



THE UNIVERSITY  
*of* EDINBURGH

# MoE-CAP: Benchmarking Cost, Accuracy and Performance of Sparse Mixture-of-Experts Systems

Yinsicheng Jiang<sup>\*1</sup>, Yao Fu<sup>\*1</sup>, Yeqi Huang<sup>\*1</sup>, Ping Nie<sup>3</sup>, Zhan Lu<sup>1</sup>, Leyang Xue<sup>1</sup>, Congjie He<sup>1</sup>, Man-kit Sit<sup>1</sup>, Jilong Xue<sup>2</sup>, Li Dong<sup>2</sup>, Ziming Miao<sup>2</sup>, Dayou Du<sup>1</sup>, Tairan Xu<sup>1</sup>, Kai Zou<sup>4</sup>, Edoardo Ponti<sup>1 5</sup>, Luo Mai<sup>1</sup>

University of Edinburgh<sup>1</sup>, Microsoft Research<sup>2</sup>, Peking University<sup>3</sup>, NetMind.AI<sup>4</sup>, NVIDIA<sup>5</sup>

# Challenges for Deploying MoE Systems

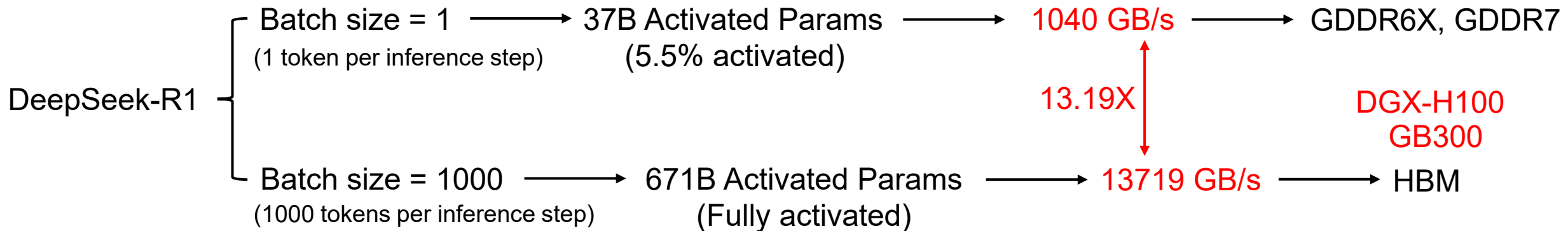
**What is MoE system?** MoE system is software that incorporates specific optimizations for MoE inference:  
- Expert Parallelism, Quantization, Expert offloading (from GPU VRAM to CPU DRAM), ...

## 1. Sparsity largely impacts bandwidth requirements for inference

e.g. DeepSeek-R1 – 671B total params – 37B activated params per token

Target 10 tokens/s

RTX 4090  
RTX 5090



## 2. Too many deployment scenarios

- single-user vs. multi-user
- small vs. large batch sizes

**Seeking benchmarks to answer this question**

**Given a model and a deployment scenario, what is the most suitable hardware and MoE system?**

# Existing Benchmarks Limitations

## **Inadequate insights to guide hardware selection and further optimization**

- Open-LLM-Leaderboard
  - Accuracy, CO<sub>2</sub> Cost
  - No system performance reported
- LLM-Perf, Artificial Analysis, MLPerf
  - Prefill latency, Decoding throughput, Energy usage, Accuracy
  - Lack insights for hardware utilization and MoE system further improvement
- Databricks benchmark, LLM-Viewer
  - MBU (Model Bandwidth Utilization), MFU (Model FLOPS Utilization), Roofline model
  - Inaccurate for MoE systems.

# Our Contributions

## We are the first to introduce CAP benchmarking for MoE systems

### How to classify current MoE systems?

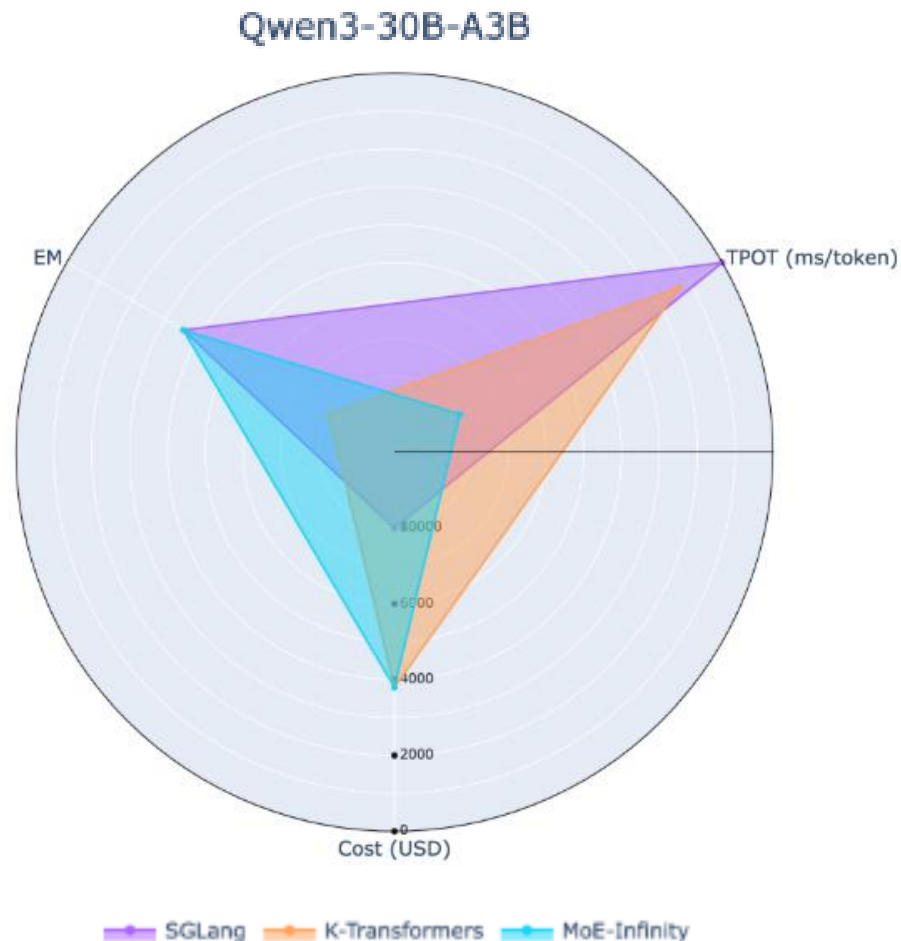
- PA systems (Expert Parallelism)
- PC systems (Quantization)
- CA systems (Offloading)

### How to measure accurate hardware utilization for MoE systems?

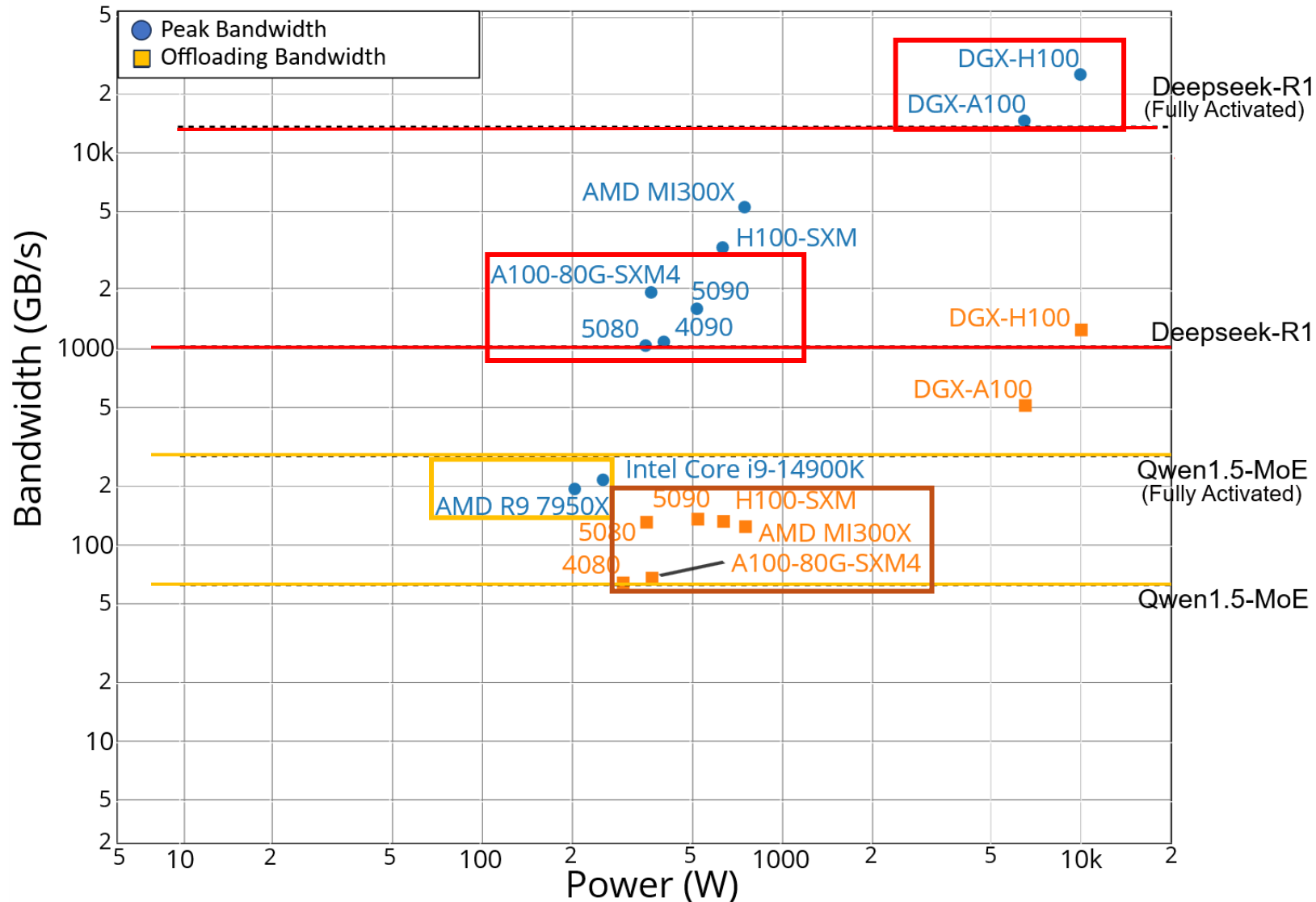
- S-MFU (Sparse Model FLOPS Utilization)
- S-MBU (Sparse Model Bandwidth Utilization)

### How to visualize the benchmark?

- CAP Radar diagram



# Benchmark and Takeaways



**MoE systems enable a broader range of devices to perform inference.**  
- Personal machines can serve large MoE models at small batch size

**Hybrid computing will become more prevalent.**  
- CPU matrix computation capabilities are increasing and DRAM can be easily scaled.

**MoE systems should be co-designed to align with specific applications and deployment scenarios.**

**New benchmarking and design principles are needed for emerging sparse AI systems.**  
- MoE is not everything for sparse AI systems. What about sparse KV Cache?

# Epilogue

MoE-CAP can

- Offer accurate hardware utilization metrics
- Guide further improvements in the system efficiency
- Help choose the best hardware based on the use case for cost savings
- Analyze the trade-off of Cost, Accuracy and Performance among different MoE systems

More details of MoE-CAP, please refer to our paper and code:

<https://arxiv.org/abs/2412.07067>

<https://github.com/sparse-generative-ai/MoE-CAP>

Let's work together to enhance MoE efficiency!

