### 8 Immediate Experts + 0 Deferred Experts

GPU
CPU

Timeline (µs)
0    500    1000    1500    2000    2500

CPU GPU timelines in the MoE layer of DeepSeek V3

# Design – Expert Deferral

❑**How to decide the expert deferral configuration?**

Heuristic ways:
1. Achieve full CPU utilization
2. Ensure 2 immediate experts to maintain model accuracy

Get the best number of experts. ⇐



CPU GPU timelines in the MoE layer of DeepSeek V3

# Implementation – Flexible Module Injection

❑ **Build on HuggingFace Transformer:**

    ❖ **Lightweight injection framework**

        ➢ **Use a YAML file to drive the substitution**

    ❖ **Expose pybind11 to expose CPU kernels**

```
 1  - match:
 2      class: modeling_deepseek_v3.DeepseekV3MoE
 3    replace:
 4      class: operators.experts.FusedMoE
 5      device: "cpu"
 6      kwargs:
 7        backend: "hybrid_AMX_AVX512"
 8        data_type: "Int4"
 9        n_deferred_experts: 6
10
11  - match:
12      name: "^model\\.layers\\..*\\.self_attn$"
13    replace:
14      class: operators.attention.FlashInferMLA
15      device: "cuda:0"
16
17  - match:
18      name: "^(?!lm_head$).*"
19      class: torch.nn.Linear
20    replace:
21      class: operators.linear.MarlinLinear
22      device: "cuda:0"
23      kwargs:
24        data_type: "Int4"
```

Example configuration for adapting DeepSeek-V3
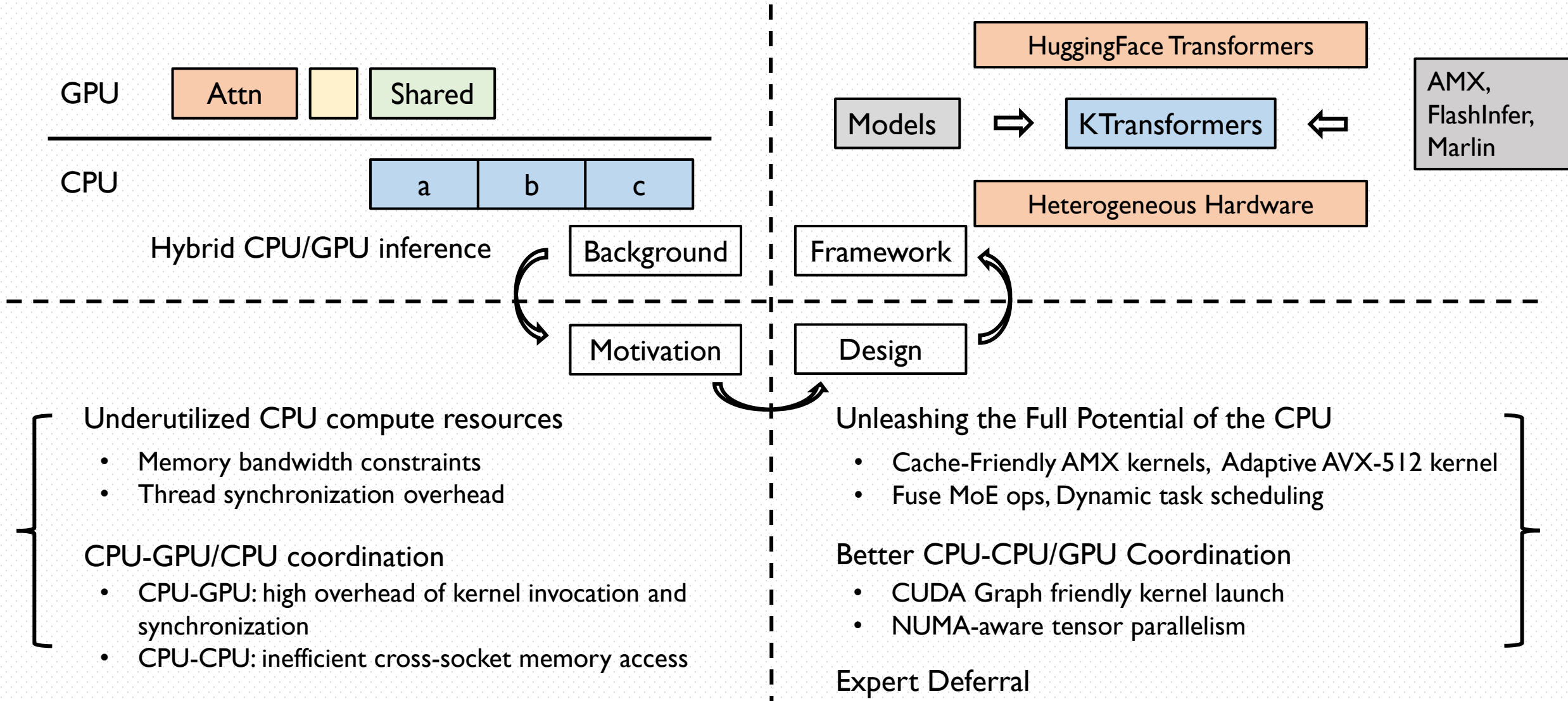
# Outline

❑ **Background**

❑ **Motivation**

❑ **Challenge**

❑ **Design**
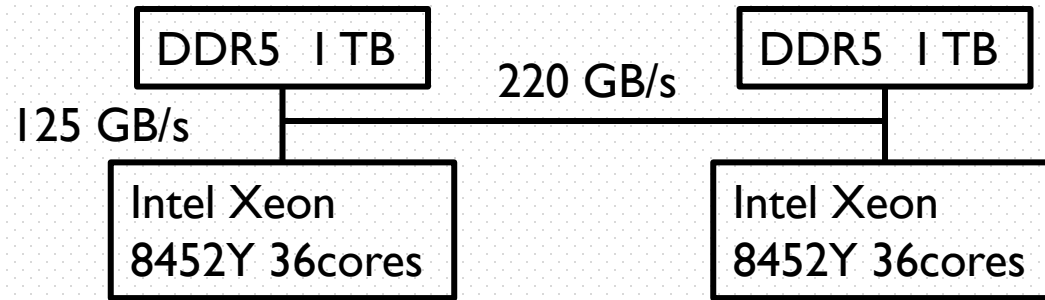
❑ **Evaluation**

❑ **Conclusion**

# Summary

GPU

| Attn | | Shared |

CPU

| a | b | c |

Hybrid CPU/GPU inference

Background

HuggingFace Transformers

Models ⇨ KTransformers ⇦ AMX, FlashInfer, Marlin

Heterogeneous Hardware

Framework

Motivation

Design

**Underutilized CPU compute resources**
- Memory bandwidth constraints
- Thread synchronization overhead

**CPU-GPU/CPU coordination**
- CPU-GPU: high overhead of kernel invocation and synchronization
- CPU-CPU: inefficient cross-socket memory access

**Unleashing the Full Potential of the CPU**
- Cache-Friendly AMX kernels, Adaptive AVX-512 kernel
- Fuse MoE ops, Dynamic task scheduling

**Better CPU-CPU/GPU Coordination**
- CUDA Graph friendly kernel launch
- NUMA-aware tensor parallelism

Expert Deferral

# Evaluation - Setup

❑ **Hardware**

   ❖ **CPU:**

| | | |
|---|---|---|
| DDR5  1 TB | 220 GB/s | DDR5  1 TB |

125 GB/s

| | |
|---|---|
| Intel Xeon 8452Y 36cores | Intel Xeon 8452Y 36cores |

   ❖ **GPU: a NVIDIA A100, a RTX 4080, Pcie 4.0**

❑ **Models:**

   ❖ **DeepSeek-V3-0324, DeepSeek V2.5-1210, Qwen2-57B-A14B**

❑ **Datasets:**

   ❖ **HumanEval, MBPP, GSM8K, StrategyQA, LiveBench**
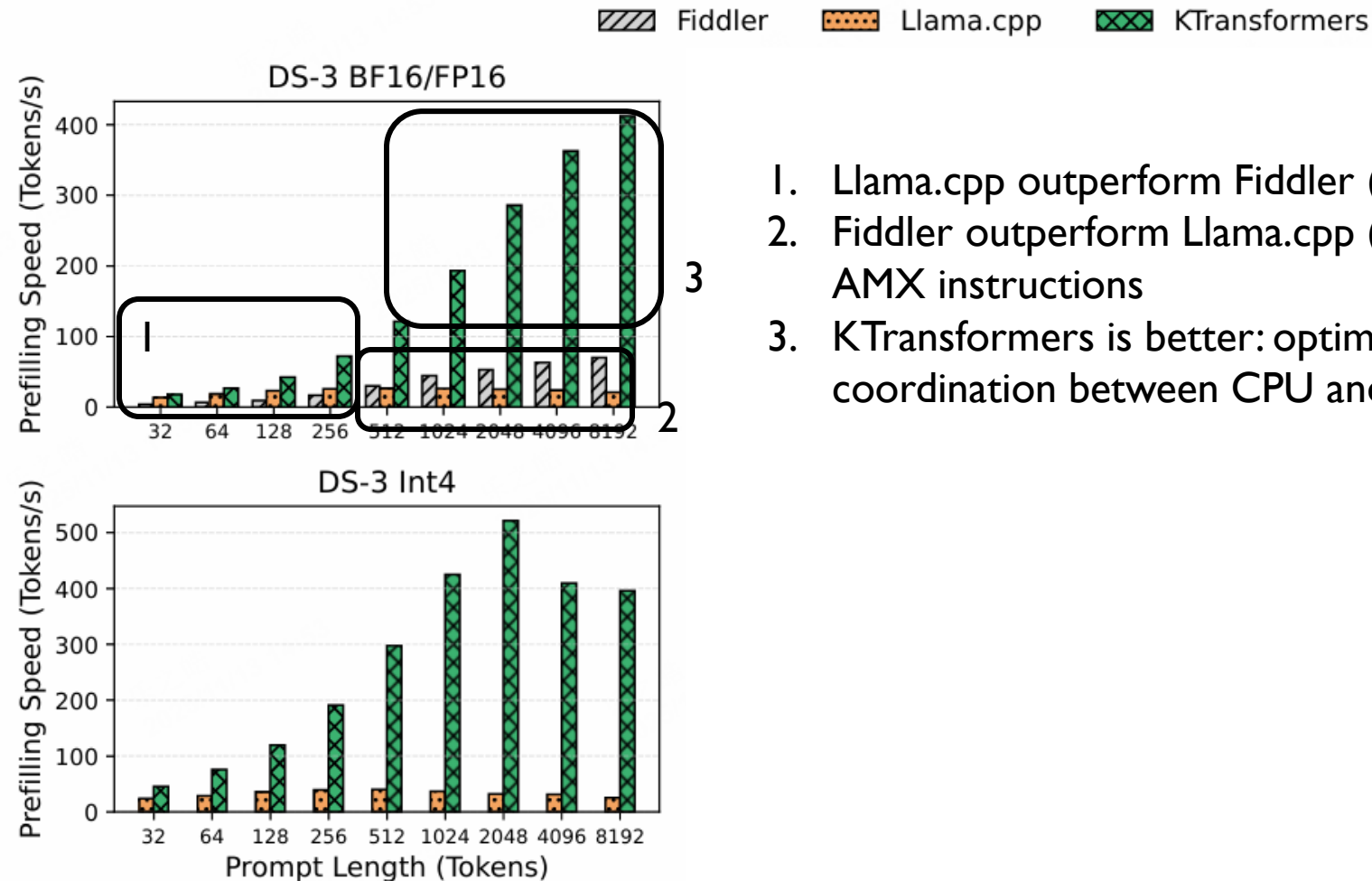
❑ **Baselines:**

   ❖ **Fiddler, Llama.cpp**

# Evaluation – End2End Performance



Comparison of prefilling speed between KTransformers and the state-of-the-art baselines
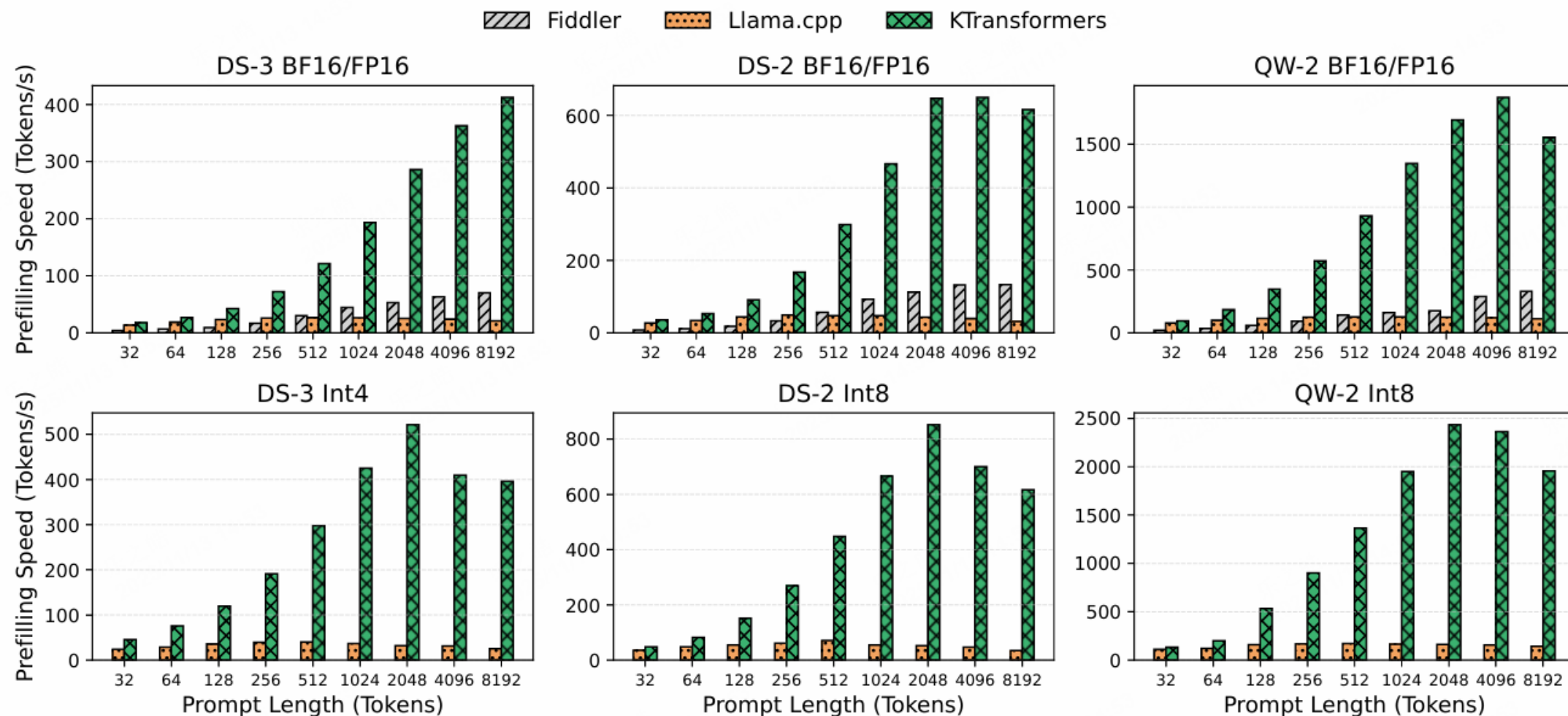
# Evaluation – End2End Performance



1. Llama.cpp outperform Fiddler (short prompt): superior fusion ops
2. Fiddler outperform Llama.cpp (long prompt): better utilization of AMX instructions
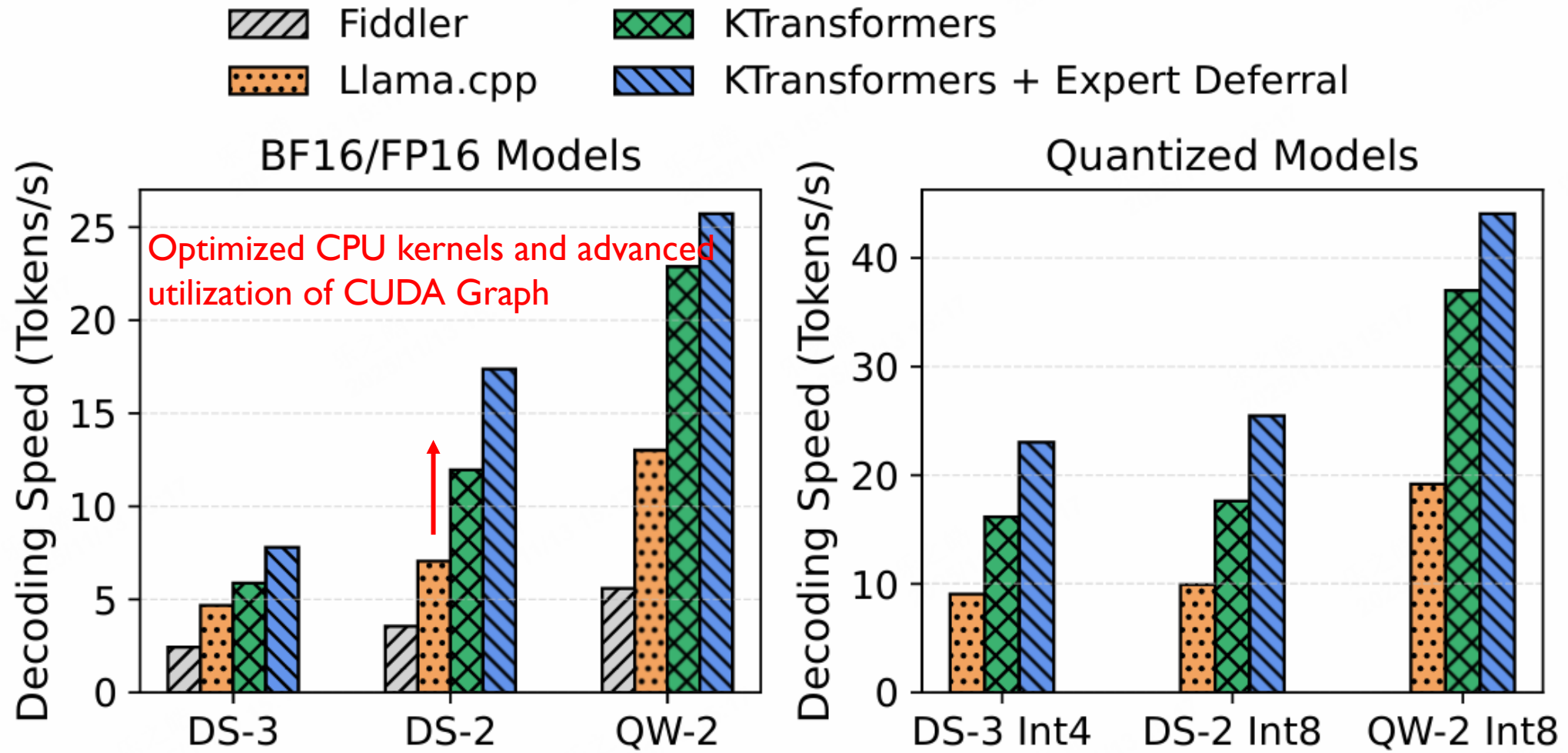3. KTransformers is better: optimized CPU kernels and improved coordination between CPU and GPU

Comparison of prefilling speed between KTransformer and the state-of-the-art baselines
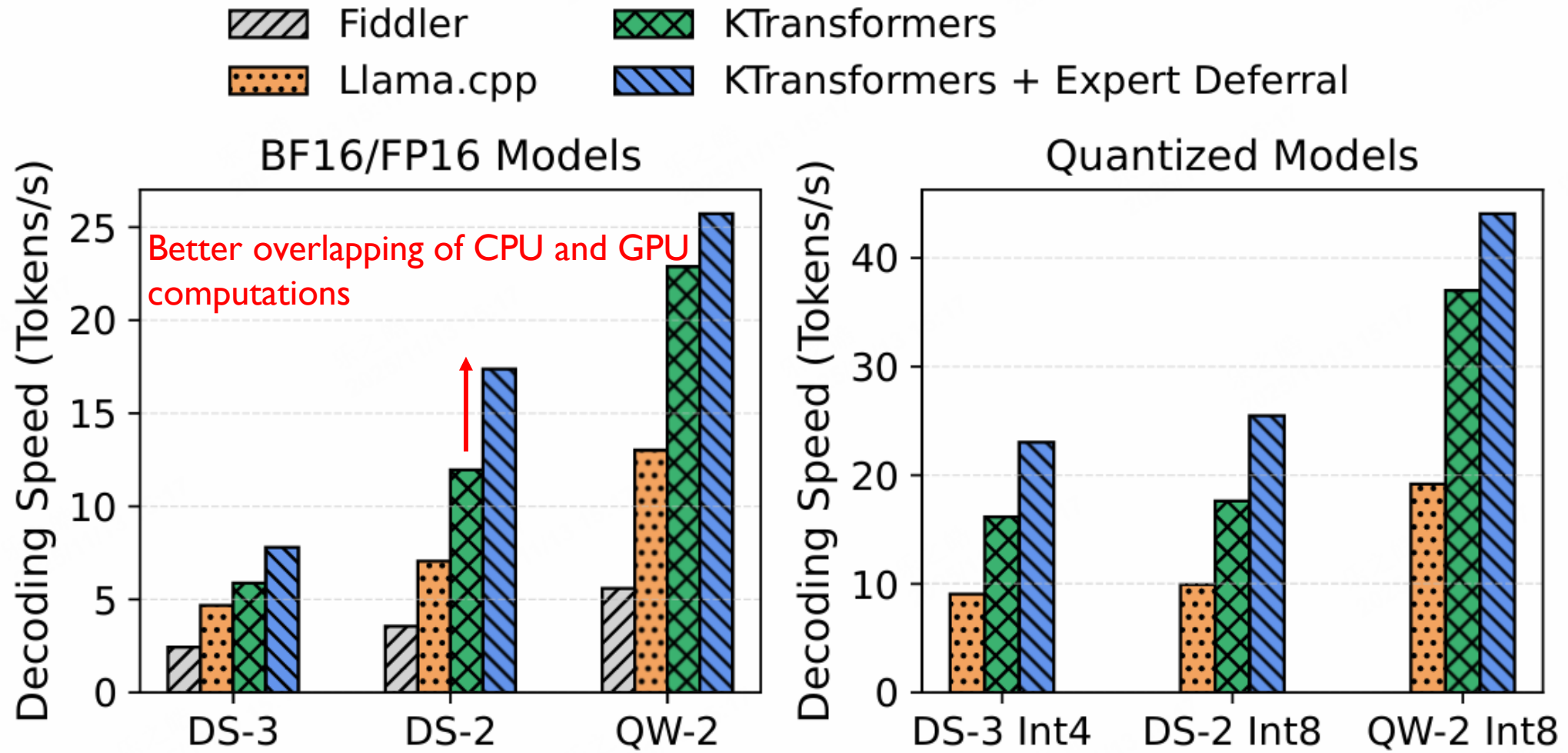
# Evaluation – End2End Performance



Comparison of prefilling speed between KTransformers and the state-of-the-art baselines

# Evaluation – End2End Performance



Comparison of decoding speed between KTransformers and the state-of-the-art baselines

# Evaluation – End2End Performance



Comparison of decoding speed between KTransformers and the state-of-the-art baselines

# Evaluation – Expert Deferral

KTransformers keeps good accuracy with the number of deferred experts less than 6.

| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| Coding | 68.8 | -0.1% | +0.1% | +0.4% | +1.1% | +0.7% | -0.6% | -4.7% | -11.2% |
| Data Analysis | 57.8 | +0.2% | +0.3% | +0.4% | +0.0% | +0.6% | +1.0% | +1.0% | -2.4% |
| Instruction Following | 82.6 | +0.3% | -0.0% | +0.2% | -0.2% | +0.1% | -0.2% | -0.5% | -1.8% |
| Language | 46.3 | -0.0% | -0.1% | +1.2% | +1.0% | +0.4% | +0.5% | +0.9% | +0.4% |
| Math | 71.7 | +0.2% | +0.4% | +0.1% | -0.1% | -0.2% | -1.6% | -4.7% | -13.4% |
| Reasoning | 71.5 | +0.4% | -0.2% | -0.4% | -0.3% | -0.9% | -1.6% | -2.2% | -9.4% |
| Average | 66.4 | +0.2% | +0.1% | +0.2% | +0.2% | +0.1% | -0.5% | -1.9% | -6.7% |

Number of Deferred Experts

DeepSeek-V3 accuracy on LiveBench under Expert Deferral

# Evaluation – Expert Deferral

KTransformers keeps good accuracy with Expert Deferral compared to with Expert Skipping.



| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| Coding | 68.8 | -0.1% | +0.1% | +0.4% | +1.1% | +0.7% | -0.6% | -4.7% | -11.2% |
| Data Analysis | 57.8 | +0.2% | +0.3% | +0.4% | +0.0% | +0.6% | +1.0% | +1.0% | -2.4% |
| Instruction Following | 82.6 | +0.3% | -0.0% | +0.2% | -0.2% | +0.1% | -0.2% | -0.5% | -1.8% |
| Language | 46.3 | -0.0% | -0.1% | +1.2% | +1.0% | +0.4% | +0.5% | +0.9% | +0.4% |
| Math | 71.7 | +0.2% | +0.4% | +0.1% | -0.1% | -0.2% | -1.6% | -4.7% | -13.4% |
| Reasoning | 71.5 | +0.4% | -0.2% | -0.4% | -0.3% | -0.9% | -1.6% | -2.2% | -9.4% |
| Average | 66.4 | +0.2% | +0.1% | +0.2% | +0.2% | +0.1% | -0.5% | -1.9% | -6.7% |

Number of Deferred Experts

DeepSeek-V3 accuracy on LiveBench under Expert Deferral



| | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
|---|---|---|---|---|---|---|---|---|---|
| Coding | 68.8 | -0.2% | +1.5% | -0.1% | -3.2% | -5.2% | -15.9% | -46.7% | -100.0% |
| Data Analysis | 57.8 | +0.2% | -1.7% | -2.4% | -6.7% | -18.8% | -38.5% | -52.3% | -67.8% |
| Instruction Following | 82.6 | +0.8% | +0.2% | -0.0% | +0.1% | -0.3% | -2.2% | -8.3% | -75.2% |
| Language | 46.3 | -0.2% | +1.4% | +1.5% | +1.3% | -1.3% | -3.0% | -17.5% | -93.3% |
| Math | 71.7 | +0.3% | -0.4% | -0.6% | -2.5% | -6.2% | -11.9% | -30.9% | -95.9% |
| Reasoning | 71.5 | +0.2% | -0.2% | -1.8% | -3.0% | -5.3% | -11.2% | -20.0% | -100.0% |
| Average | 66.4 | +0.2% | +0.1% | -0.6% | -2.3% | -5.9% | -13.3% | -28.6% | -88.7% |

Number of Skipped Experts

DeepSeek-V3 accuracy on LiveBench under Expert Skipping

# Evaluation - Breakdown

+v: AVX-512 instructions
+m: AMX instructions
+d: dynamic work scheduling
+n: NUMA-aware tensor parallelism
+c: CUDA Graph

Breakdown of prefill phase

Breakdown of decode phase

# Evaluation - Breakdown

1. AMX better in prefill: prefill is computation heavy.
2. Dynamic work scheduling is more efficient in prefill: decode is more load balanced.
3. NUMA-aware is efficient in decode phases: decode is more memory bound.
4. CUDA Graph is efficient in decode: in prefill phase, the overhead of CUDA launch is amortized into a large number of tokens.
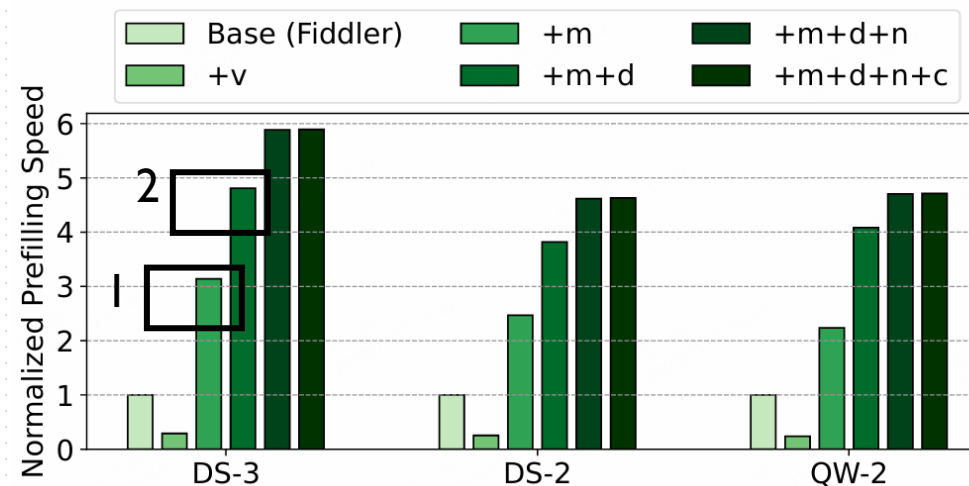
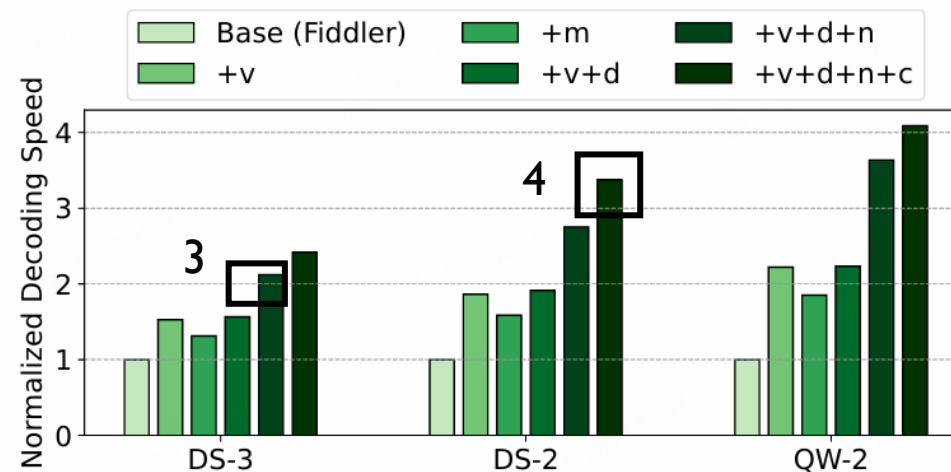+v: AVX-512 instructions
+m: AMX instructions
+d: dynamic work scheduling
+n: NUMA-aware tensor parallelism
+c: CUDA Graph



Breakdown of prefill phase



Breakdown of decode phase

# Outline

❑ **Background**

❑ **Motivation**

❑ **Challenge**

❑ **Design**

❑ **Evaluation**

❑ **Conclusion**

# Conclusion

□ **KTransformers, a system that enables efficient local inference for large MoE models on hybrid CPU/GPU platforms.**

  ❖ **Optimize CPU ops by combining AMX-optimized kernels for better utilization of CPU.**

  ❖ **Use CPU-GPU asynchronous scheduling and NUMA-aware TP for better CPU-CPU/GPU coordination.**

  ❖ **Use the Expert Deferral strategy to maximize the utilization of hardware.**