# Toward Efficient Inference for Mixture of Experts

**Haiyang Huang, Newsha Ardalani, Anna Sun, Liu Ke**
**Hsien-Hsin S. Lee, Shruti Bhosale, Carole-Jean Wu, Benjamin Lee**

NeurIPS, Dec 2024
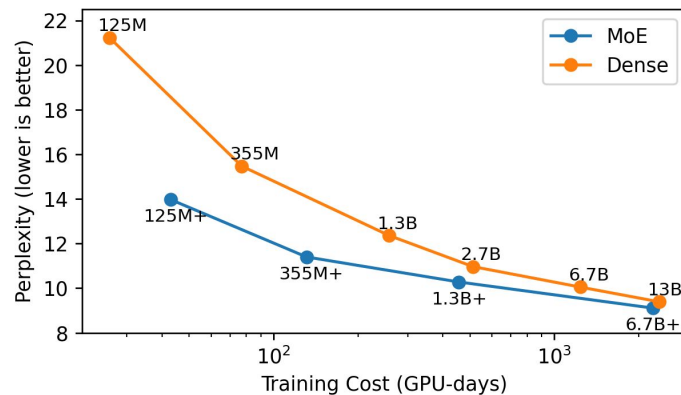
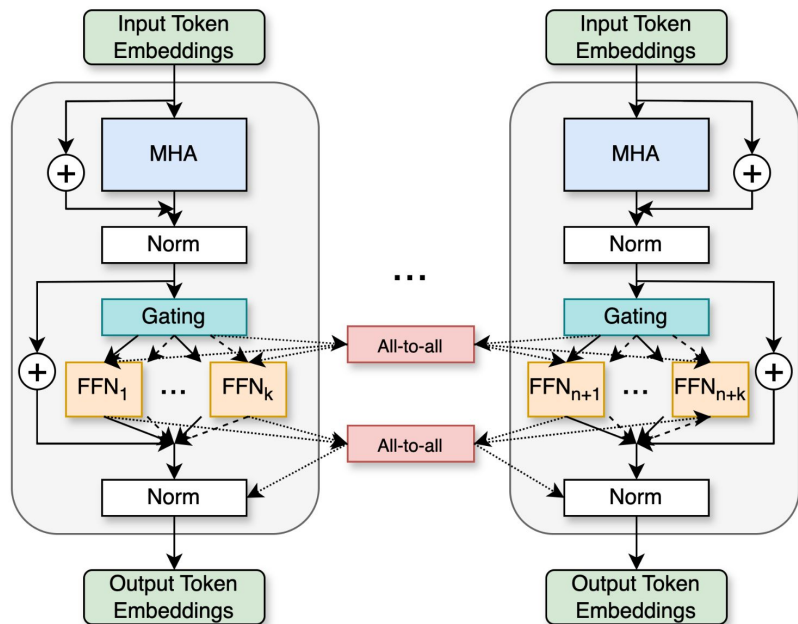Meta    Duke    DEPARTMENT *of* COMPUTER SCIENCE    Penn UNIVERSITY *of* PENNSYLVANIA    WashU

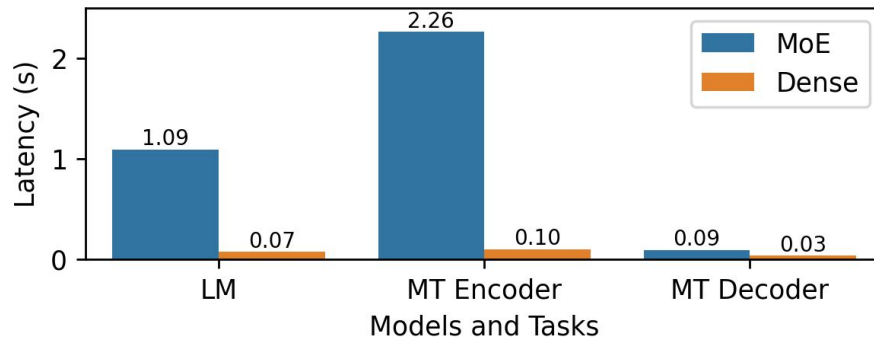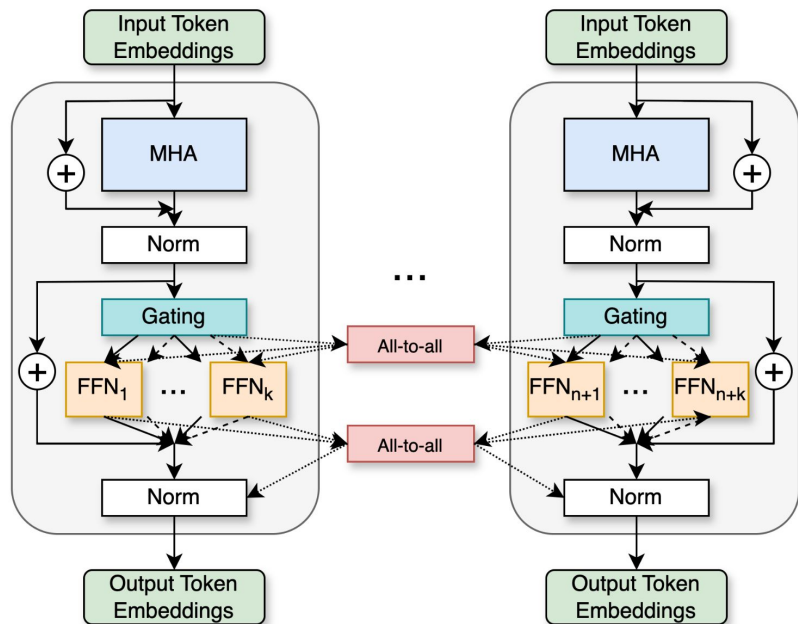# MoE Performance in Training Scenario

**Mixture of Experts** (MoE) models with expert parallelism



Lower **training** cost by up to 5x compared to dense Transformer

# MoE Performance in Inference Scenario

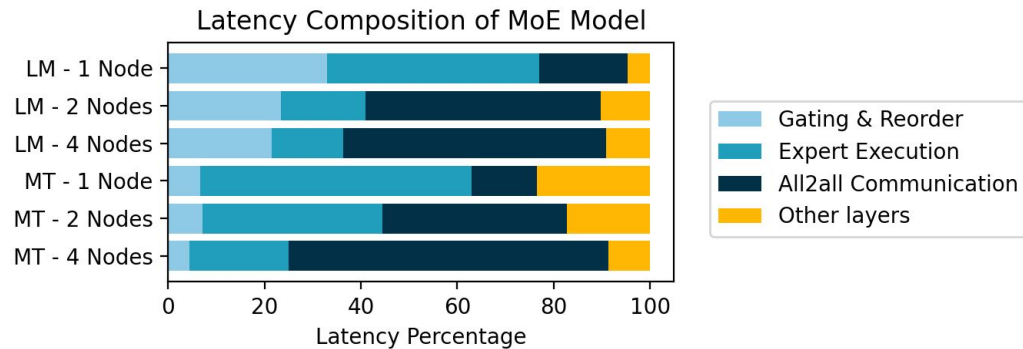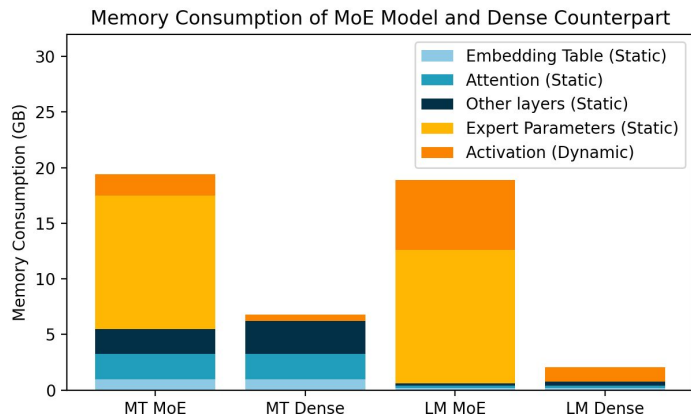**Mixture of Experts** (MoE) models with expert parallelism



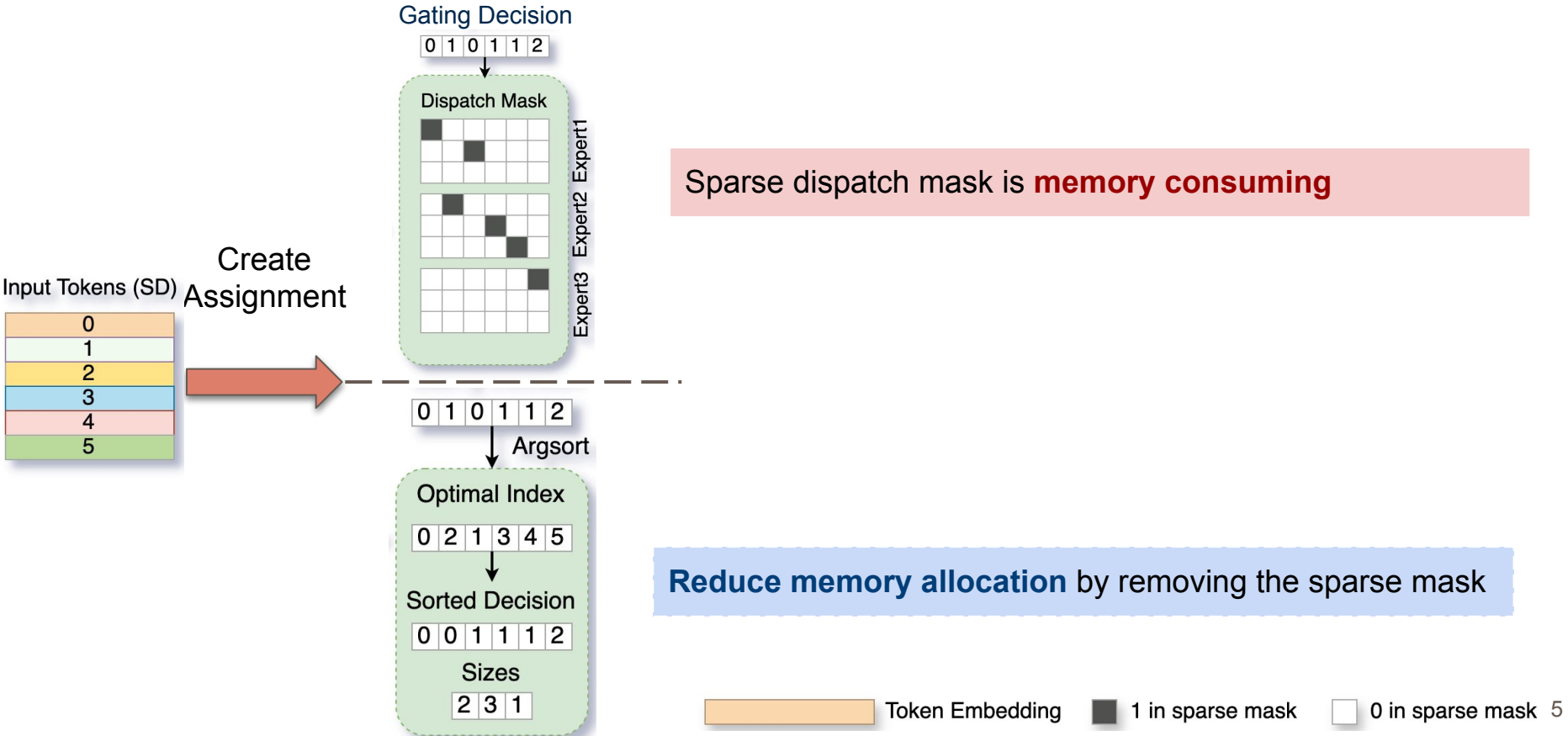Lower **training** cost by up to 5x compared to dense Transformer.

Higher **inference** cost, 15x slower for language models and >3x slower for machine translation.

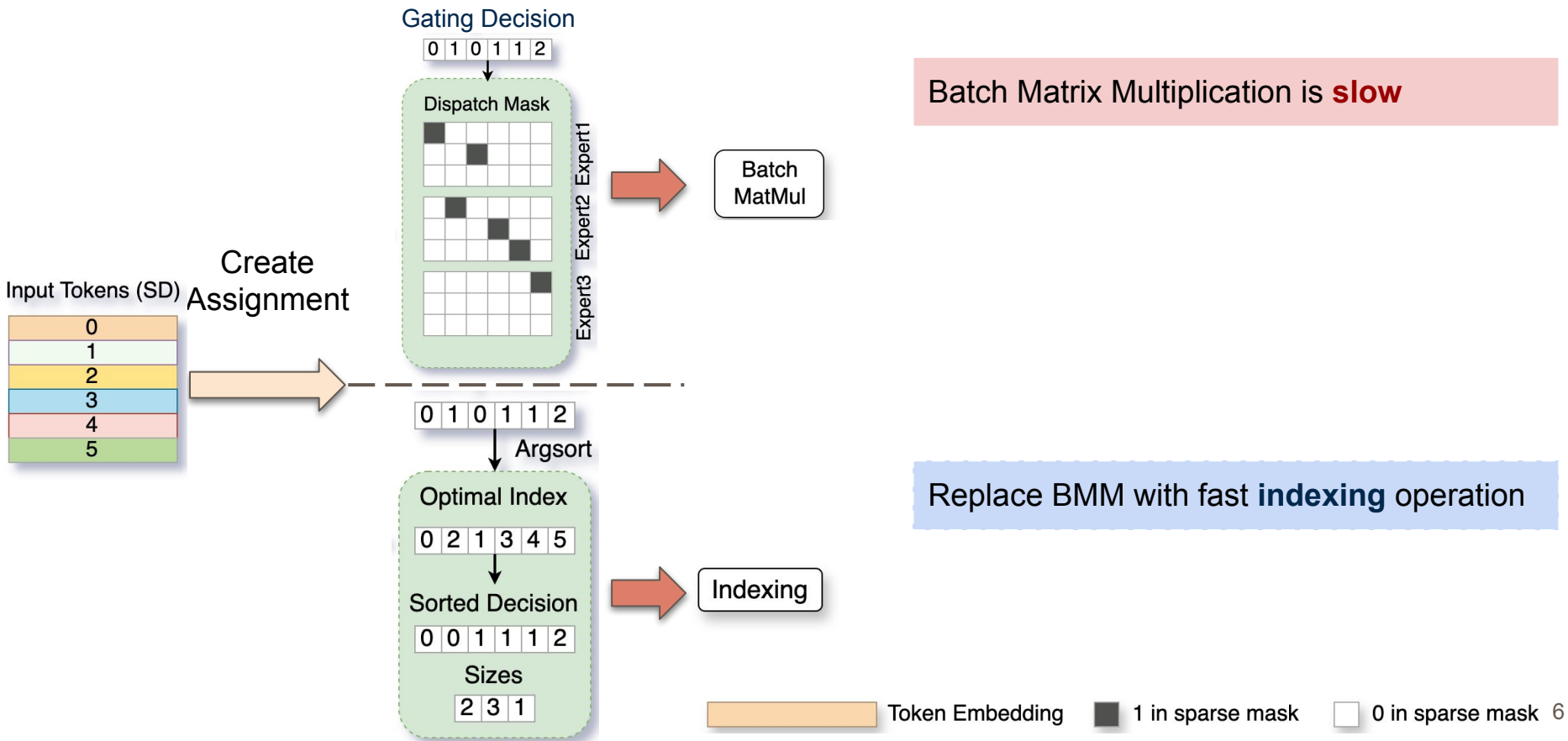# MoE Inference Latency and Memory Characterizations

- **slower** than equivalent* dense counterparts
- consumes **more memory** than equivalent* dense counterparts



Memory Consumption of MoE Model and Dense Counterpart



Latency Composition of MoE Model

# Dynamic Gating: Less Memory

Gating Decision

| 0 | 1 | 0 | 1 | 1 | 2 |

**Dispatch Mask**

Expert1 Expert2 Expert3

**Input Tokens (SD)**

| 0 |
| 1 |
| 2 |
| 3 |
| 4 |
| 5 |

Create Assignment

Sparse dispatch mask is **memory consuming**

| 0 | 1 | 0 | 1 | 1 | 2 |

Argsort

**Optimal Index**

| 0 | 2 | 1 | 3 | 4 | 5 |

**Sorted Decision**

| 0 | 0 | 1 | 1 | 1 | 2 |

**Sizes**

| 2 | 3 | 1 |

**Reduce memory allocation** by removing the sparse mask

Token Embedding    ■ 1 in sparse mask    □ 0 in sparse mask   5

# Dynamic Gating: Less Latency



Gating Decision

Dispatch Mask

Expert1 Expert2 Expert3

Batch MatMul

Create Assignment

Input Tokens (SD)

Batch Matrix Multiplication is **slow**

Argsort

Optimal Index

Sorted Decision

Sizes

Indexing

Replace BMM with fast **indexing** operation

Token Embedding   1 in sparse mask   0 in sparse mask   6
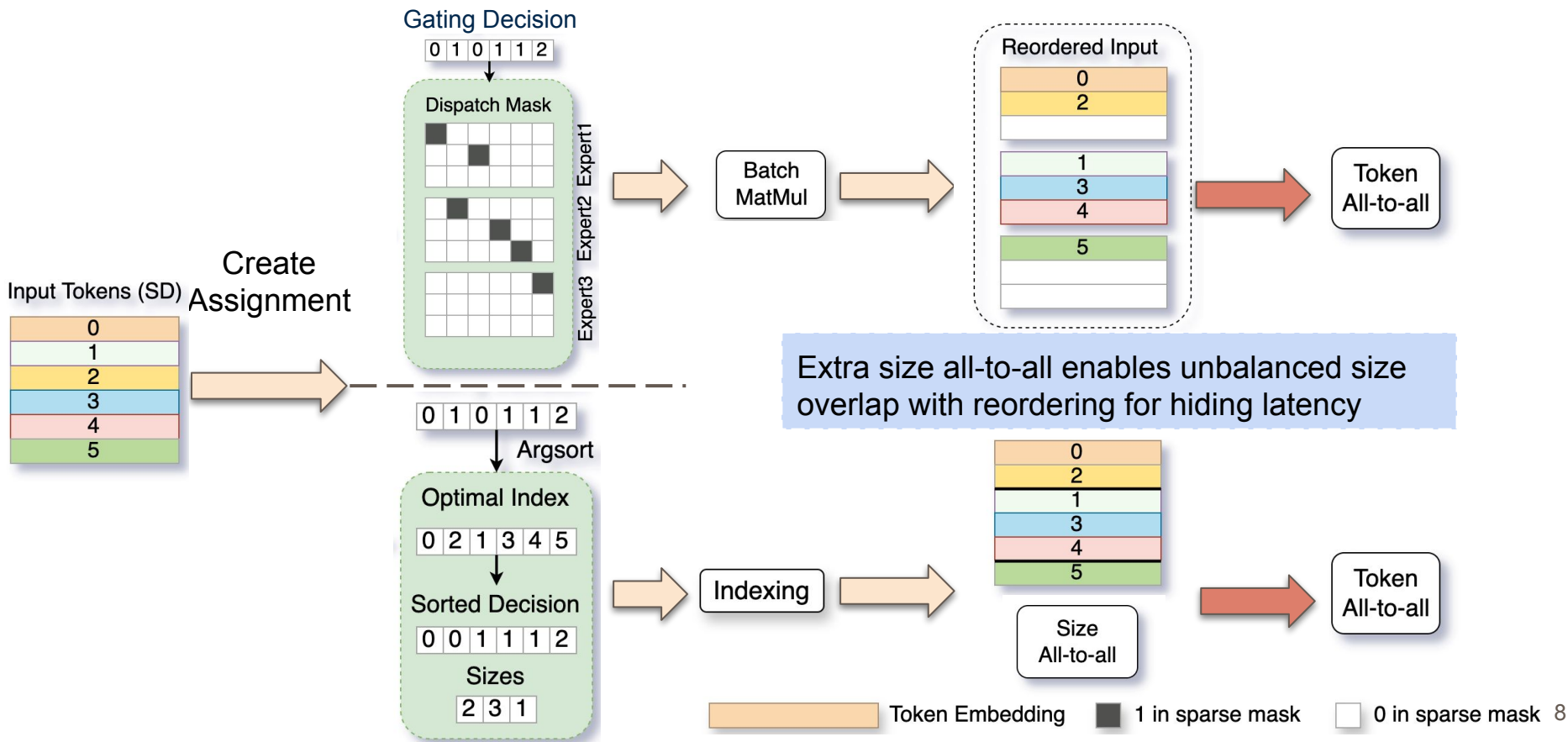
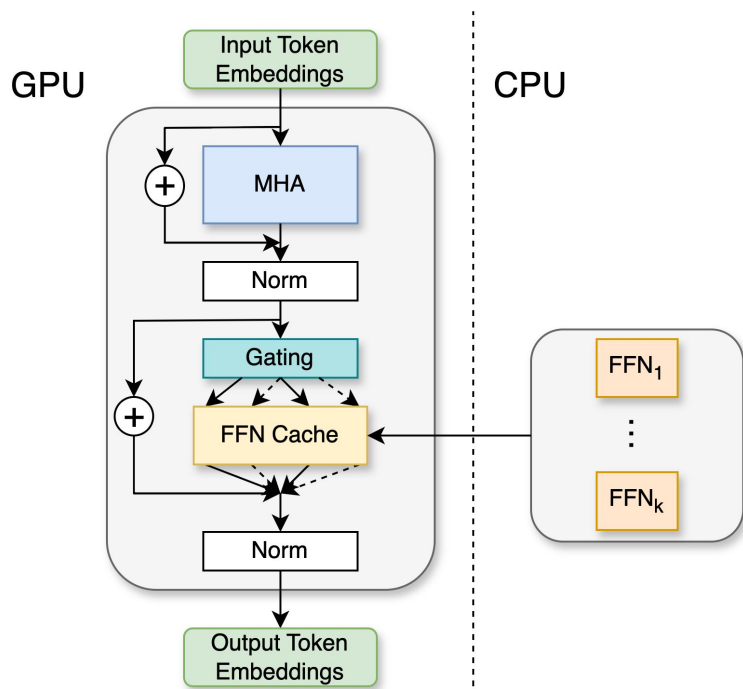# Dynamic Gating: No More Placeholders

# Dynamic Gating: Dynamic All-to-all
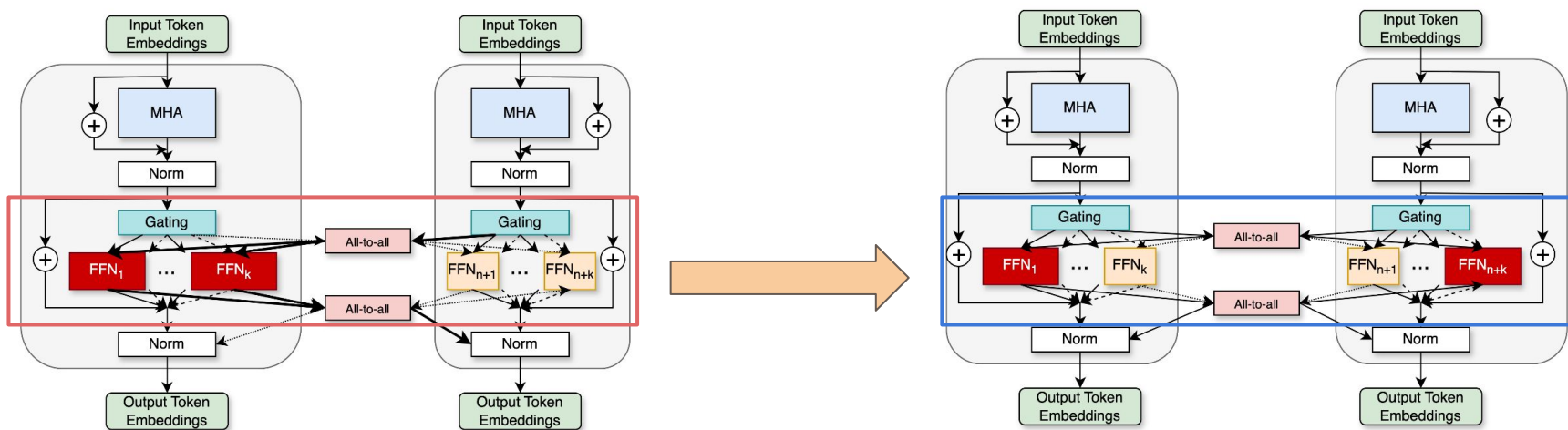
# Expert Buffering: Reduce Static Memory Usage

Only **a subset of experts** is activated in each batch



- Store bulk of the expert parameters in main memory
  → **Reduce Parameter Memory**
- Maintain *LIFO* cache on GPU that stores activated experts
  → **Mitigate GPU-CPU Latency**
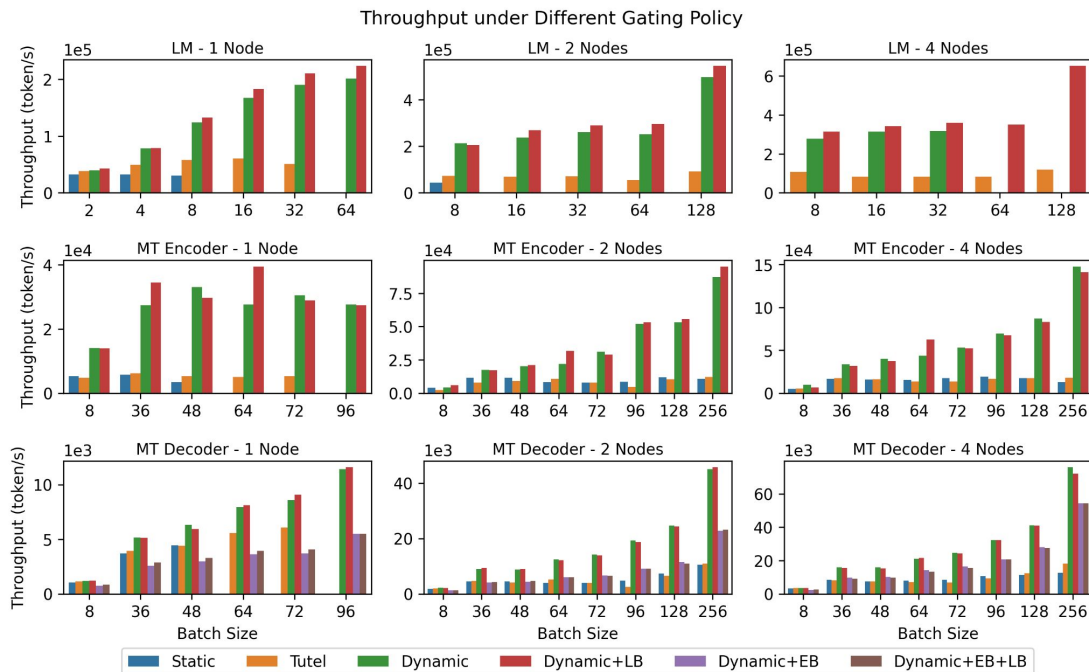
# Load Balancing: Improve Service Robustness

Imbalanced load on each GPU create **memory spikes** and **bottlenecks**



- Estimate expert load from historical activation
- Assign expert based on balancing load

# Results

Improved **throughput** and maximum **batch size**



Throughput under Different Gating Policy

# Check our paper and code for more information!

Paper

Paper (arXiv version)

Code