

# SiDA-MoE: Sparsity-inspired Data-Aware Serving for Efficient and Scalable Large Mixture-of-Experts Models

Zhixu Du<sup>1</sup> Shiyu Li<sup>1</sup> Yuhao Wu<sup>1</sup> Xiangyu Jiang<sup>2</sup> Jingwei Sun<sup>1</sup> Qilin Zheng<sup>1</sup>  
Yongkai Wu<sup>2</sup>, Ang Li<sup>3</sup>, Hai “Helen” Li<sup>1</sup> Yiran Chen<sup>1</sup>

<sup>1</sup>Duke University

<sup>2</sup>Clemson University

<sup>3</sup>University of Maryland, College Park

Duke

Duke  
UNIVERSITY



# Era of Large Models

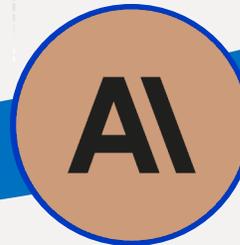
**Nov. 2022**  
ChatGPT  
GPT-3.5

**Mar. 2023**  
Gemini  
GPT-4  
LLaMA

**Jan. 2024**  
Mistral  
Mixtral

**Mar. 2024**  
Claude 3  
Opus  
Sonnet  
Haiku

**Apr. 2024**  
LLaMA3-8B  
LLaMA3-70B  
LLaMA3-400B



# Era of Large MoE Models

Nov. 2022  
ChatGPT  
GPT-3.5

Mar. 2023  
**Gemini**  
**GPT-4**  
LLaMA

Jan. 2024  
Mistral  
**Mixtral**

Mar. 2024  
Claude 3  
Haiku  
Sonnet  
Opus

Apr. 2024  
LLaMA3-8B  
LLaMA3-70B  
LLaMA3-400B

Many large models employ the Mixture-of-Experts (MoE) techniques.



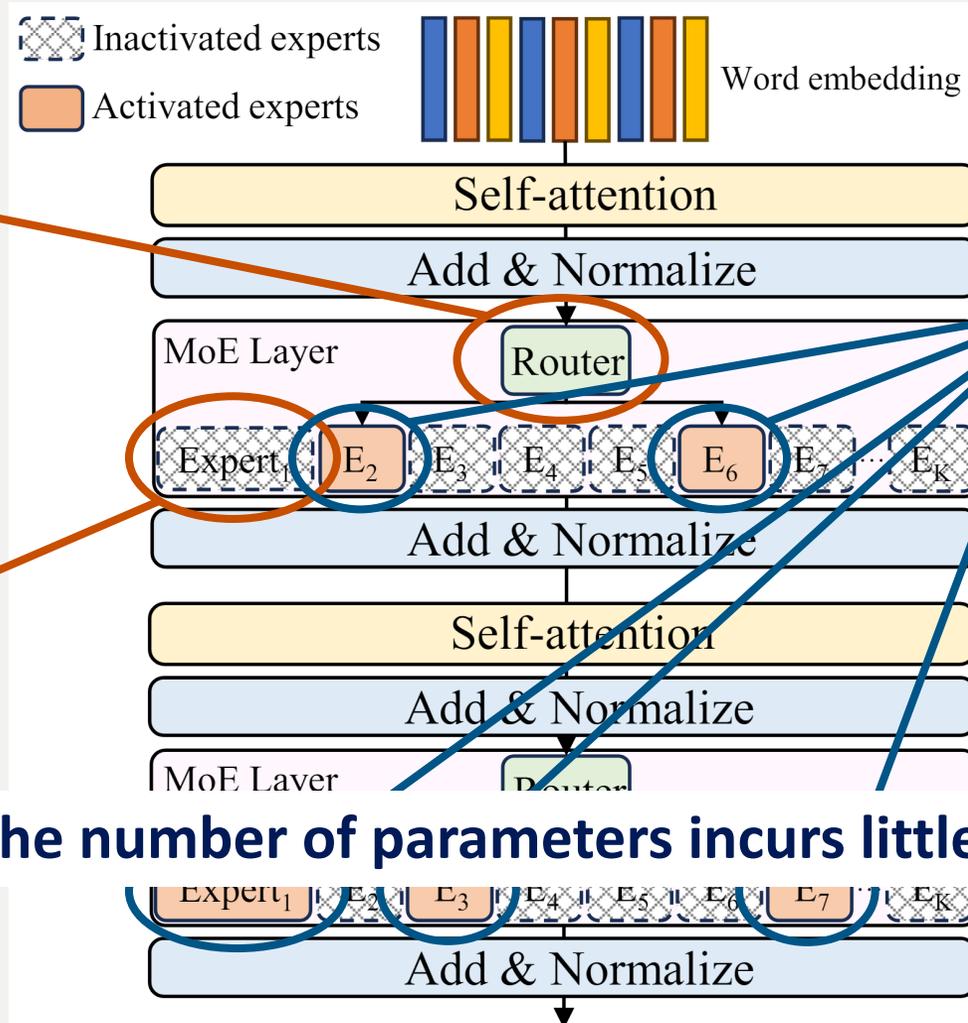
# Era of Large MoE Models

## Router:

Usually, a linear classifier to decide which expert will be selected.

## Expert:

Usually, a MLP model.



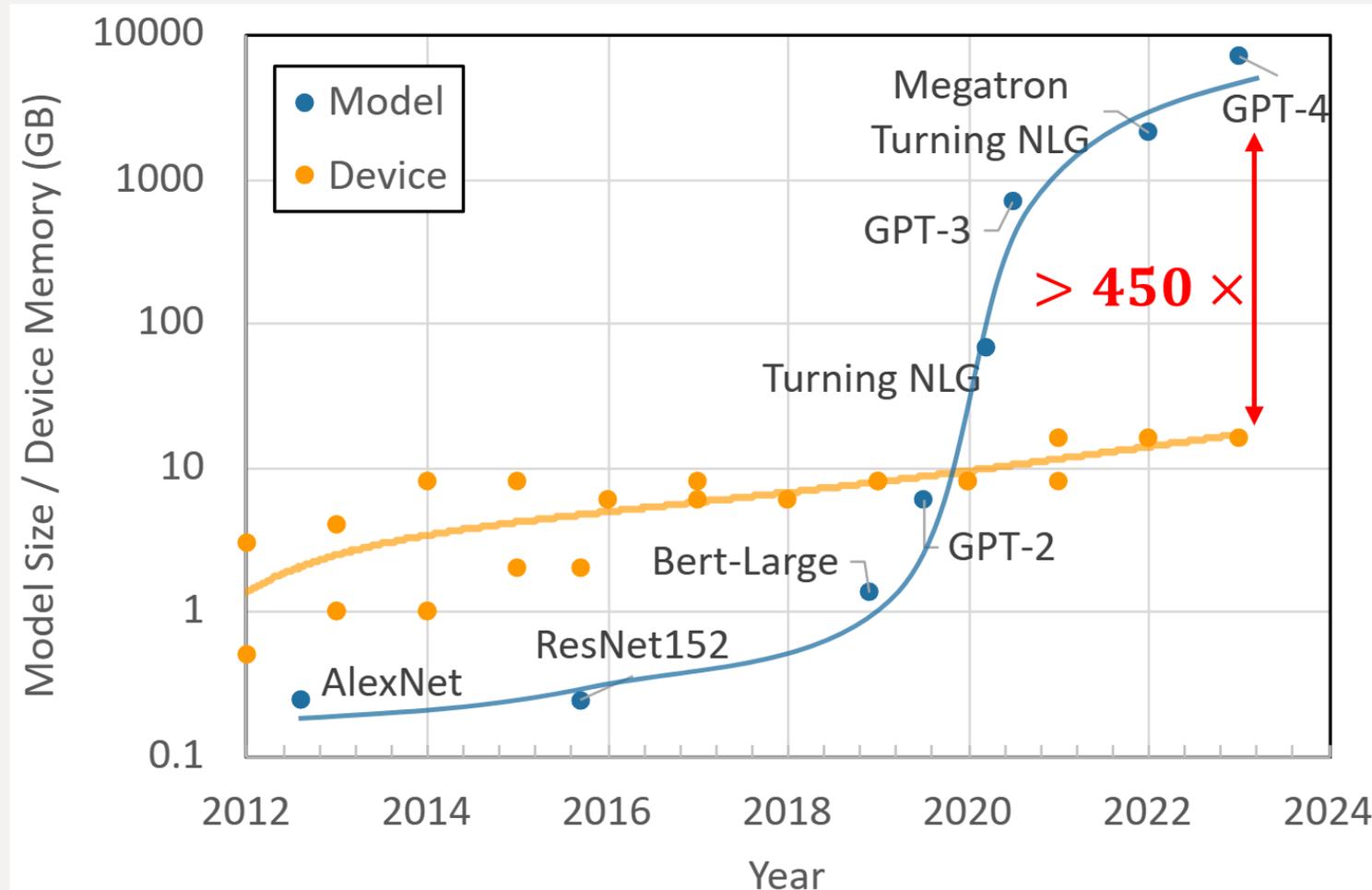
## Activated experts:

Only selected experts will participate in the current round of inference, combined by a weighted sum.

$$\sum_{i \in \mathbb{I}} \alpha_i(\mathbf{x}) f_i(\mathbf{x}; \theta_i)$$

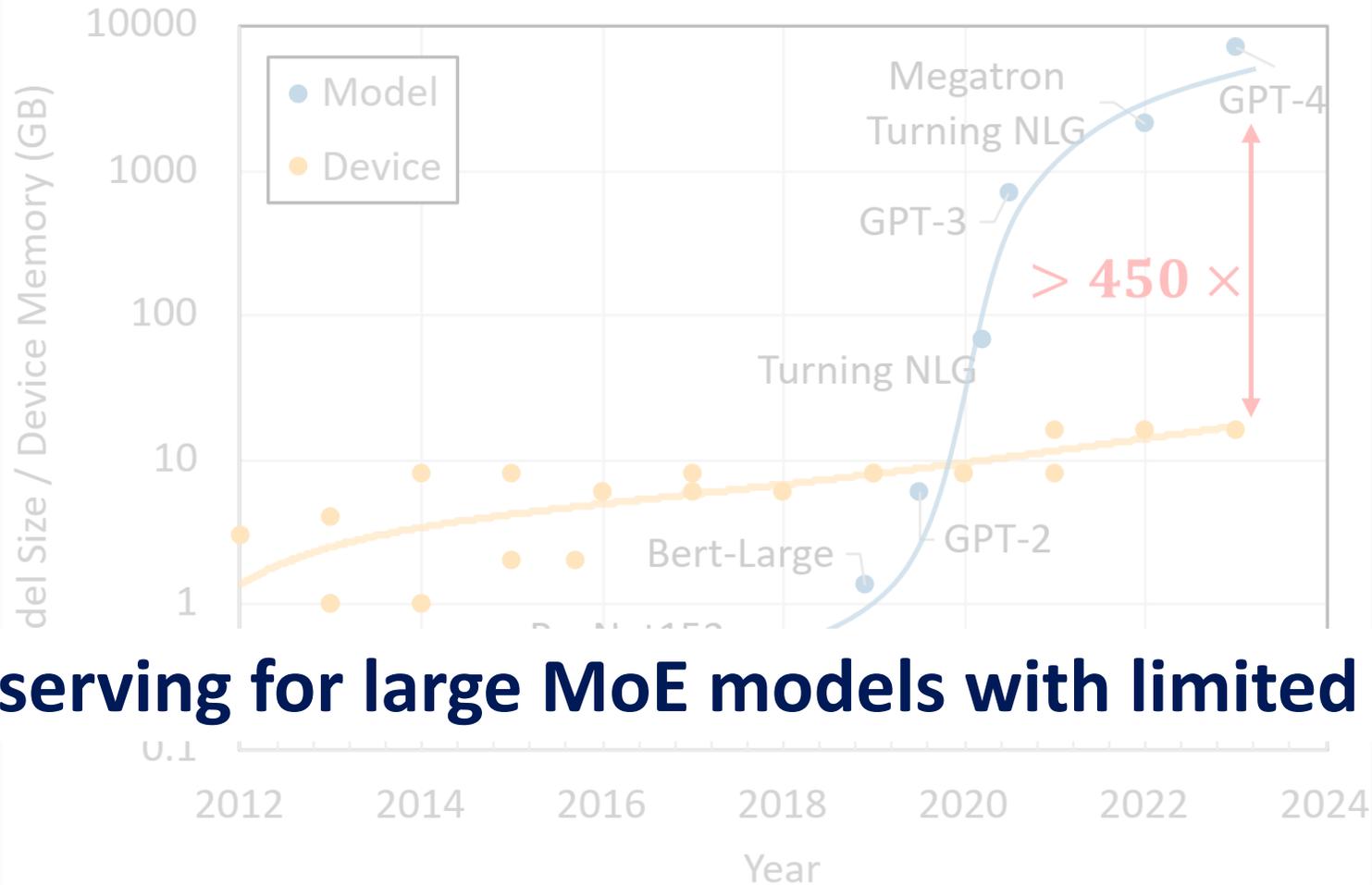
**Drastically increase the number of parameters incurs little computational overhead.**

# Big Models and Small Devices



Devices: Geforce GTX Mobile Series, Apple A Series, and Qualcomm Snapdragon.

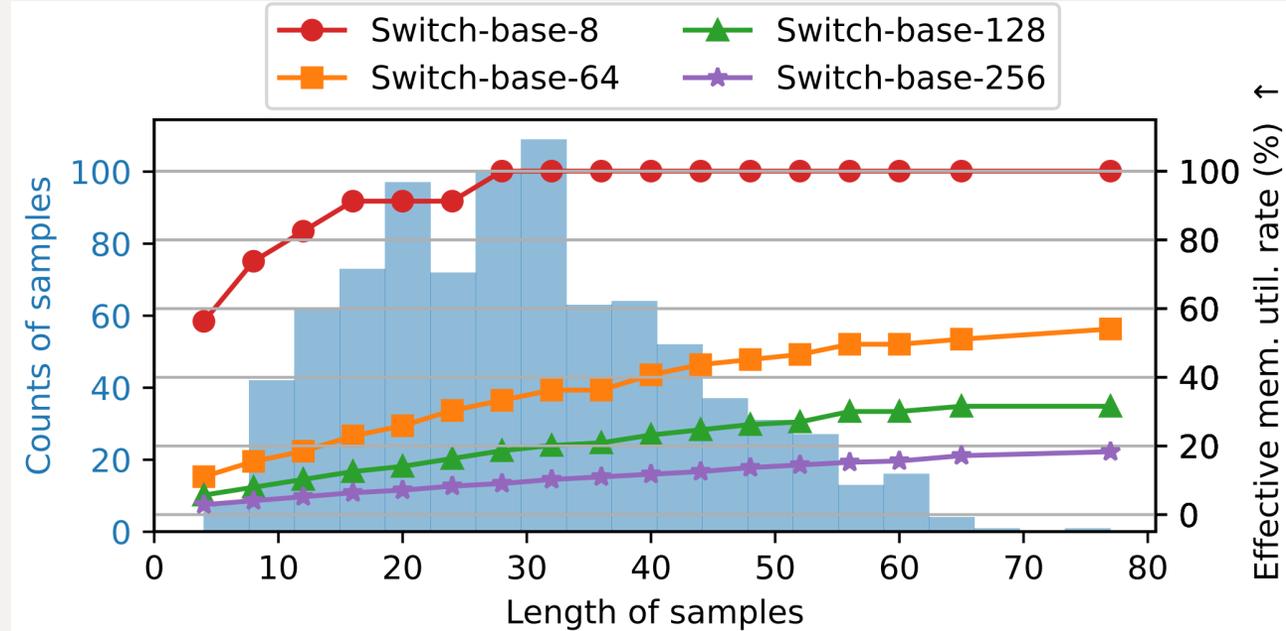
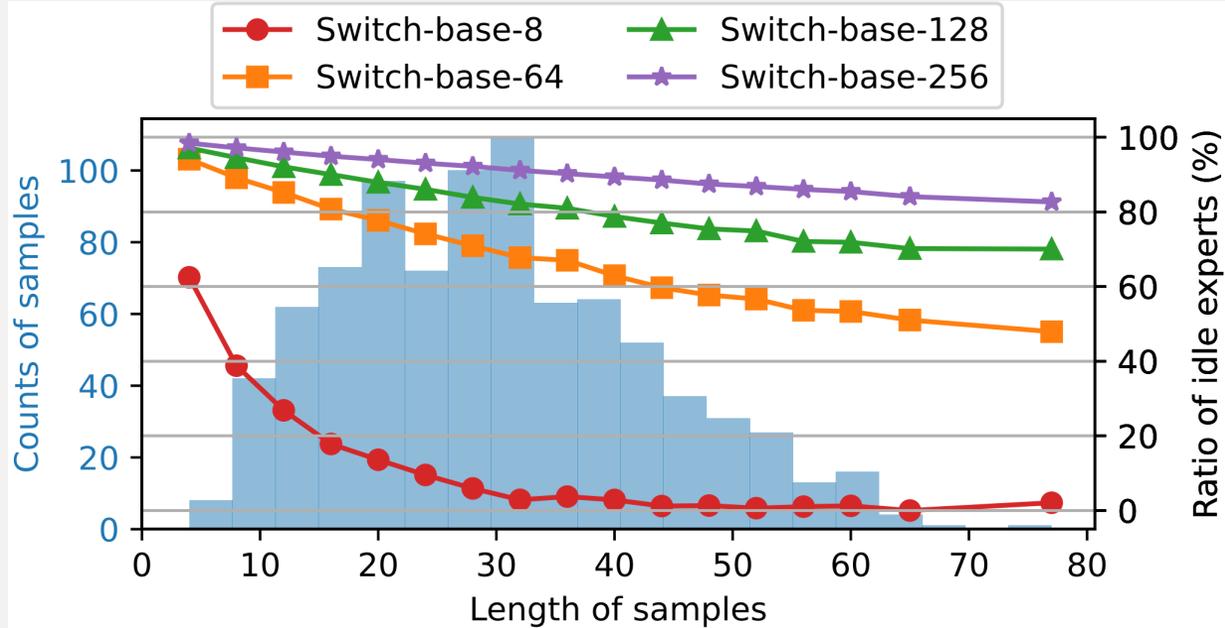
# Big Models and Small Devices



**Efficient serving for large MoE models with limited resources.**

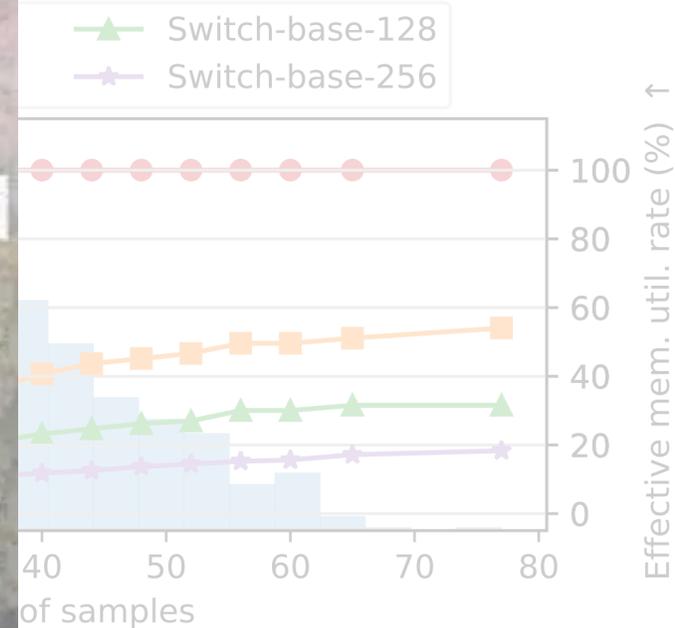
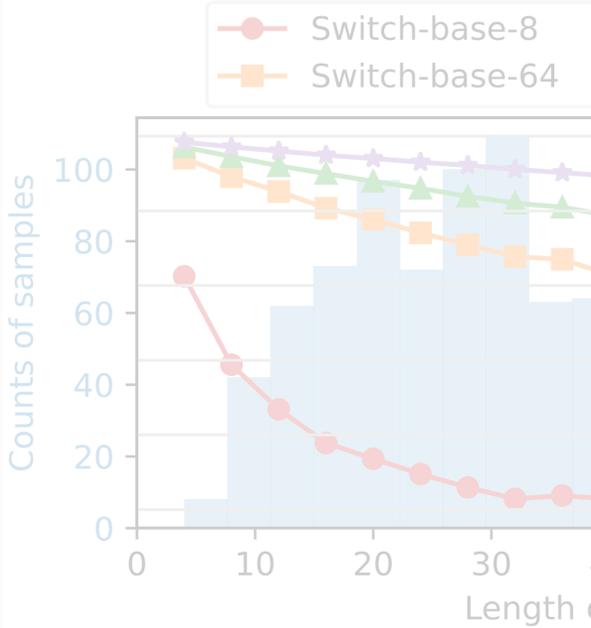
# Motivations

# Sparsely-activated Experts



- Teaser on SST2.
- Up to 80% experts are idle on Switch-base-256.
- Up to 80% GPU memory are ineffective for the inference on Switch-base-256.

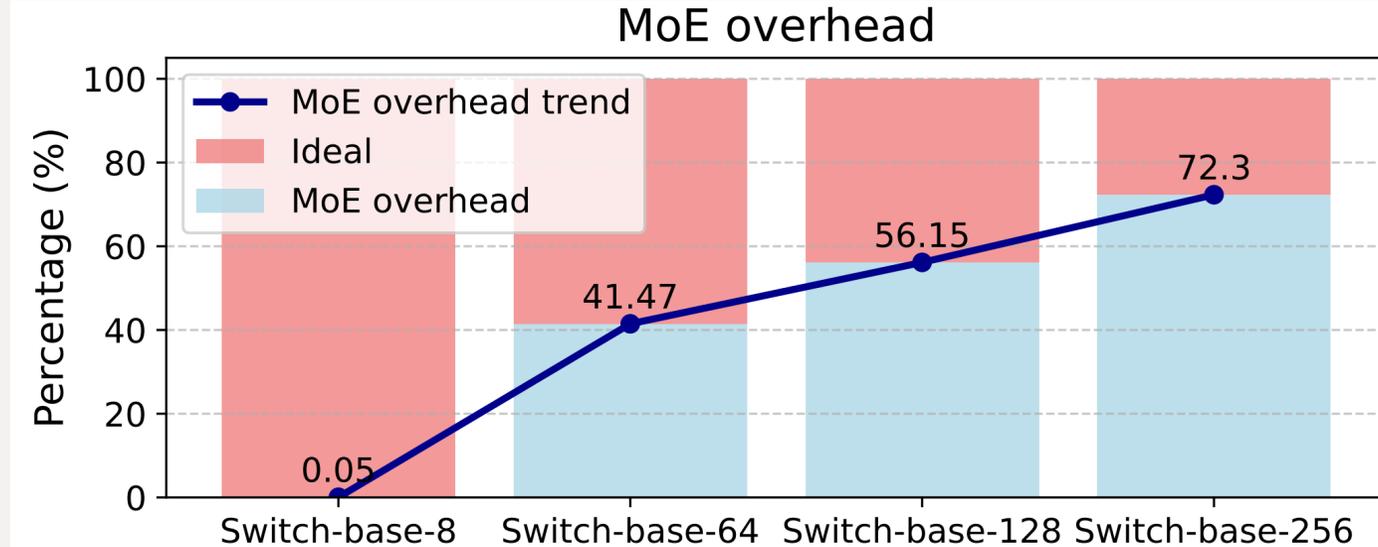
# Sparsely-activated Experts



- Teaser on SS
- Up to 80% ex
- Up to 80% G

wicth-base-256.

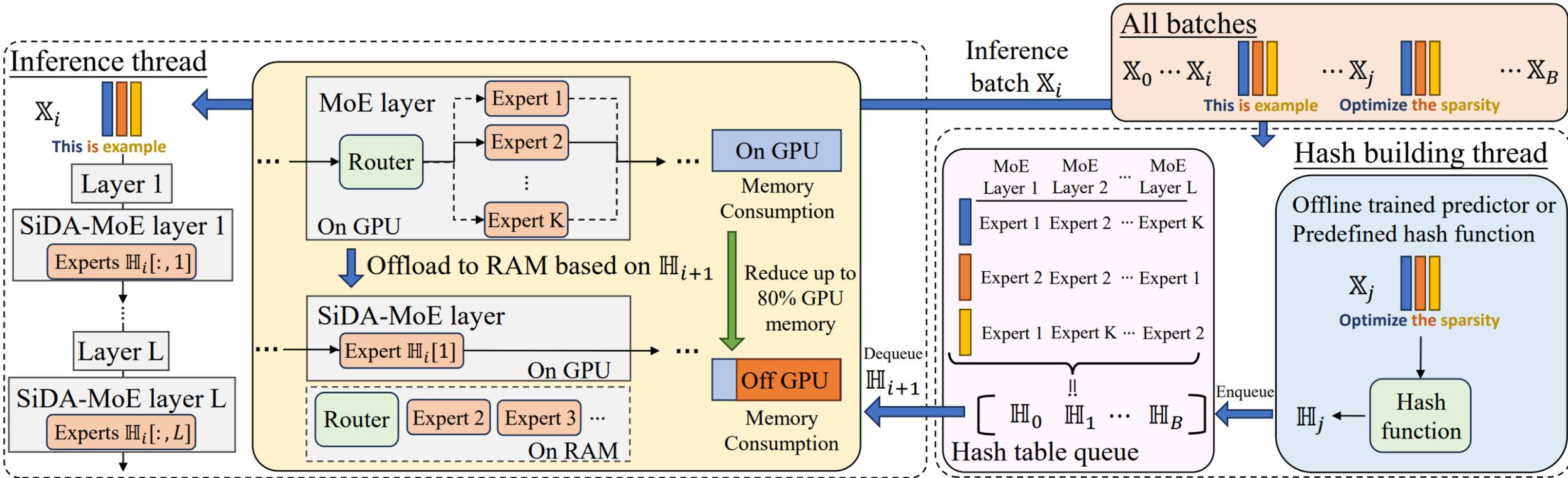
# MoE Overhead on Resource Limited Devices



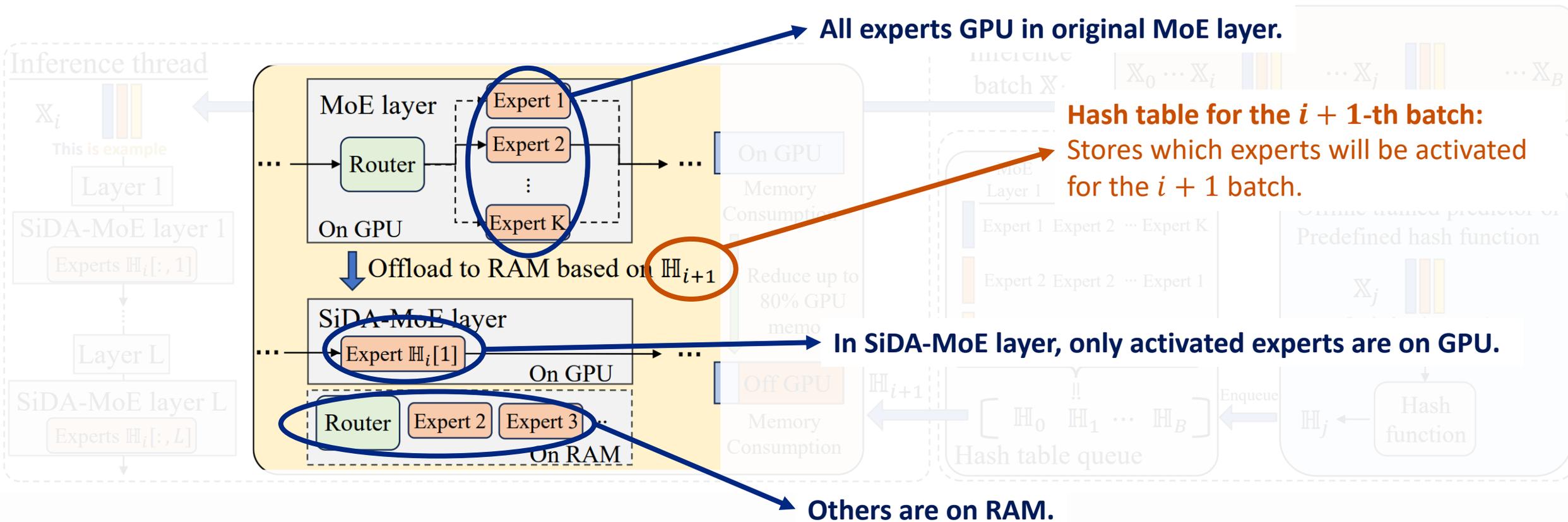
- Teaser on SST2
- MoE Overhead: expert selection, expert invocation, additional communication costs.
- In resource constrained scenarios, the invocation overhead surpasses the computation, meaning the number of experts called dominates the overall inference time.

SiDA-MoE

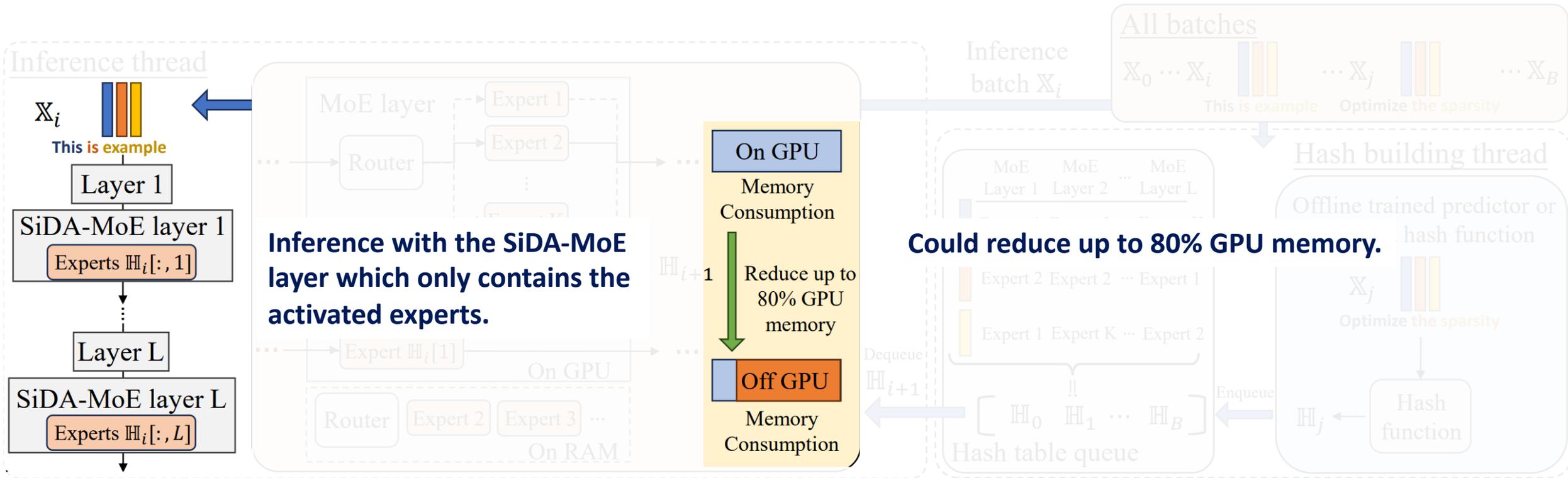
# SiDA: Sparsity-inspired Data-Aware



# SiDA: Sparsity-inspired Data-Aware



# SiDA: Sparsity-inspired Data-Aware



# SiDA: Sparsity-inspired Data-Aware

## The hash function:

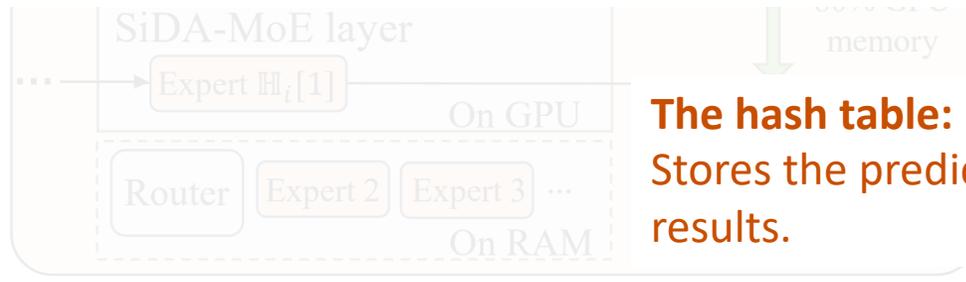
Takes batch of data as input.

Output the index of expert to be activated for each token in each layer.

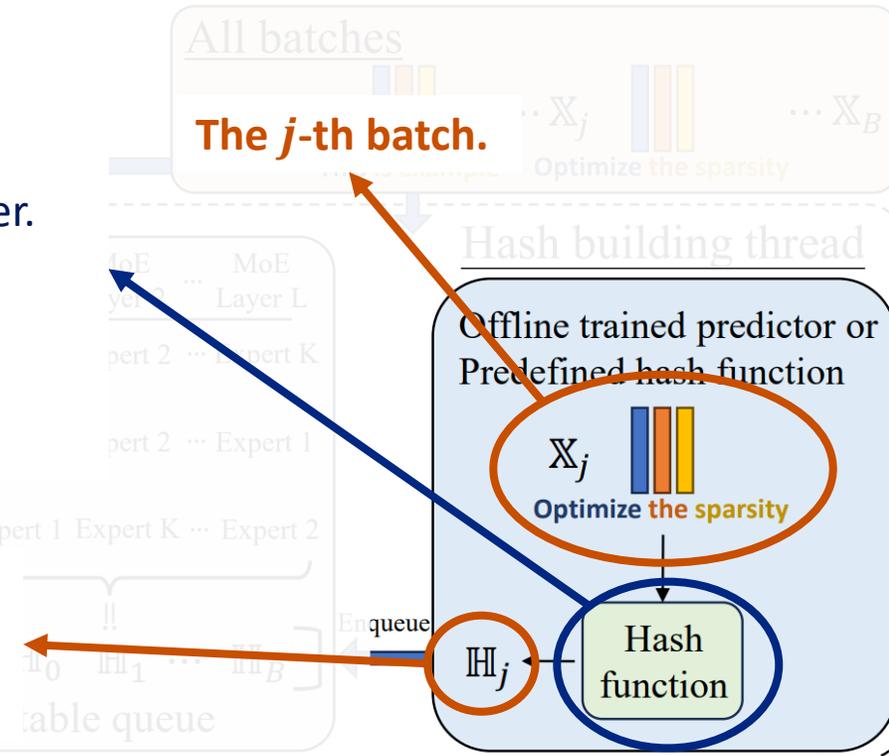
Can be predefined (such as the Hash Layer MoE) or offline trained.

We offline trained a LSTM with  $L$  classification heads.

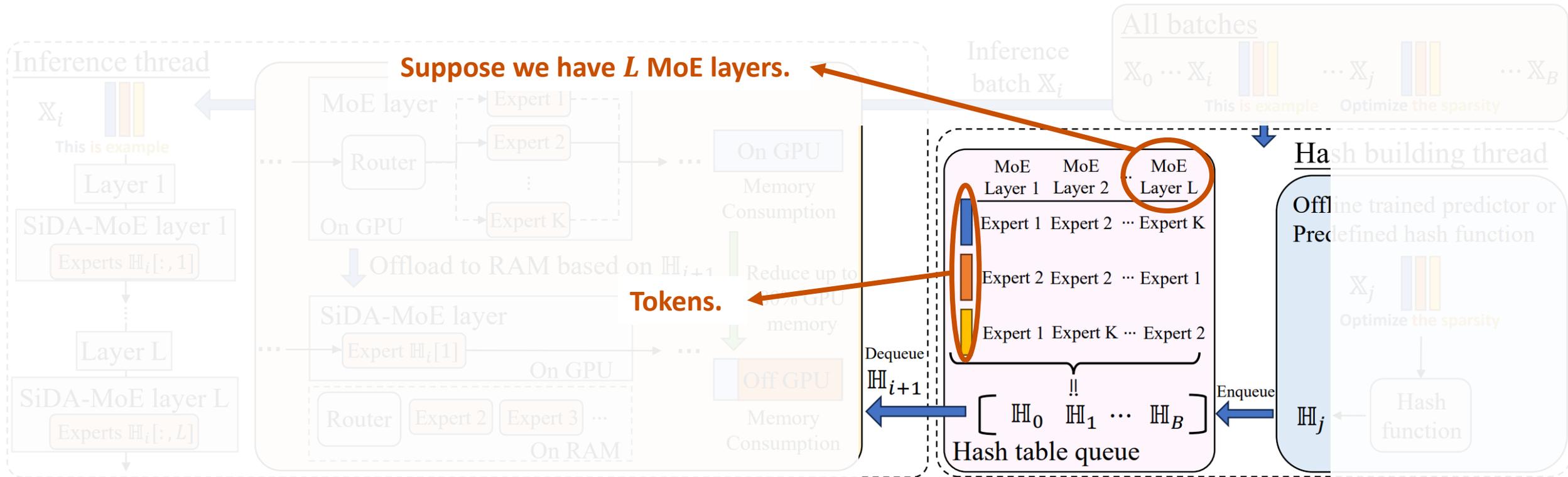
Trained on (sample, expert activation pattern) pairs.



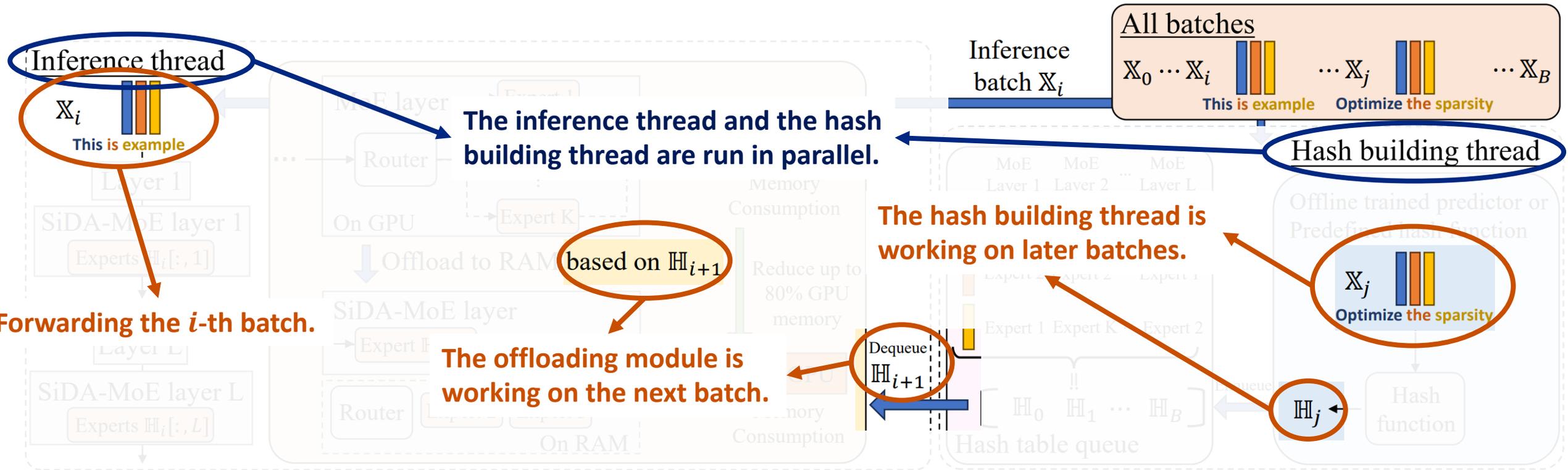
**The hash table:**  
Stores the prediction results.



# SiDA: Sparsity-inspired Data-Aware

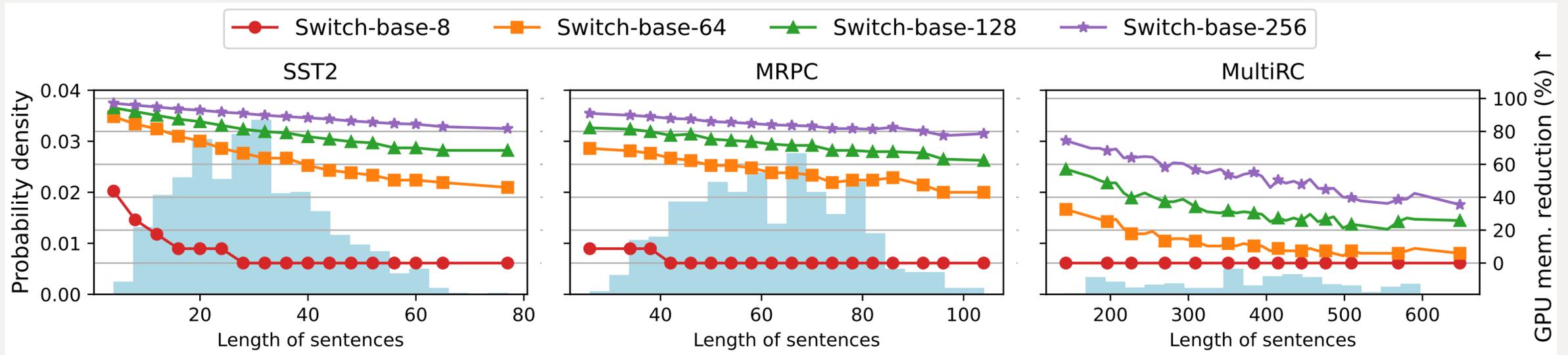


# Our System



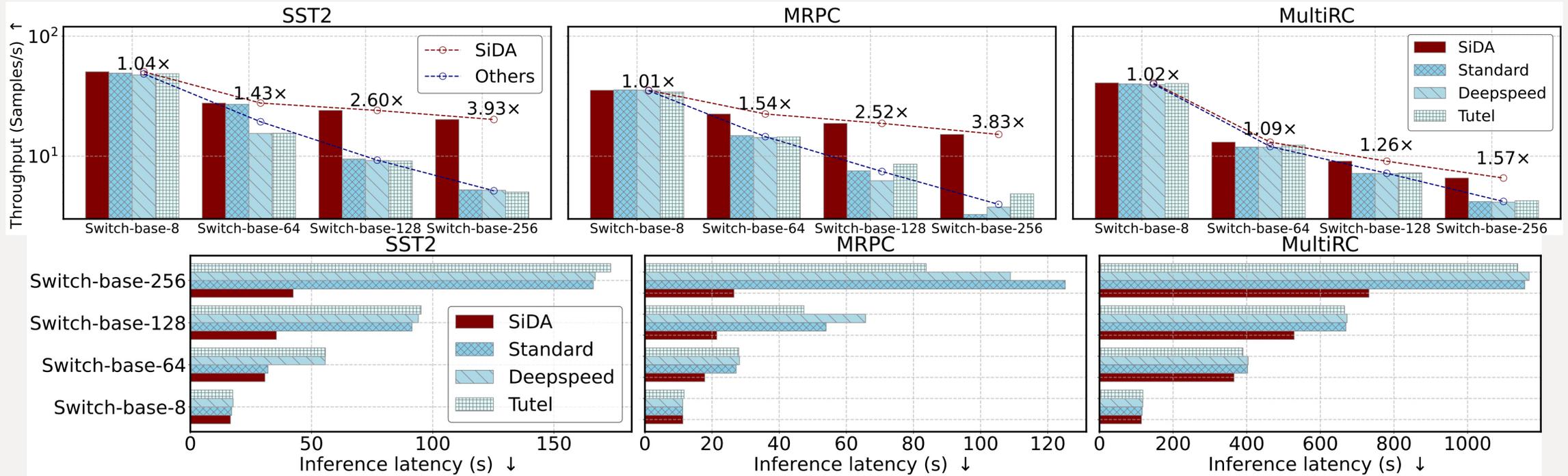
# Experiments

# Memory



- SiDA saves up to 80% memory for Switch-base-256 on short sentences.
- SiDA saves over 40% memory for Switch-base-256 on Long sentences.

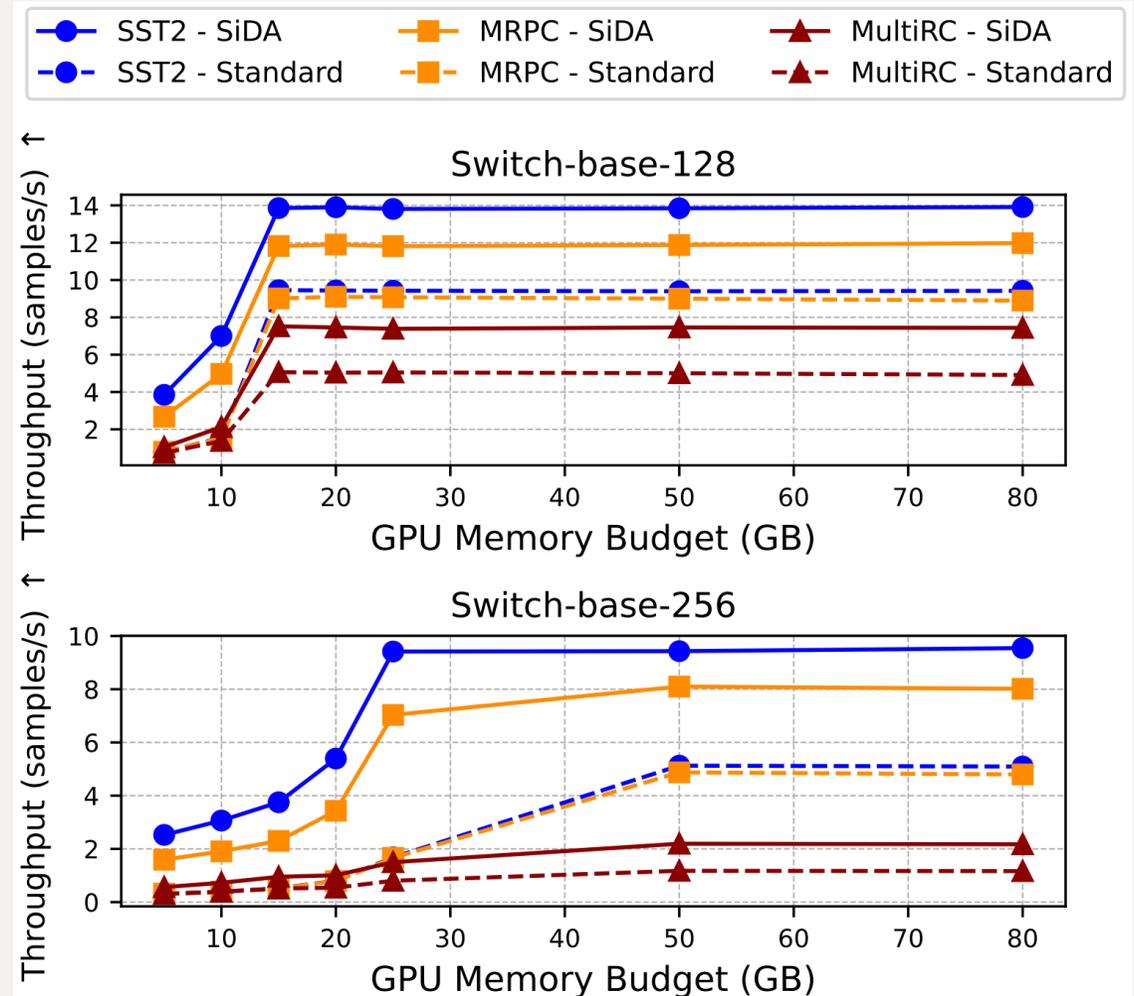
# Throughput and Latency



- SiDA achieves up to 3.93x speed up on throughput for short sentences.
- SiDA achieves at least 1.57x speed up on throughput for short sentences.

# Memory-limited Scenarios

SiDA consistently outperforms baselines given different GPU memory budgets.



# Accuracy Preservation

Backbone	SST2	MRPC	MultiRC
Switch-base-8	99.00%	97.41%	91.74%
Switch-base-128	98.78%	98.65%	90.49%

**Top-3 hash hits rate (prediction accuracy of the hash function).**

Backbone	Pretrained ppl. ( $\downarrow$ )	SiDA-MoE ppl. ( $\downarrow$ )
Switch-base-8	6.68	18.49
Switch-base-64	4.93	11.84
Switch-base-128	4.86	11.73
Switch-base-256	4.59	8.11

**Accuracy compromise as a pretrained model.**

Backbone		SST2	MRPC	MultiRC
Switch-base-8	Finetuned	92.20	89.14	56.70
	SiDA-MoE	90.59	86.91	56.11
	Fidelity	98.25%	97.49%	98.95%
Switch-base-128	Finetuned	93.57	89.66	59.95
	SiDA-MoE	87.04	83.01	55.49
	Fidelity	93.02%	92.59%	92.56%

**Accuracy compromise on downstream tasks.**

# Summary

- The Mixture-of-Experts technique is getting popular in the era of large models.
- We propose SiDA-MoE, a sparsity-inspired data-aware serving for MoE models.
- SiDA-MoE achieves huge amount of GPU memory save and significant speed up, while compromising little accuracy.

# Thanks for your attentions!

If you have further questions, please feel free to contact me:  
[zhixu.du@duke.edu](mailto:zhixu.du@duke.edu)