

FlexInfer: Flexible LLM Inference Using CPU Computation



Seonjin Na¹, Geonhwa Jeong², Byung Hoon Ahn³, Aaron Jezghani¹, Jeffery Young¹, Christopher J. Huges⁴, Tushar Krishna¹, Hyesoon Kim¹

UC San Diego

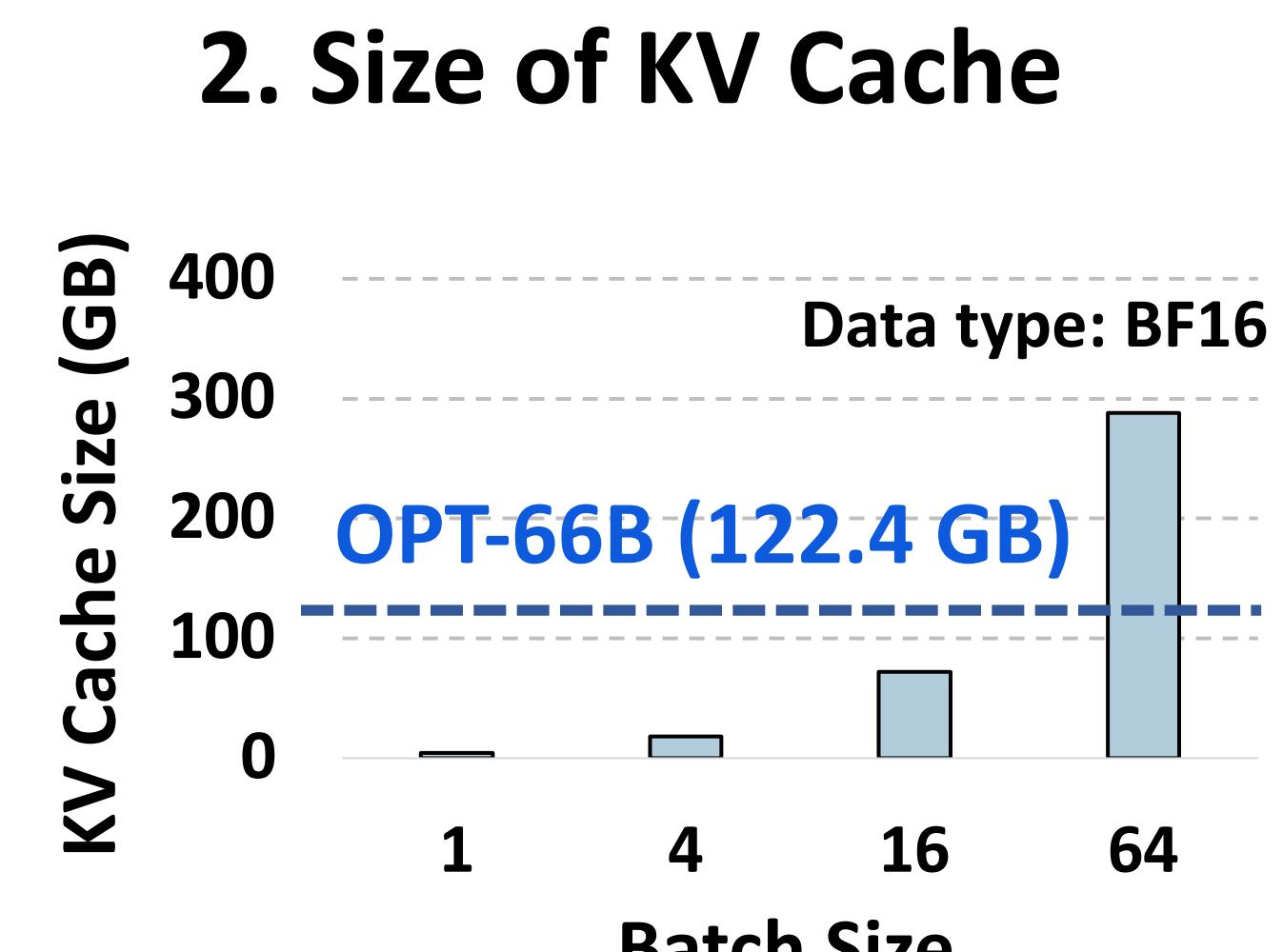
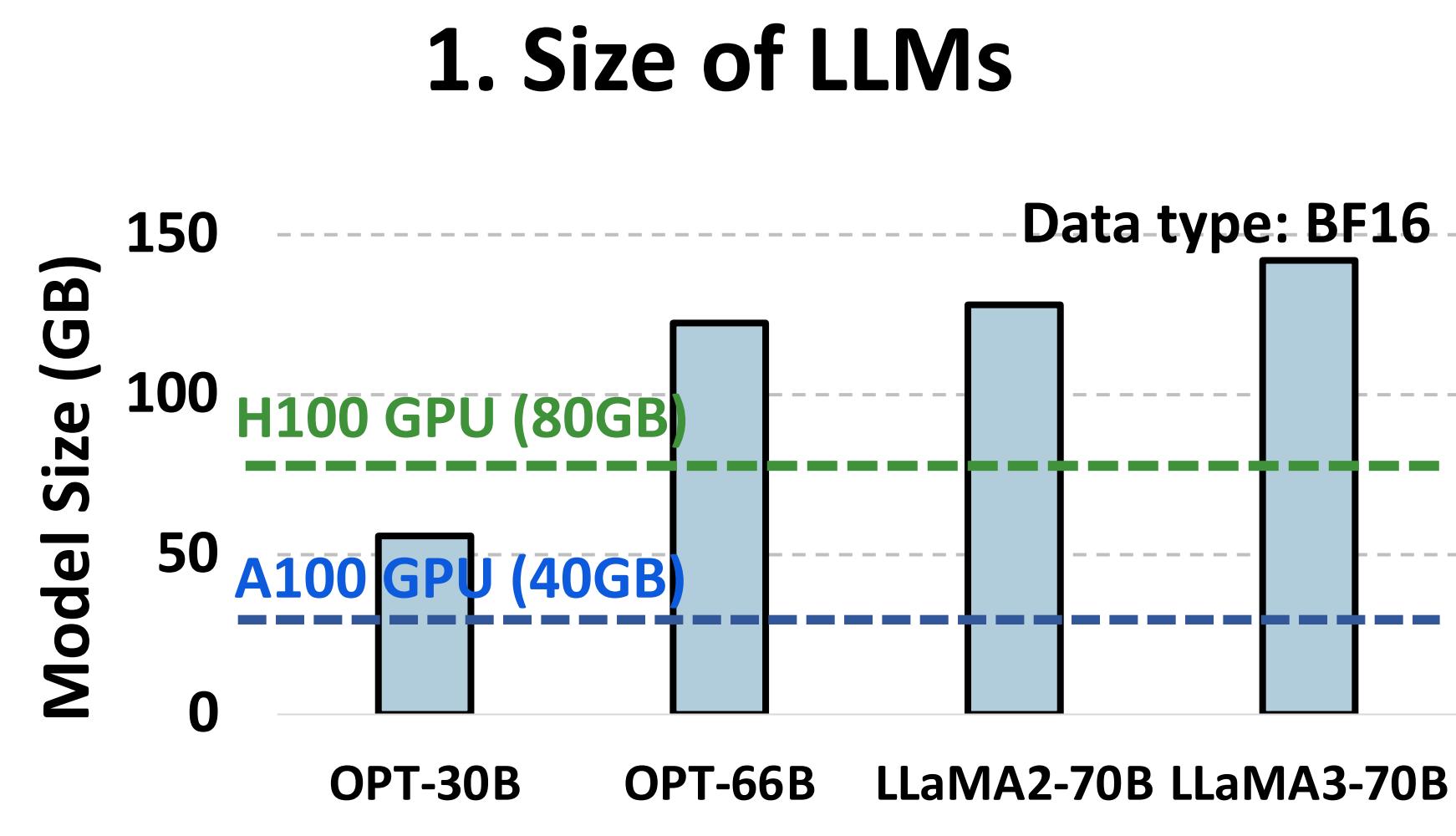


¹Georgia Institute of Technology, ²Meta, ³University of California San Diego, ⁴Intel Labs



Problem

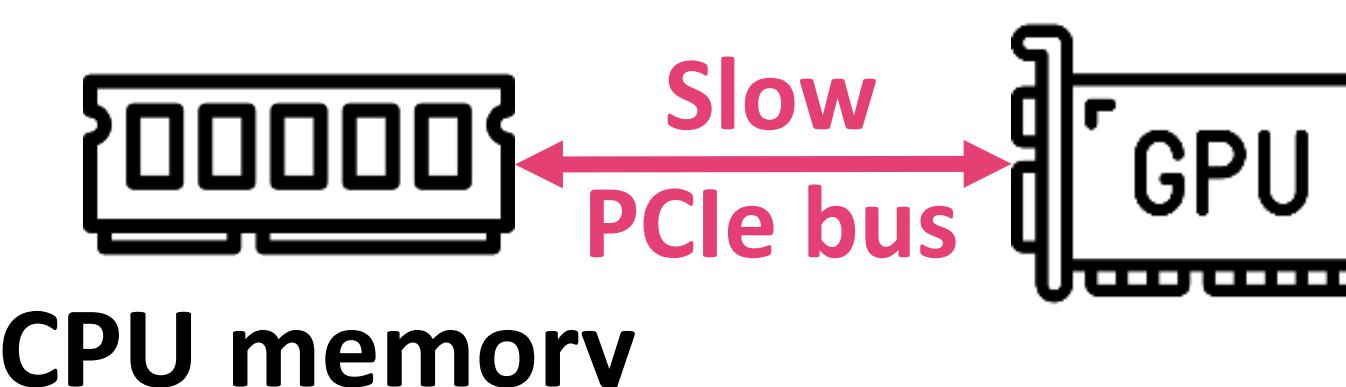
LLM Inference requires significant amount of memory



Offloading-based LLM Inference suffer from PCIe data transfer overhead

Offloading-based LLM Inference

Offload data to CPU memory
(model weights, act., KV cache)



H100 with Offloading Exec. Breakdown

(Input length: 1024, output length: 32)

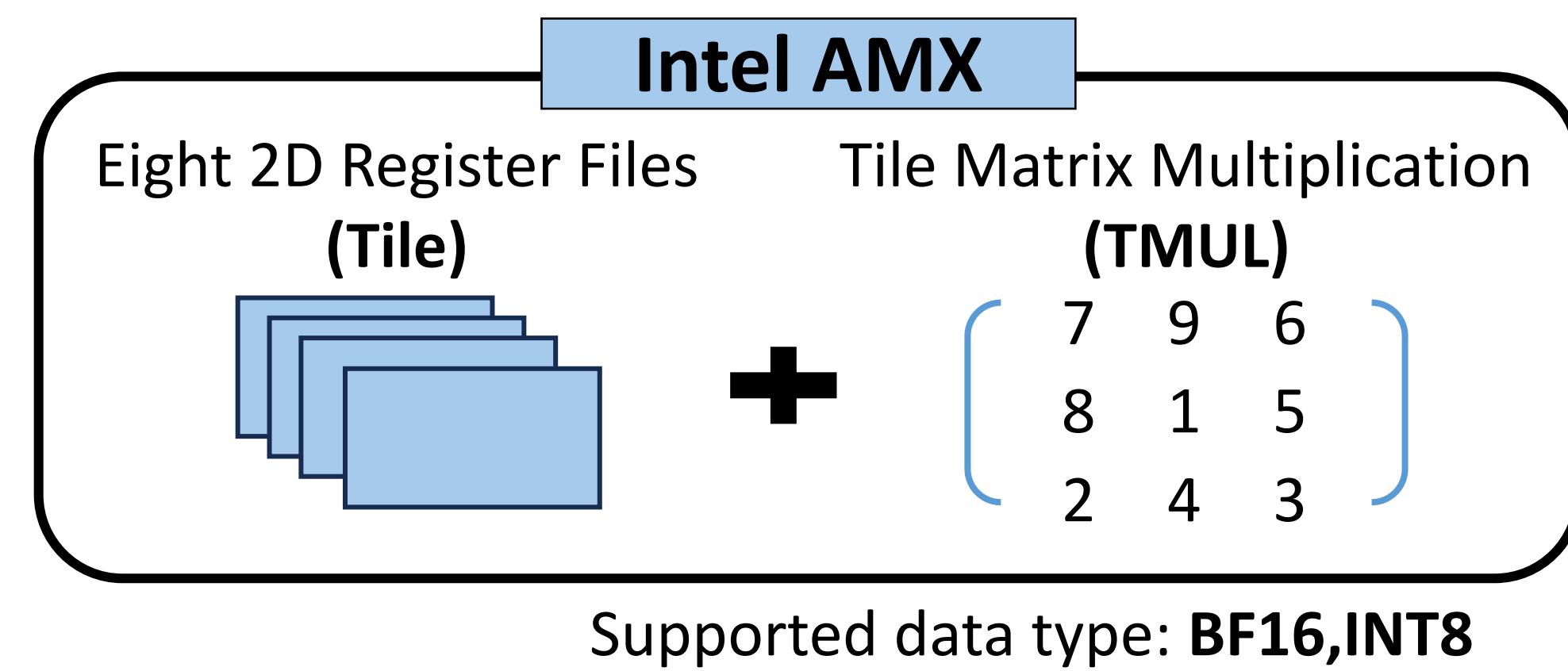
Batch Size	Data Load (PCIe)	GPU Compute
1	95.8	4.2
4	94.4	5.6
16	89.3	10.7

Significant data transfer overhead

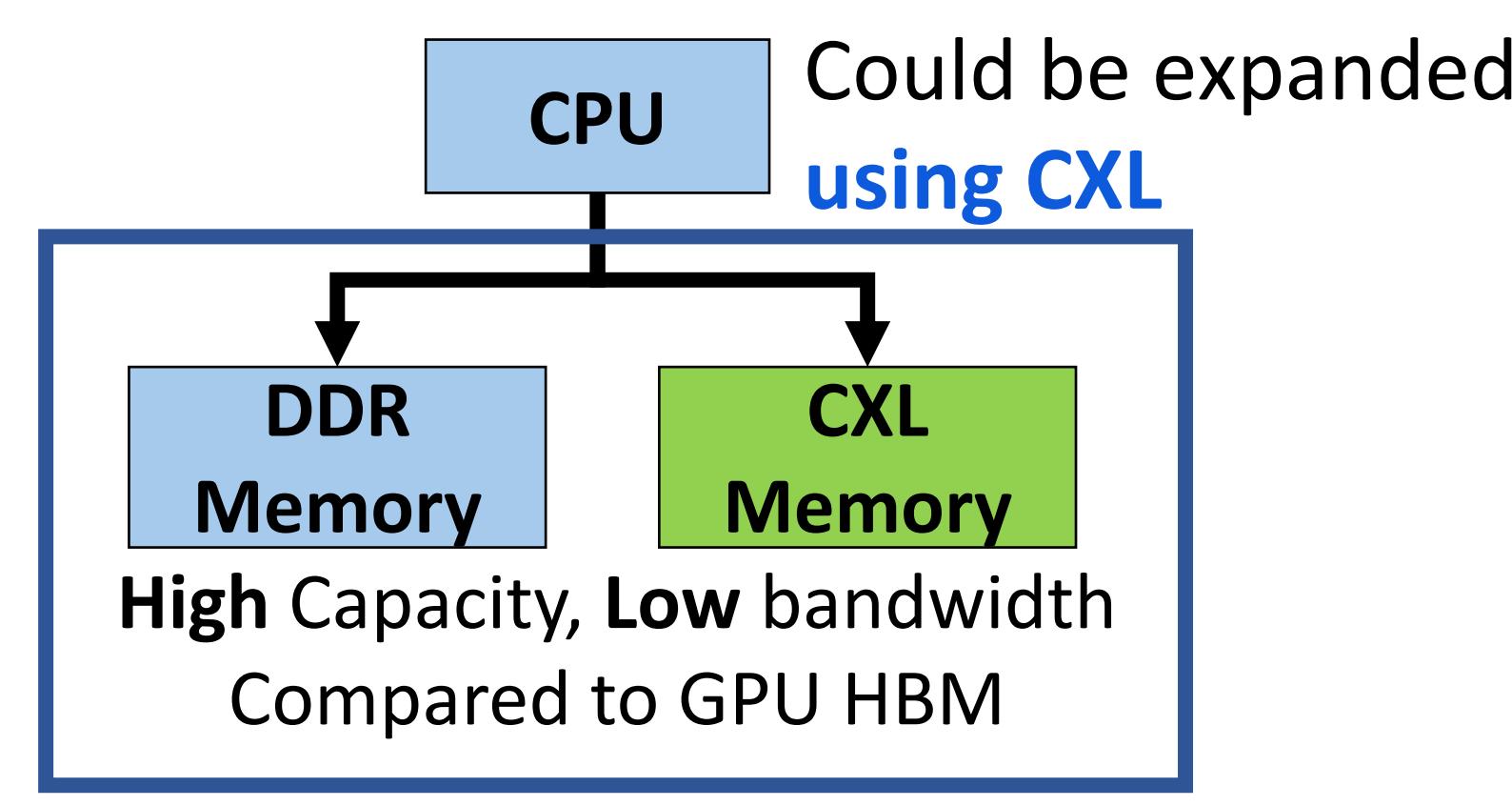
Motivation

Key Opportunities in Recent CPUs

1. Dedicated GEMM Accelerator



2. Larger Memory Capacity



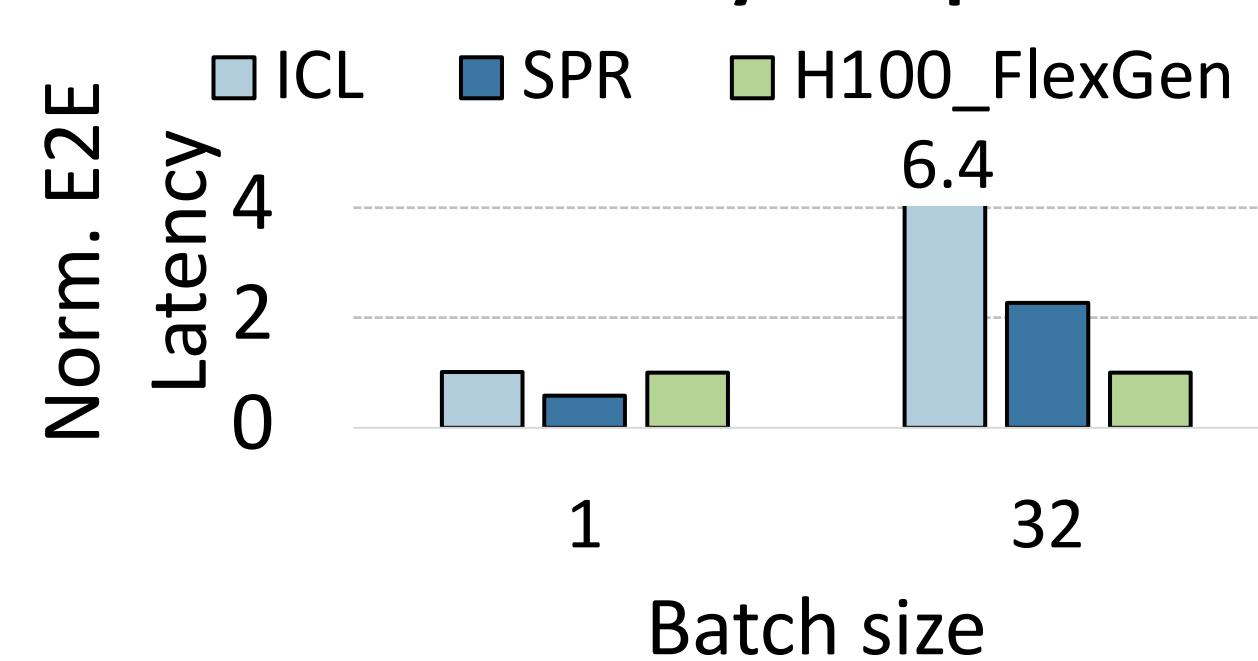
CPU Computation vs GPU with Offloading

CPU Compute Limited in Prefill with longer seq. larger batch, Beneficial in Decode

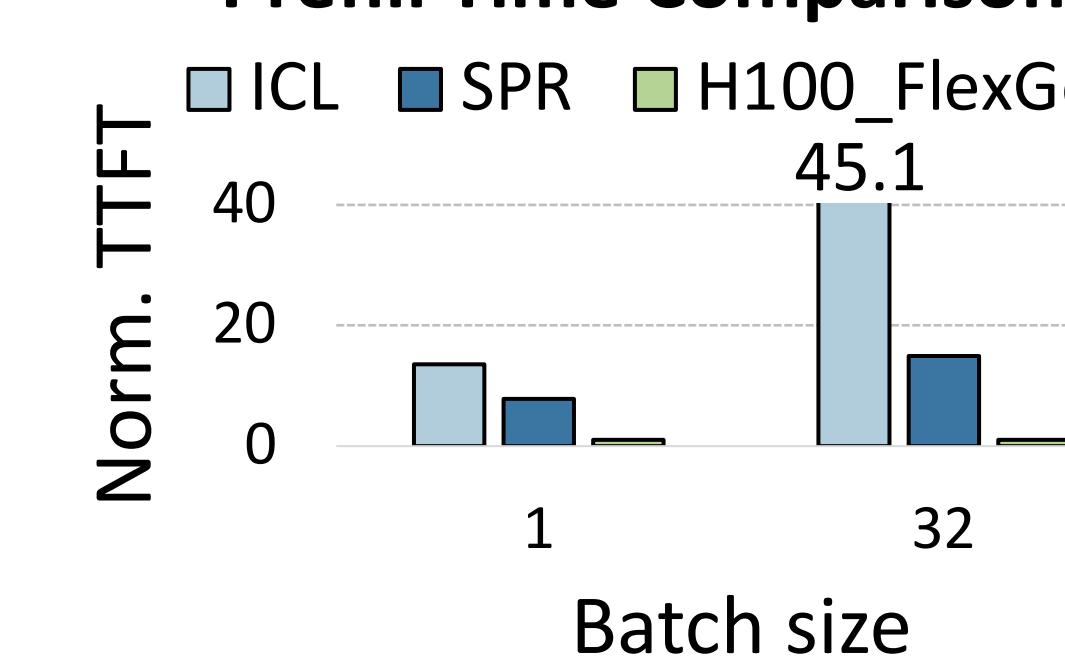
Config: ICL CPU (without AMX), SPR CPU (with AMX), H100 FlexGen

Model: OPT-66B Input length: 1024, output length: 32

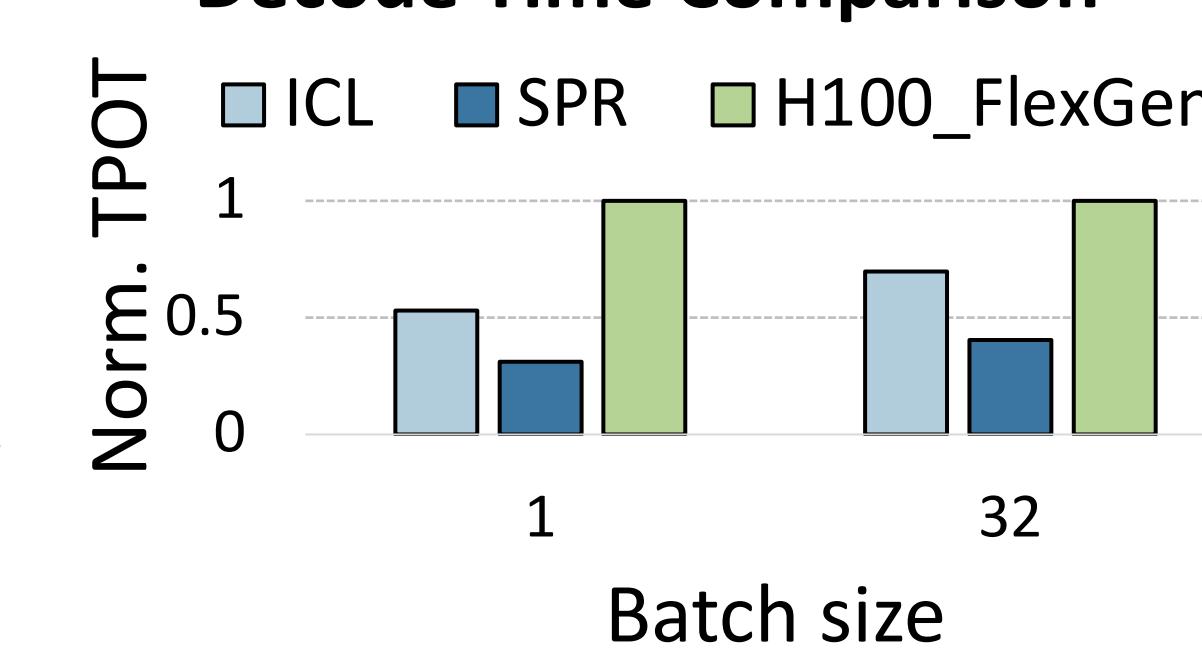
End-to-End Latency Comparison



Prefill Time Comparison



Decode Time Comparison

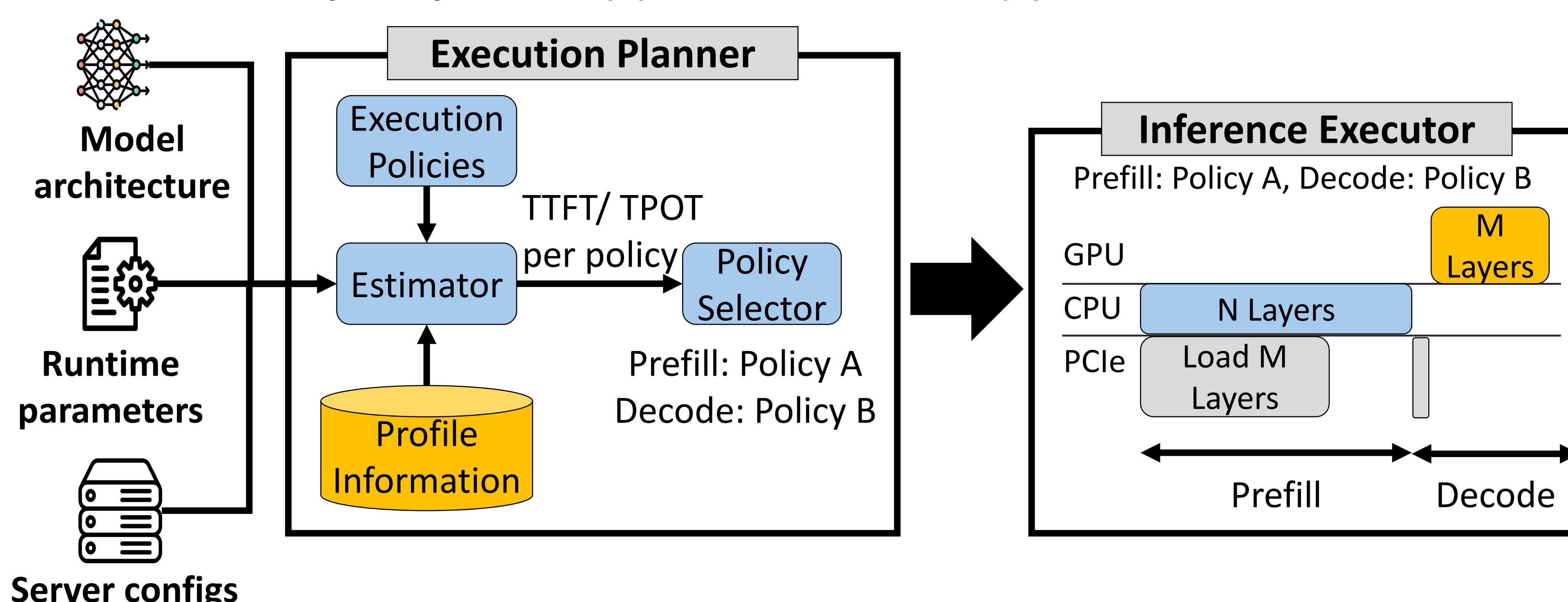


FlexInfer

Overview of FlexInfer

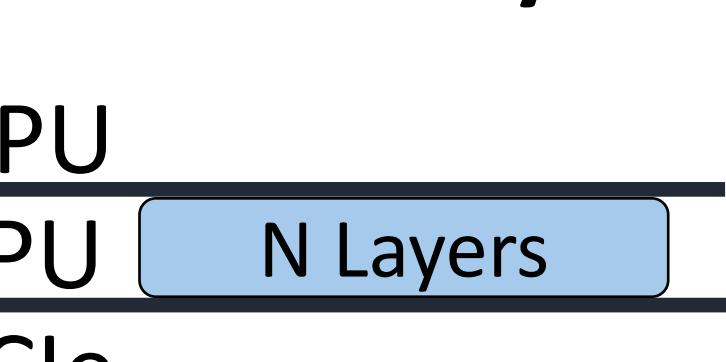
Goal: FlexInfer dynamically select execution policies based on perf. estimator

Two key components: (1) Execution Planner (2) Inference Executor

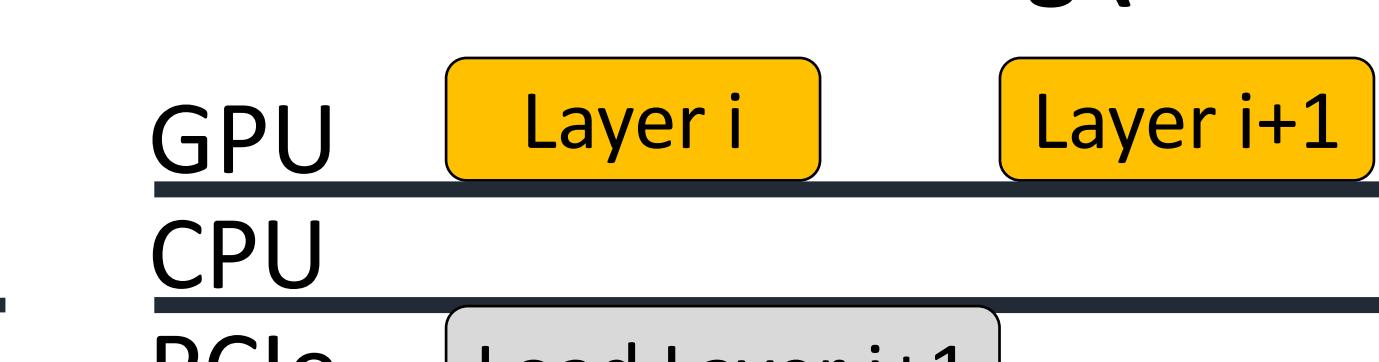


Baseline Execution Policies in FlexInfer and Concrete Examples

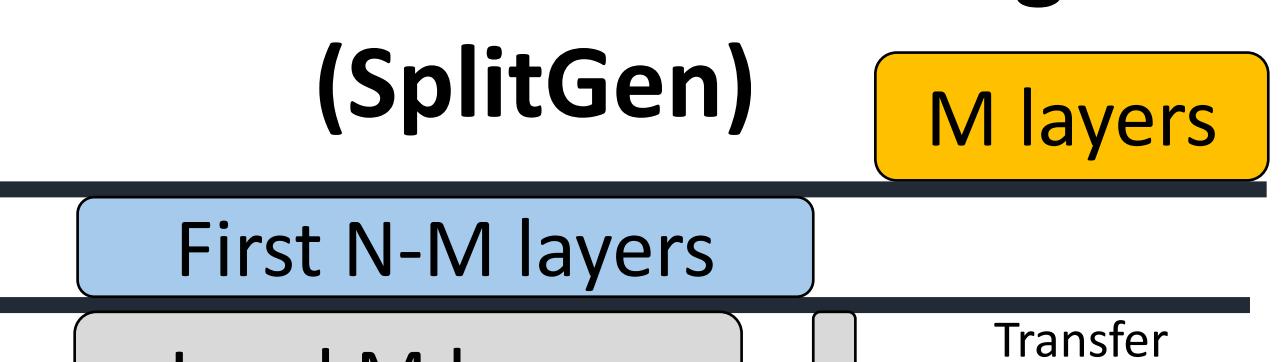
1. CPU-Only



2. GPU with Offloading (FlexGen)



3. CPU-GPU Partitioning (SplitGen)



- + No need to transfer data
- Limited compute throughput for longer seq. & larger batch

- Efficiently process longer seq. & larger batch in prefill
- Slowdown due to data transfer

- Low data transfer volume, exploit both CPU and GPU
- Slowdown due to CPU compute bottleneck

Scenario1

Long input seq.
Large batch

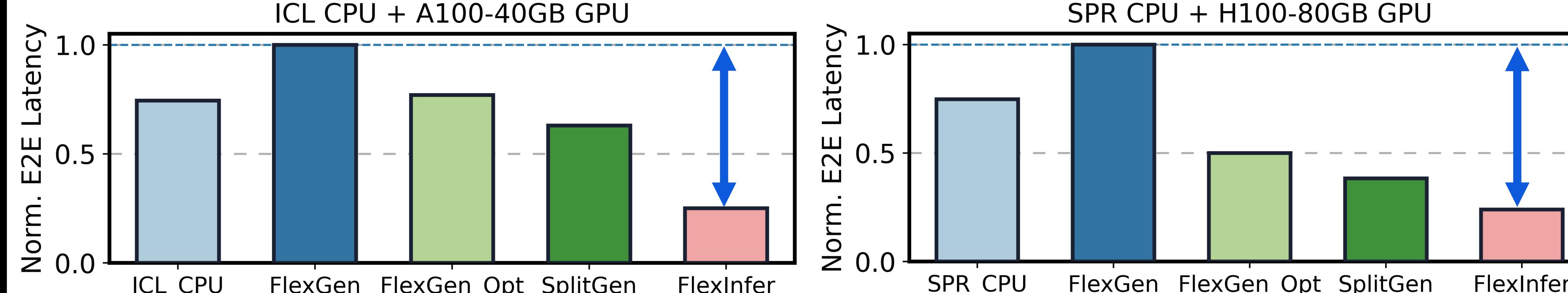
Scenario2

Short input seq.
Smaller batch

Evaluation Results and Summary

Performance Comparison Results (Input/Output length : 512/32)

FlexInfer reduces end-to-end latency by 75%, 76% on two different servers



Summary

Problem

- LLM inference requires significant memory exceeding recent GPUs memory capacity
- Offloading-based LLM inference suffers from significant data transfer overhead

Key observation & Idea

- CPU computation offers potential but struggles with compute-heavy prefill phase
- FlexInfer dynamically select execution policies based on HW configs/ runtime params

Evaluation results

- FlexInfer reduces inference latency by 75%, 76% on average across different servers