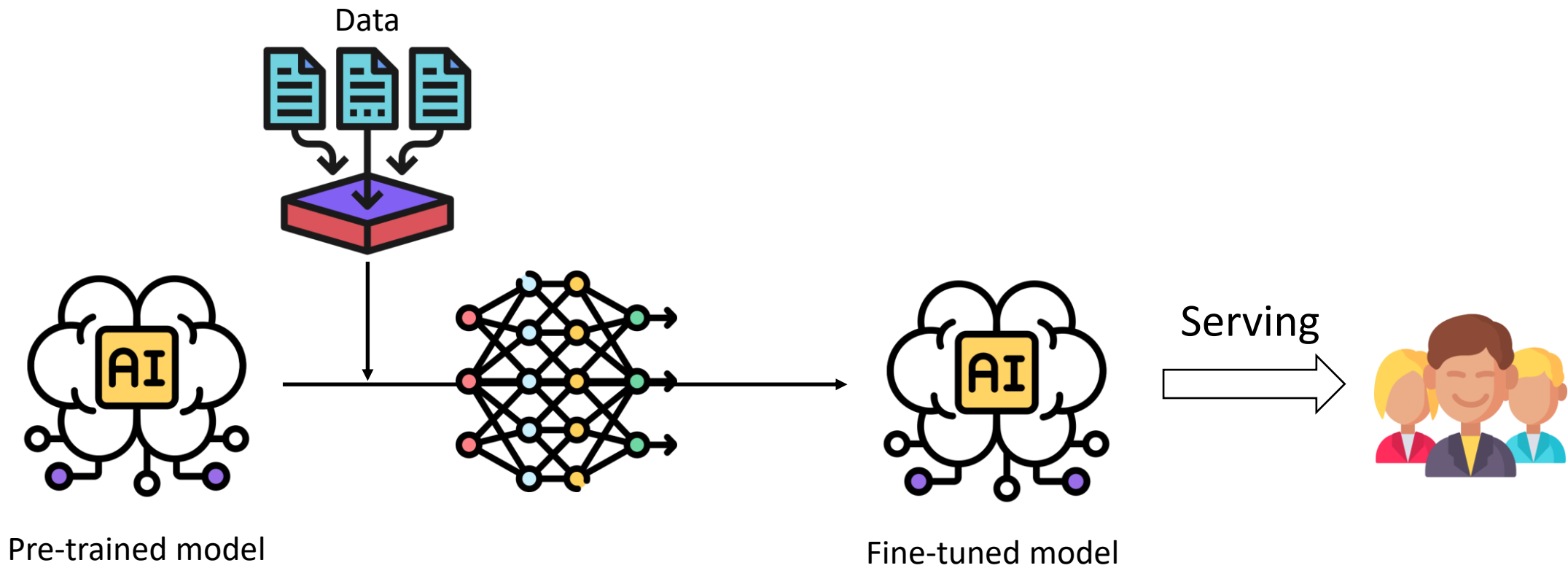# UNIVERSITY OF ILLINOIS
## URBANA-CHAMPAIGN

# CS 598
# AI Efficiency: Systems and Algorithms
## Overview & Key Challenges in AI Model Serving

Minjia Zhang

Computer Science Department

# What is Model Serving?



Data

Pre-trained model

Fine-tuned model

Serving

# Serving Scenario 1: Online ChatBot



What are some fun places to visit at Urbana Champaign?

✓ Searching for: **fun places to visit at Urbana Champaign**
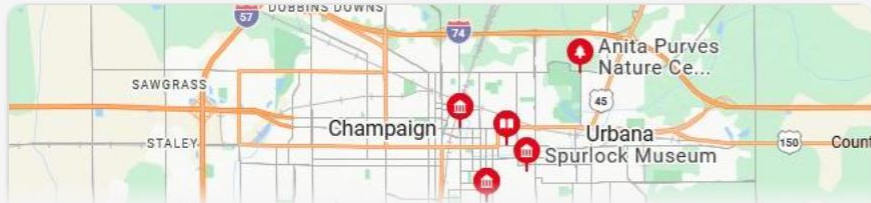
✓ Generating answers for you...

There are many fun places to visit at Urbana Champaign, depending on your interests and preferences. Here are some of the most popular ones:

- If you love nature and animals, you might enjoy visiting the **Anita Purves Nature Center** [1], where you can explore trails, feed goats, and learn about local wildlife.

- If you are into art and culture, you might want to check out the **Krannert Art Museum** [2], which has a diverse collection of artworks from various regions and periods, as well as exhibitions and events.

- If you are looking for some entertainment and history, you might like the **Virginia Theatre** [3], which is a restored historic venue that hosts concerts, movies, and shows.

- If you are feeling adventurous and sporty, you might have fun at the **University of Illinois Ice Arena**, where you can skate, play hockey, or watch games.

These are just some of the fun places to visit at Urbana Champaign. You can find more information and reviews on Tripadvisor or Bing. I hope you have a great time exploring the city! 😊

**Learn more**  ∨   1 🌐 experiencecu.org   2 🏛 kam.illinois.edu   3 🔍 bing.com

👍 👎 📋 ↓ ↗                                                          2 of 30 ●

# Serving Scenario 2: Online Image Generation

# Serving Scenario 3: Online Q&A

what are famous quote from shakespeare

**Google**

All   Images   News   Shopping   Videos   More      Settings   Tools

About 7,370,000 results (0.88 seconds)

William Shakespeare / Quotes

*Be not afraid of greatness: some are born great, some achieve greatness, and some have greatness thrust upon them.*

*To thine own self be true, and it must follow, as the night the day, thou canst not then be false to any man.*

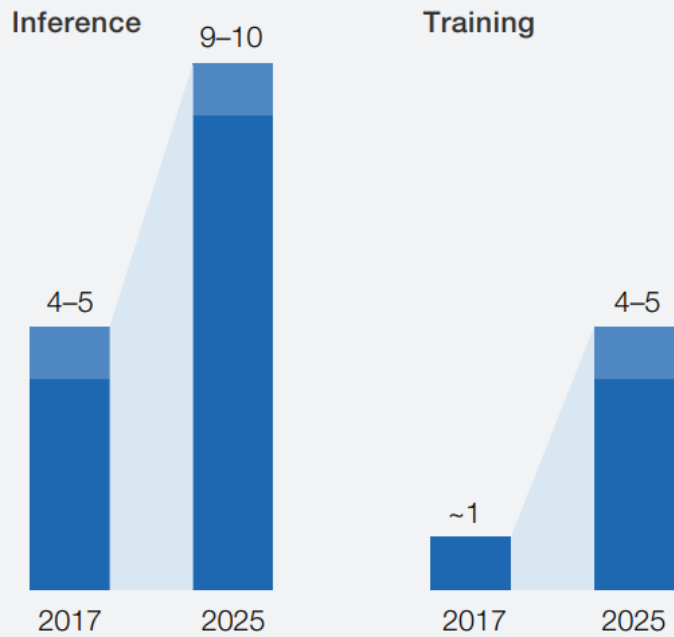*The course of true love never did run smooth.*

*Feedback*

- Question and Answer Scenario
  - Direct answer not a list of webpages
  - Good quality of answer

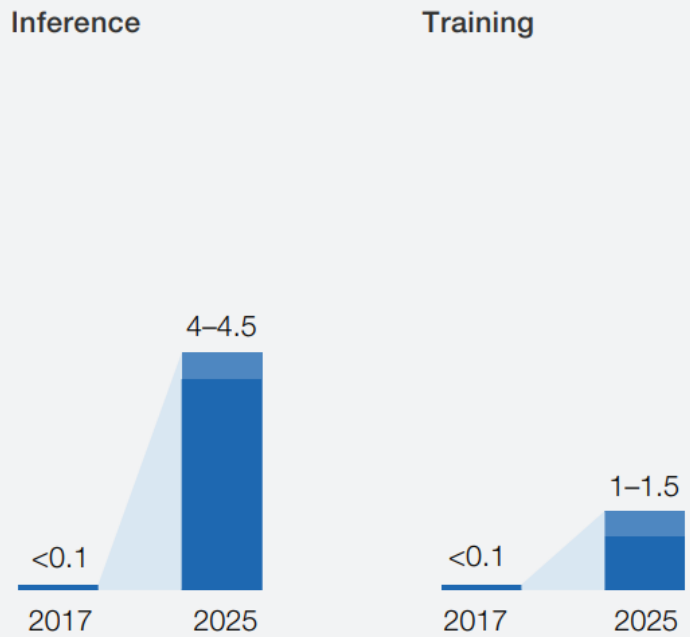# Training -> Inference



Exhibit 5

At both data centers and the edge, demand for training and inference hardware is growing.

Data center, total market, $ billion

Inference — 9–10

Training

4–5

4–5

~1

2017   2025      2017   2025

Edge, total market, $ billion

Inference

Training

4–4.5

1–1.5

<0.1

<0.1

2017   2025      2017   2025

Source: Expert interviews; McKinsey analysis

6

# Inference Challenges

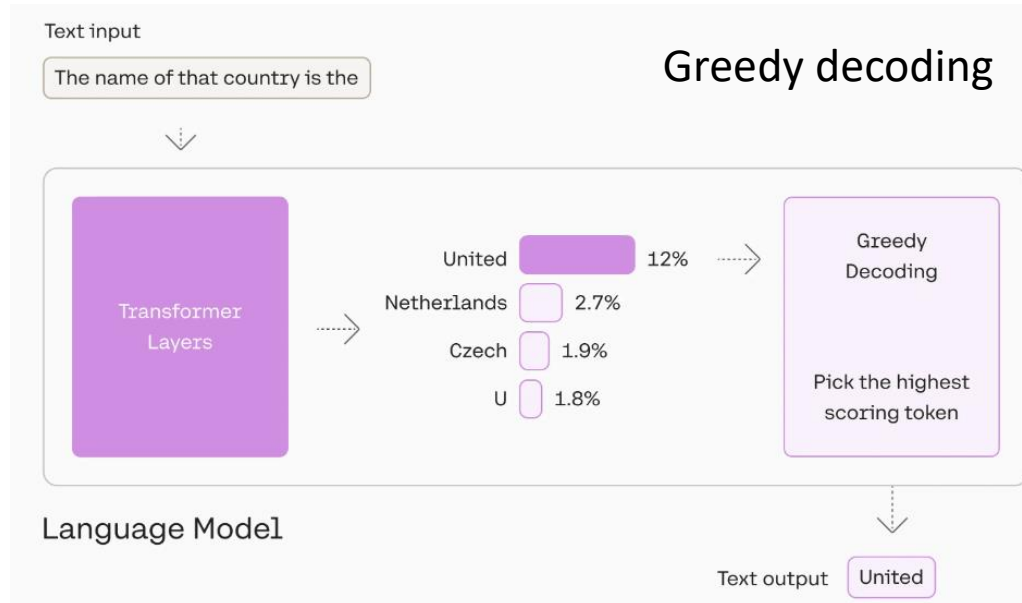|  | **Training** | vs | **Inference** |
|---|---|---|---|
| Runtime | Weeks or months | | Milliseconds or seconds |
| Challenges | TCO (Cost, Energy) | | TCO (Cost, Energy) |

**Speed** (LLM: token rates)
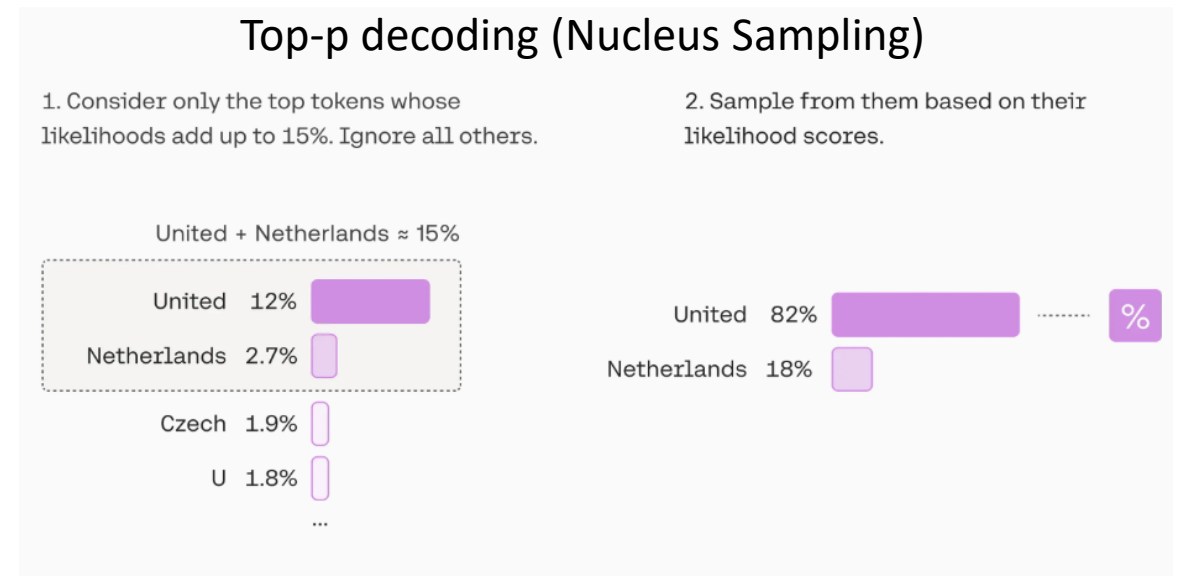**Model size**
- Parameter volume
- LLM: Context length

# LLM Autoregressive Generation

# Decoding Strategies



Greedy decoding

Top-k decoding

Top-p decoding (Nucleus Sampling)

Top-k & Top-p, Cohere

# Serving Challenge: Long Latency

- Long serving latency blocks deployment
- Support advance models while meeting latency SLA and saving cost

| DL Scenarios | Original Latency | Latency Target |
|---|---|---|
| Turing Prototype 2 | ~100ms | < 10ms |
| Turing Prototype 3 | ~107ms | < 10ms |
| Deep Query Document Similarity | 10~12ms for [query, 1 doc] x 33 docs | < 6ms |
| Malta Click Features | 10ms for [query, 1 passage] x 150 passages | < 5ms |
| Ads seq2seq model for query rewriting | ~51ms | < 5ms |

# Customized Kernels



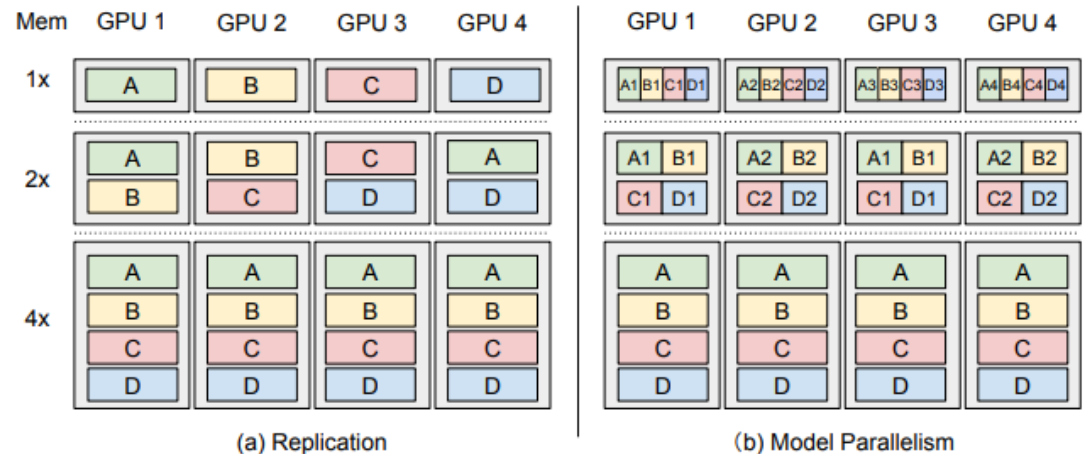ByteTransformer: A High-Performance Transformer Boosted for Variable-Length Inputs, 2023

DeepSpeed-Inference: enabling efficient inference of transformer models at unprecedented scale, SC 2022

Fig. 4: The zero padding algorithm.
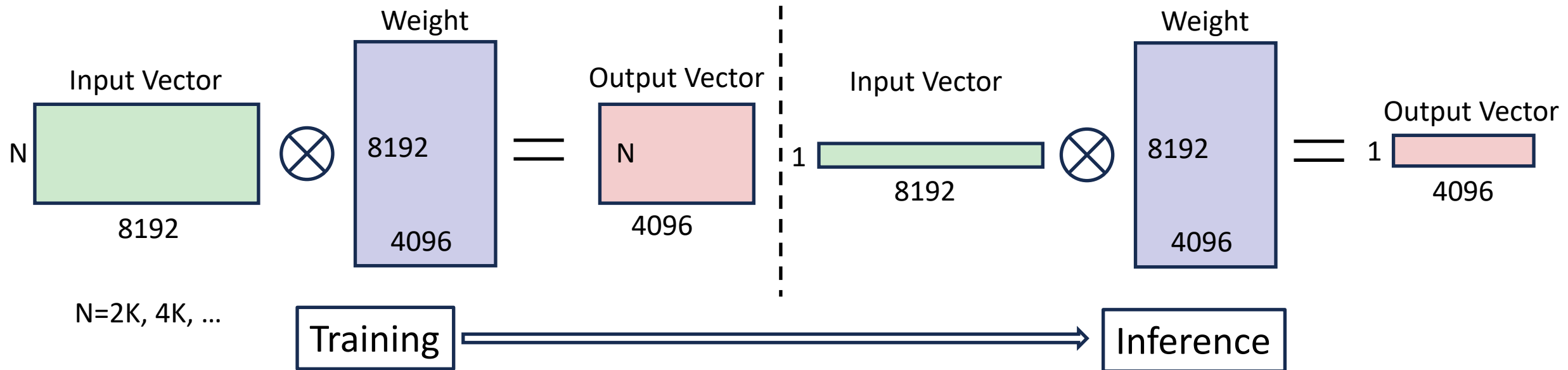
11

# Multi-GPU Inference via Partitioned Layouts



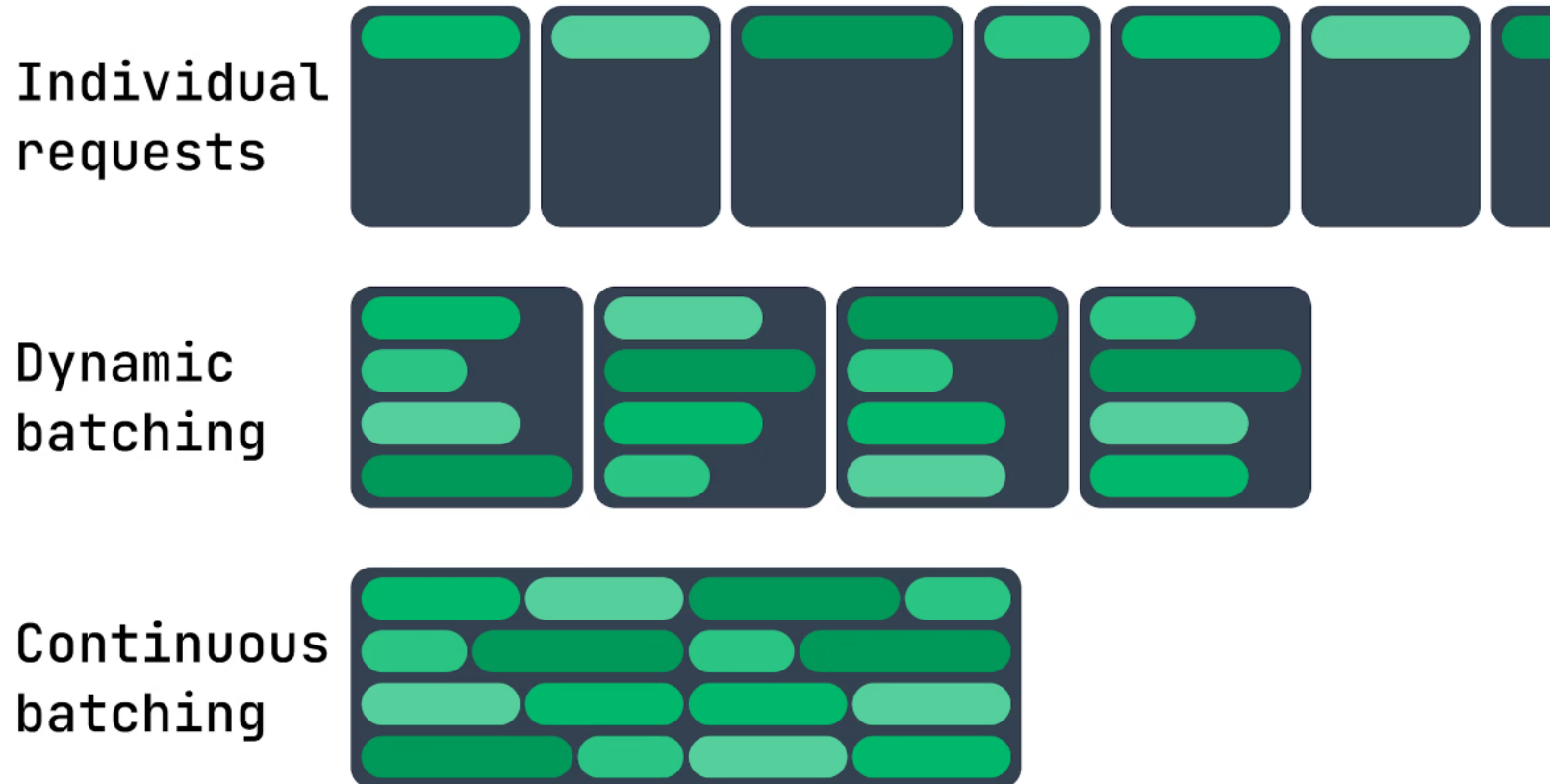AlpaServe: Statistical Multiplexing with Model Parallelism for Deep Learning Serving, OSDI 2023

Efficiently Scaling Transformer Inference, MLSys 2023

# Inference Challenge: Limited Parallelism

- Small batch size $\Longrightarrow$ Low data reuse
- Autoregressive generation $\Longrightarrow$ Sequential dependency

# Batching Strategies for LLM Inference



Individual requests

Dynamic batching

Continuous batching

Orca: A Distributed Serving System for Transformer-Based Generative Models, OSDI 2022
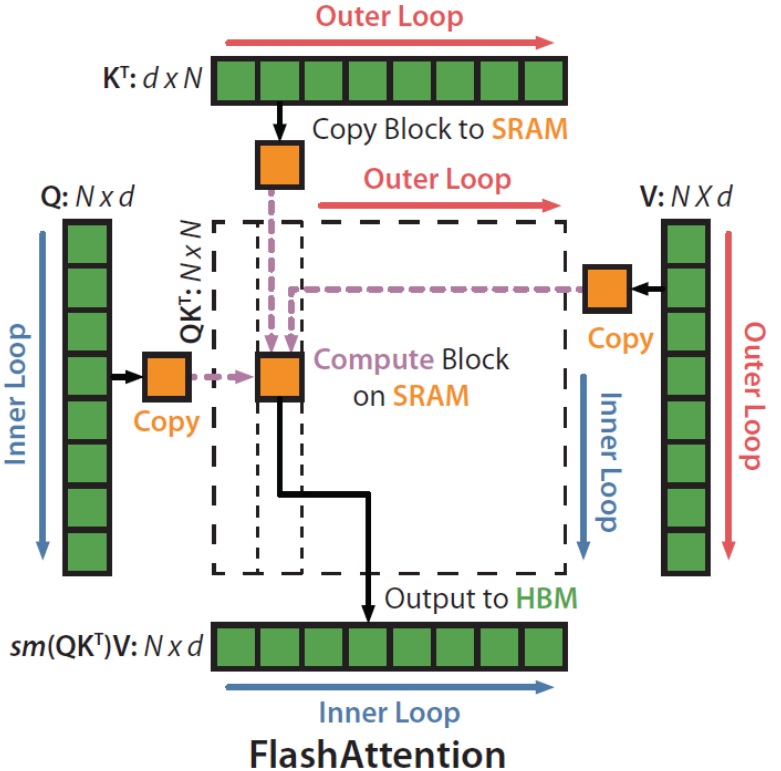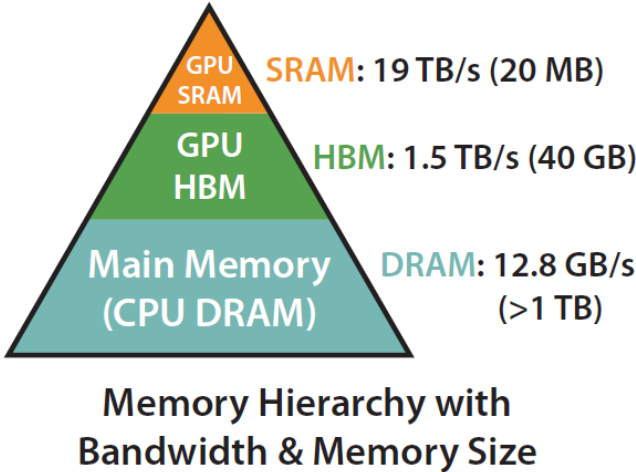
# Inference Challenge: Large Memory Footprint

- Model parameters
  - # Layers
  - # Hidden dim

- KV cache
  - Batch size
  - Sequence length
  - # Layers
  - # Hidden

- Activation and others

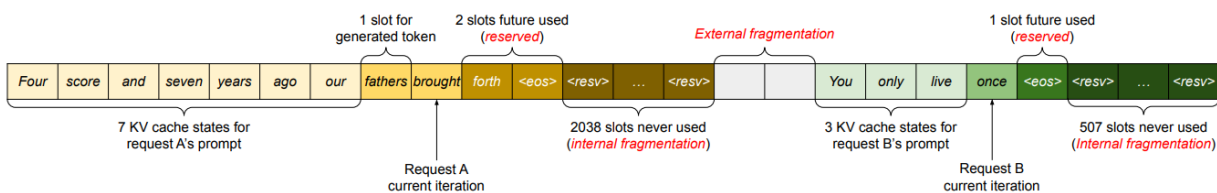| Parameters (26GB, 65%) | KV Cache (>30%) |
| | Others |

OPT-13B on A100 40 GB

Efficient Memory Management for Large Language Model Serving with PagedAttention, by Kwon et al., 2023
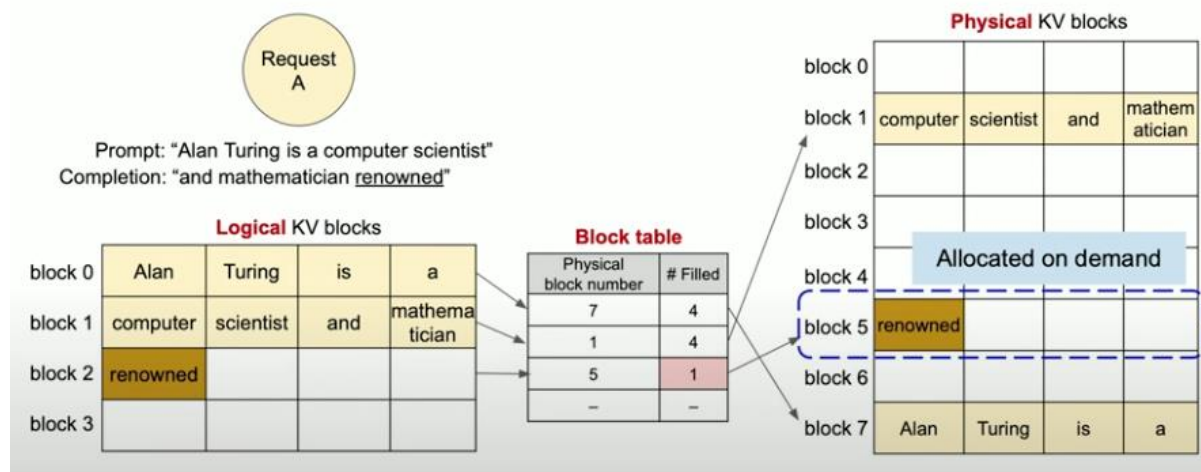
# FlashAttention



Fast and Memory-Efficient Exact Attention with IO-Awareness, 2023
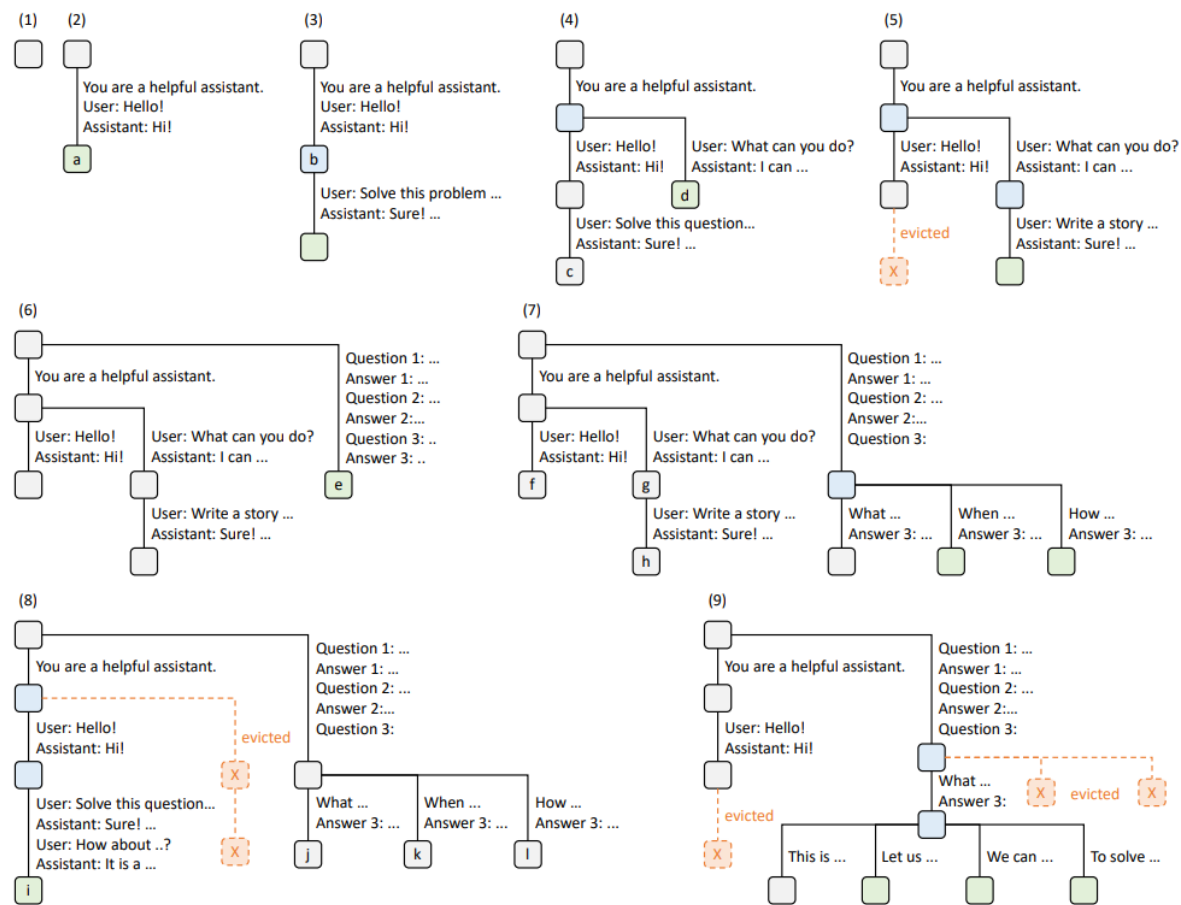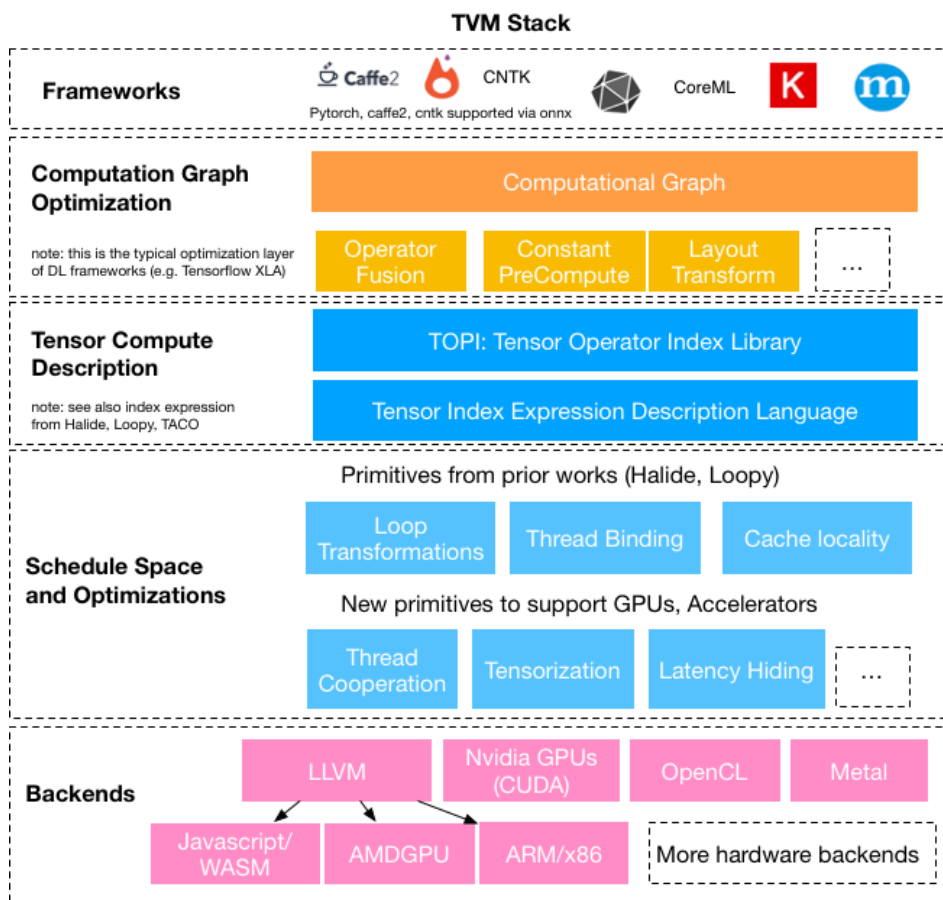
# PagedAttention

Efficient Memory Management for Large Language Model Serving with PagedAttention, 2023

17

# DL Compilation

# Class Related

- Slack Channel – Update display name

- Schedule – Newly joined students please send papers you are interested in presenting to me and the TA

- Presentation – Guidance on what to include (course website)

# QA