

Cost-Efficient LLM Training with Lifetime-Aware Tensor Offloading via GPUDirect Storage



Ziqi Yuan¹, Haoyang Zhang¹, Yirui Eric Zhou¹,
Apoorve Mohan², I-Hsin Chung², Seetharami Seelam², Jian Huang¹
¹University of Illinois Urbana-Champaign ²IBM Research



I. Large LLM Workloads Are Hungry for Memory

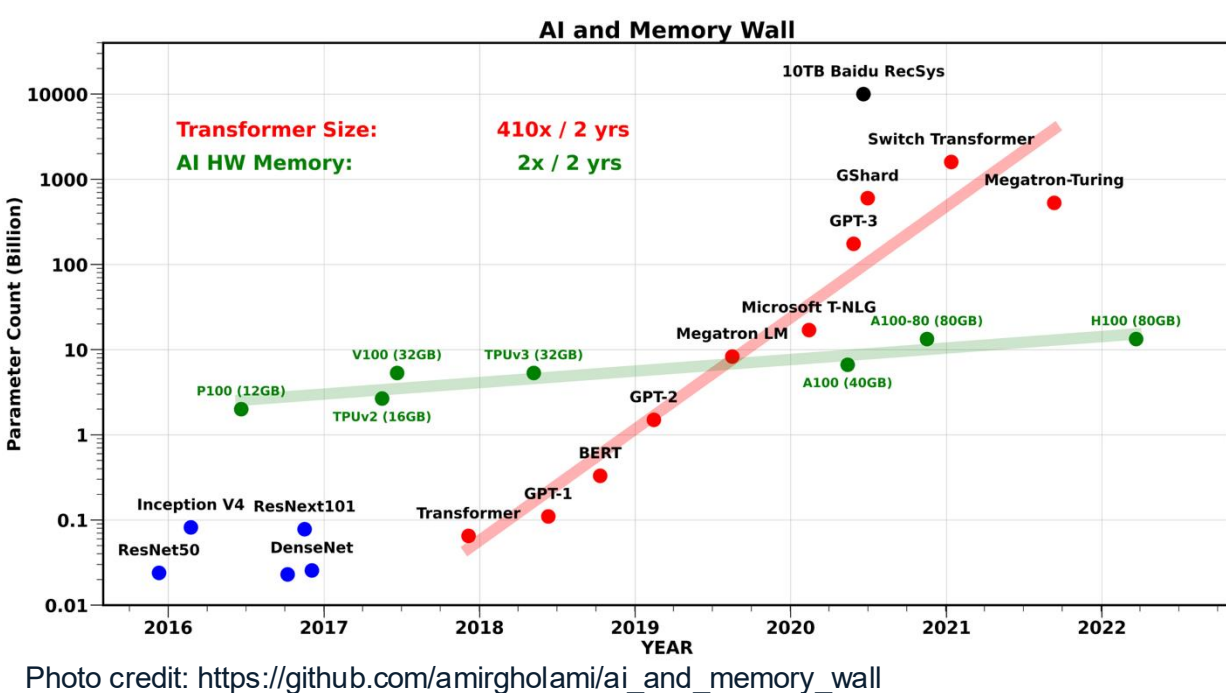
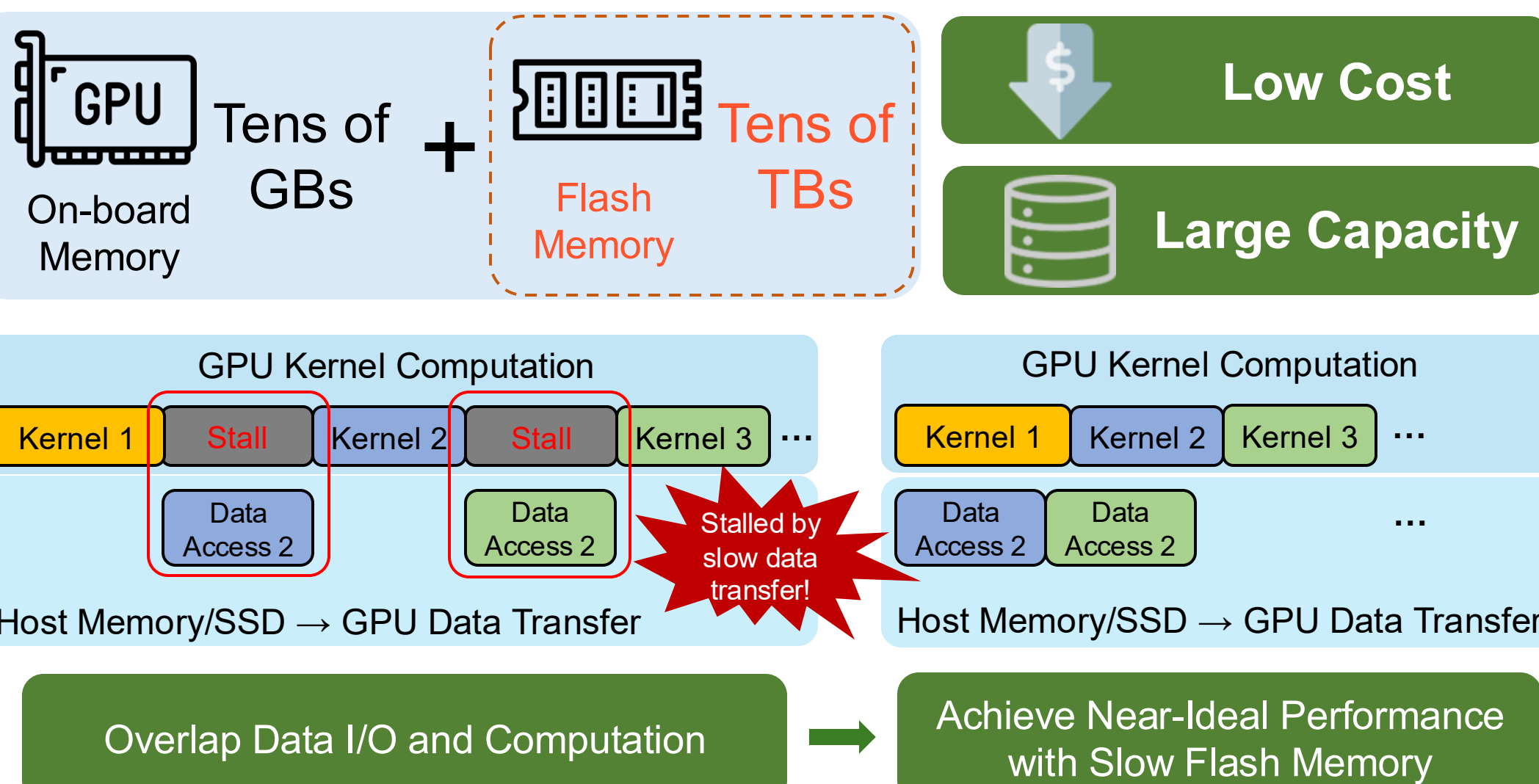


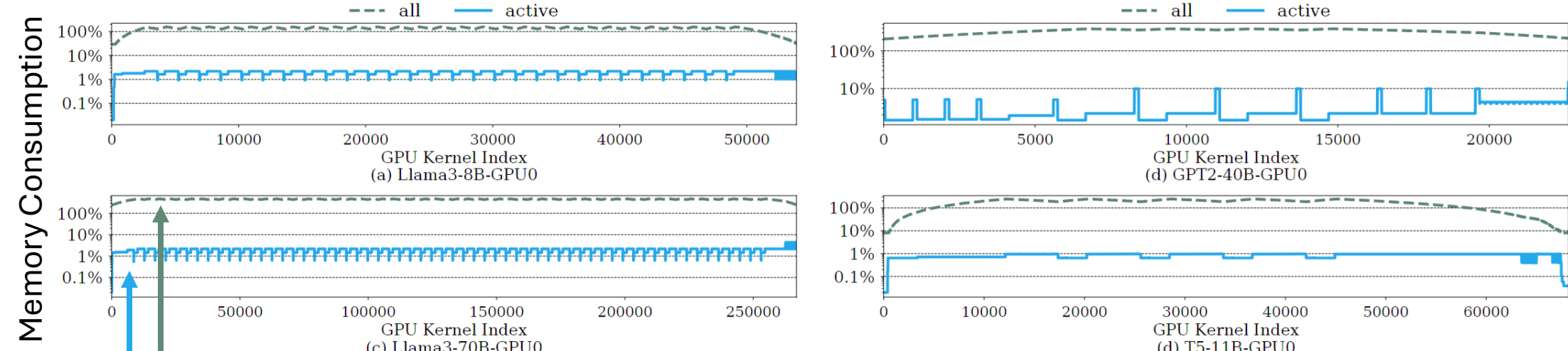
Photo credit: https://github.com/amirgholami/ai_and_memory_wall

II. Expanding GPU Memory with Flash Memory



III. LLM Training Memory Characterization Study

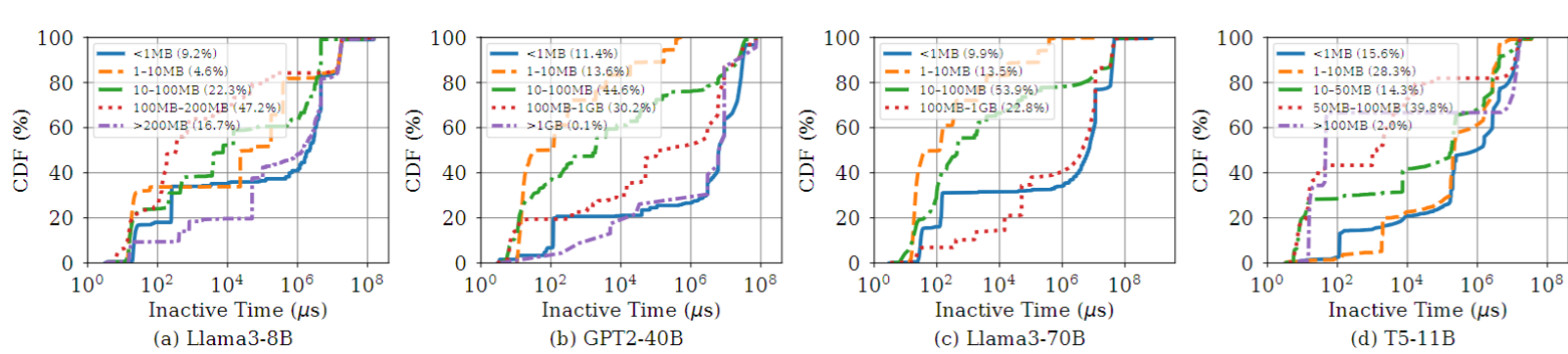
Observation 1: Active Tensors Require Only A Small Portion of GPU Memory



Active + Inactive Tensors
Active Tensors (less than 14%, 1.7% on average)

Most tensors are inactive and can be swapped out during training

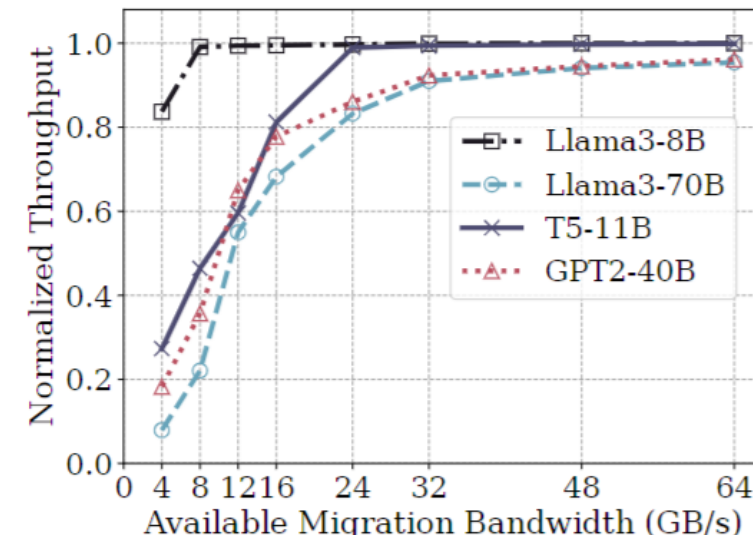
Observation 2: Many Tensors Are Unused for A Long Time



Distribution (CDF) of tensor inactive period lengths for different tensor sizes

Many tensors can be safely swapped out to slow memory

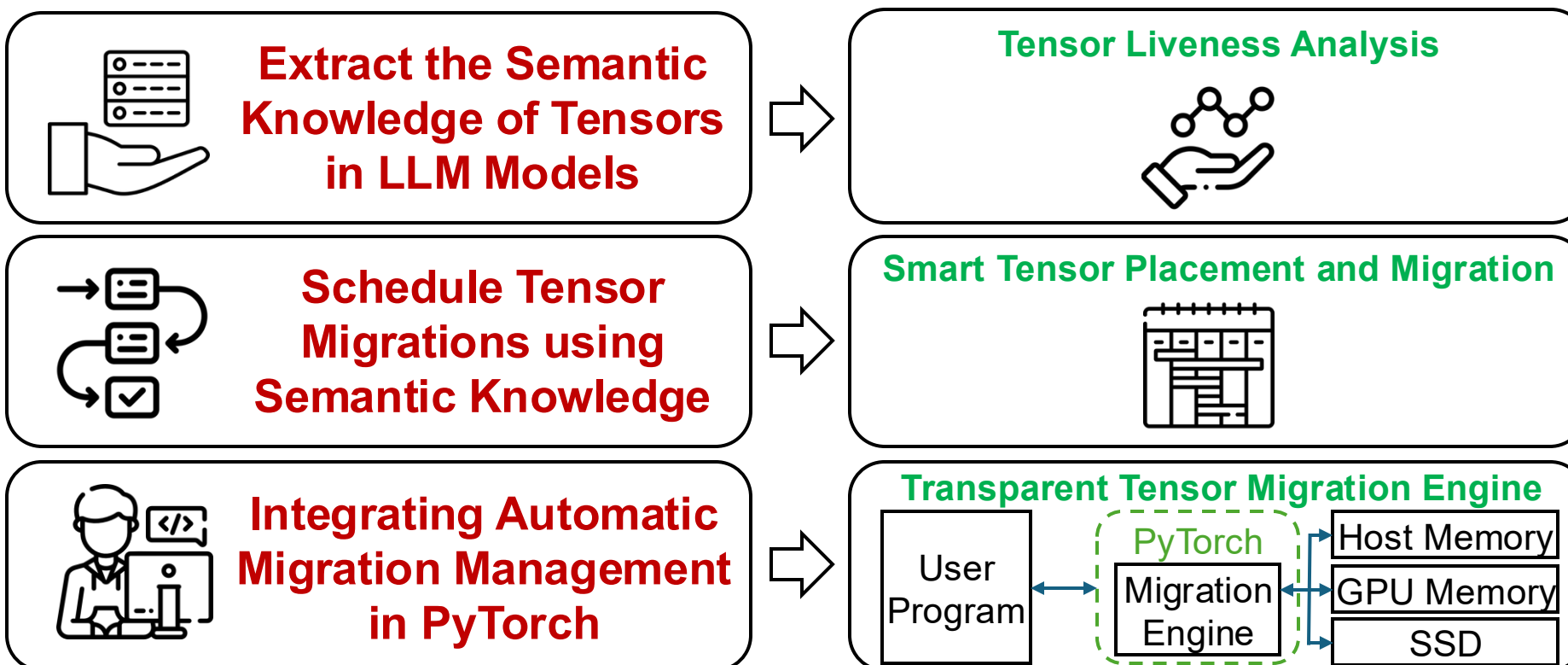
IV. Opportunities for Tensor Migration



Roofline Analysis Normalized training performance for different migration bandwidths

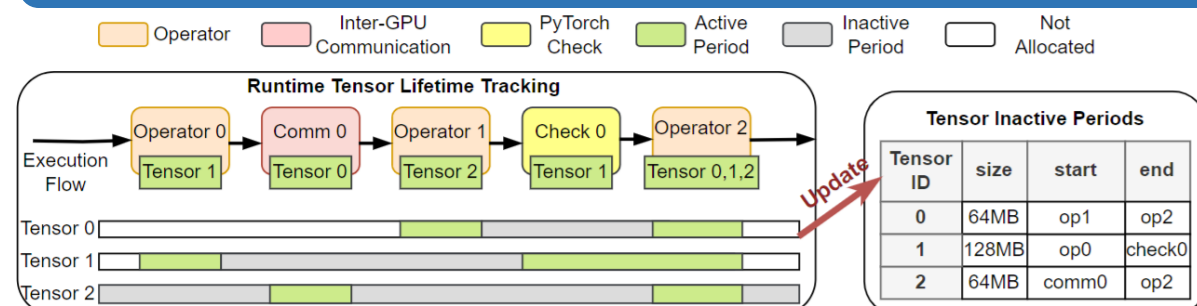
A bandwidth of 32 to 48 GB/s is sufficient to achieve near-ideal performance for large LLMs.

TeraIO



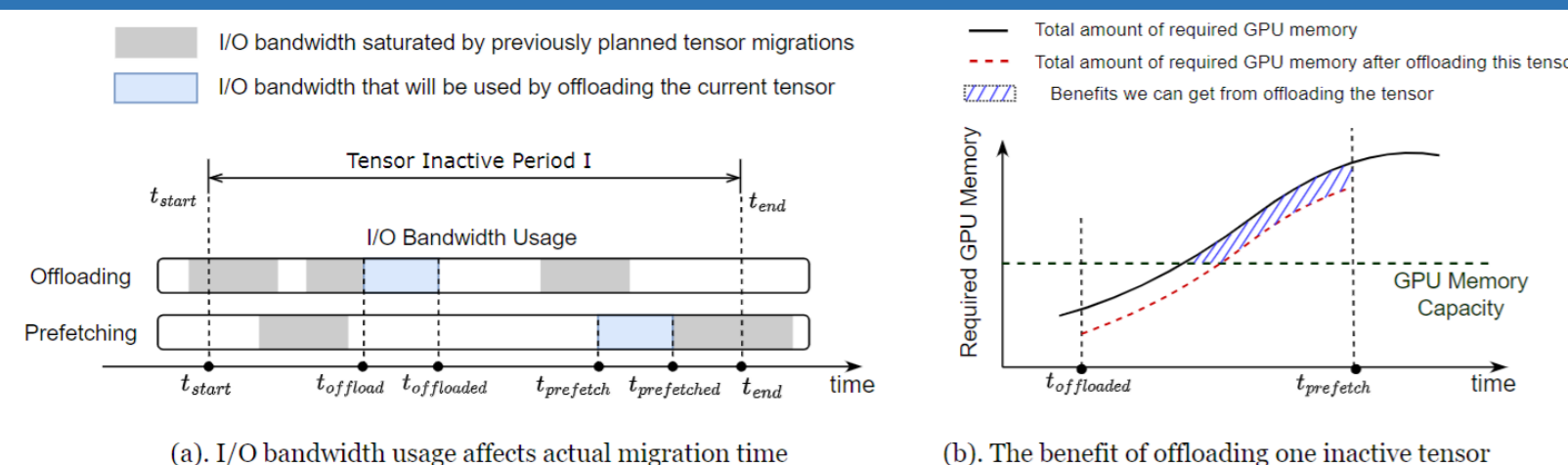
V. Enabling Smart Tensor Migrations with Rich Semantic Knowledge

Tensor Lifetime Profiler



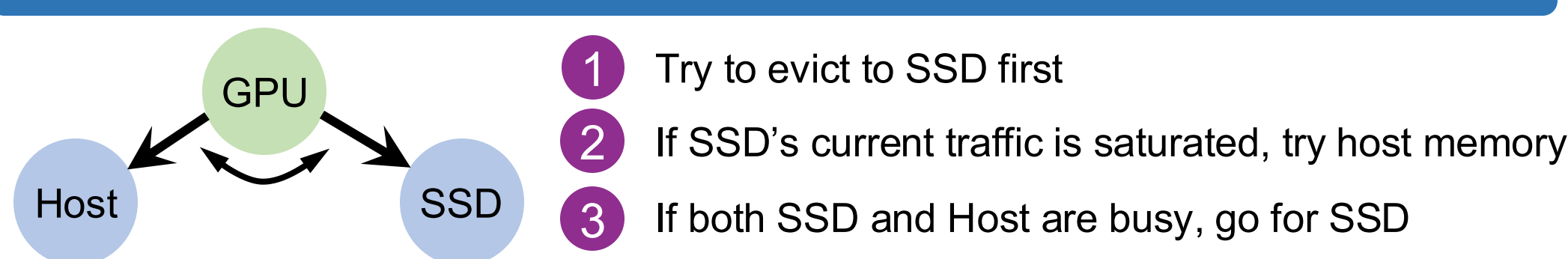
Run the model a few times and record op runtime and tensor lifetime information

I/O Bandwidth-Aware Migration Planning

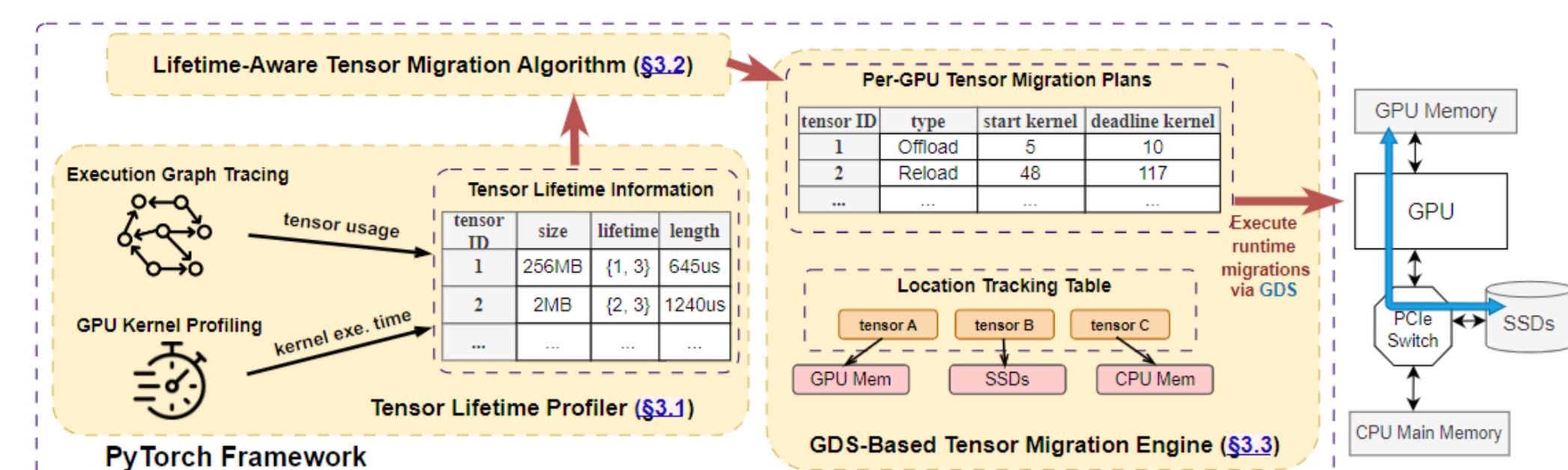


- 1 Estimate the impact of evicting an inactive tensor
- 2 Choose the one which will reduce the highest memory pressure (b), and cause the least I/O bandwidth pressure (a)
- 3 Update the "effects" of this decision

Decide Eviction Destination



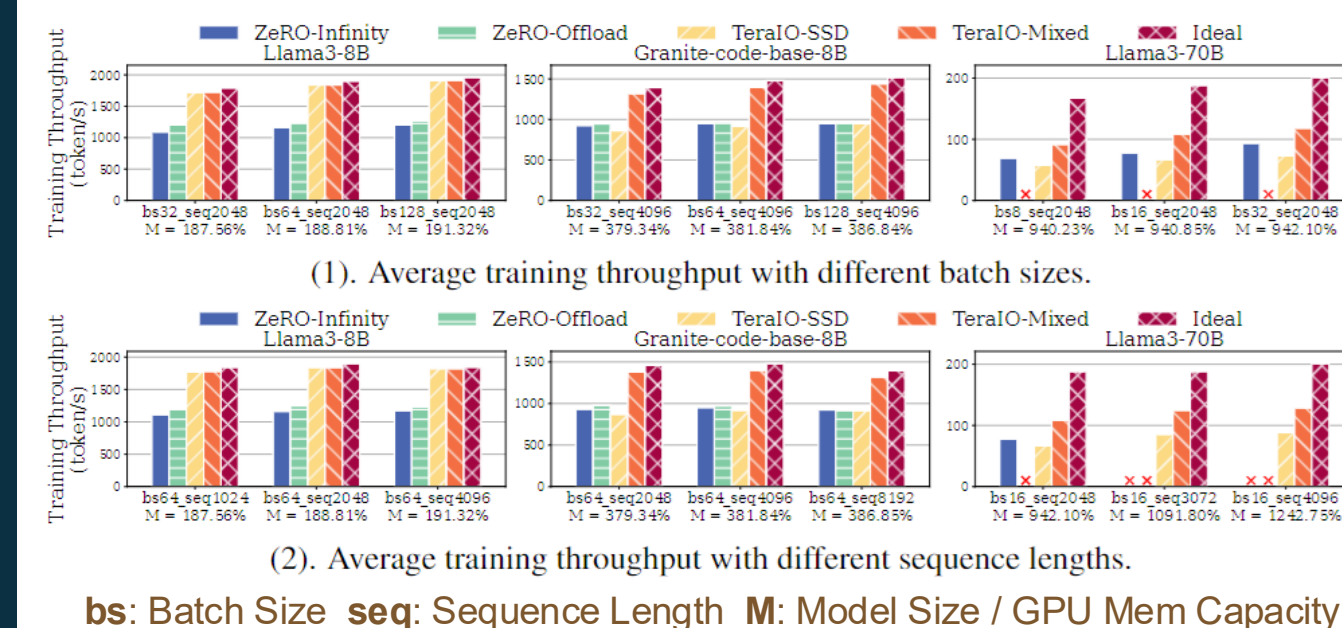
VI. Integrated Migration Engine in PyTorch



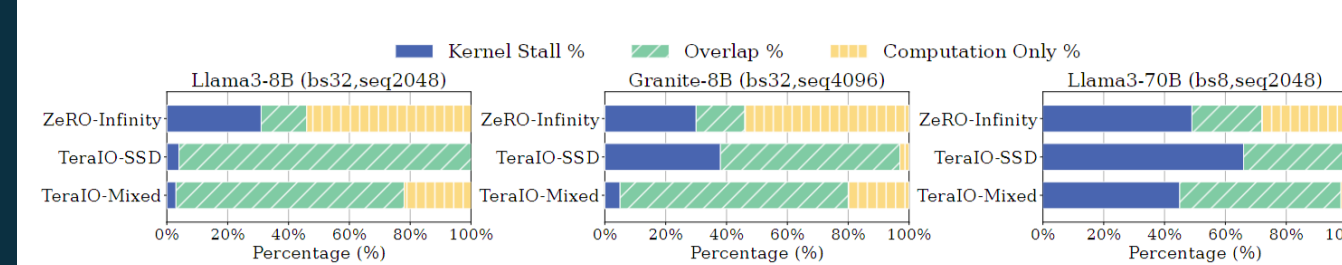
- 1 We place hooks before and after each operator to allow profiling and migration instrumentation
- 2 The migration schedule will be automatically executed before each operator is executed
- 3 We use GPUDirect Storage (GDS) when moving tensors between the GPU and SSDs

We use GDS because it provides significant scalability advantages by eliminating host resource contention

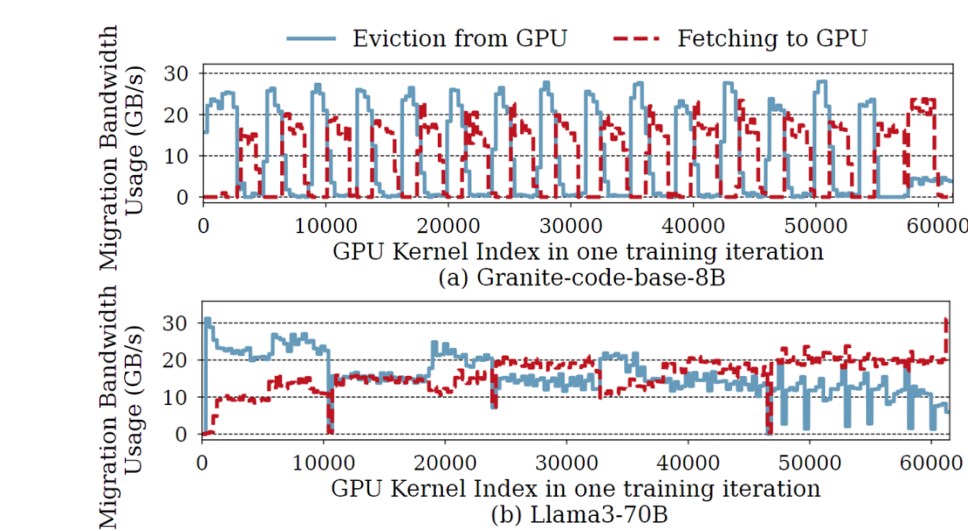
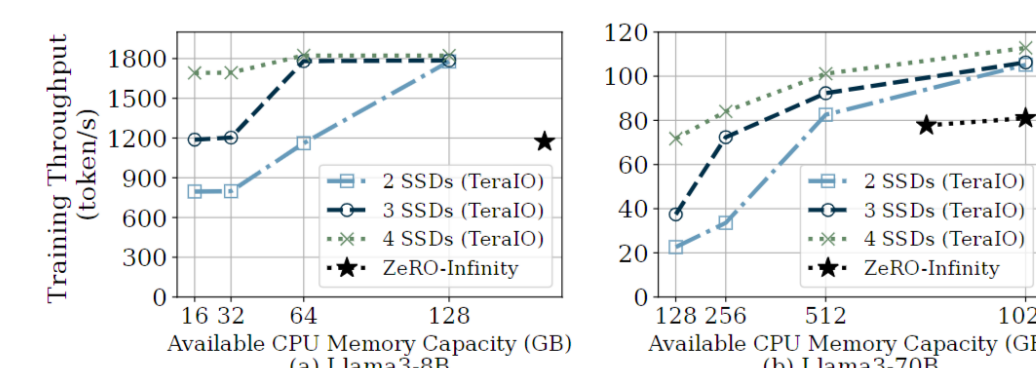
VII. Evaluation



bs: Batch Size seq: Sequence Length M: Model Size / GPU Mem Capacity



Latency breakdown of training iterations



With limited GPU memory, TeraIO achieves 80.7% of the ideal performance.

TeraIO incurs the least stall time, as it achieves better I/O and computation overlapping.

TeraIO can achieve better performance with lower CPU memory requirements.

TeraIO maintains high I/O bandwidth utilization with our proposed I/O-aware migration algorithm.