# Scaling GNN Sampling on Large-Scale Graphs with io_uring

Qixuan Chen
taliac@bu.edu

Yuhang Song
yuhangs@bu.edu

Melissa Martinez
melimtz@bu.edu

Vasiliki Kalavri
vkalavri@bu.edu

BOSTON UNIVERSITY

CASP Systems

## Graph Representation Learning and GNN Sampling

Node classification

"conservative"

"liberal"

follows

friend

Link prediction

Sampling → Aggregation → Model Training

sort & dedup

sort & dedup

epoch(target nodes for training)

| T1 | T2 | T3 | T4 | T5 | T6 | T7 | T8 | T9 | T10 | T11 | T12 | ......... | | | | | Tn |

batch 1 — batch 2 — batch 3 ......... batch n-1 — batch n

Sampling ■ Others ■

Pytorch-CPU ▮▮▮▮ 96.2
DGL-CPU ▮▮▮▮ 70.1

0 — % of Training time — 100

| Datasets | #Nodes | #Edges | Size |
|---|---|---|---|
| Papers100M | 111M | 1.62B | 70GB |
| Hyperlink | 3.5B | 128B | **3.4TB** |
| Facebook | 1.4B | 1T | **8.5TB** |

**Graph size may exceed main memory**

## Current Approaches and Limitations

### CPU-based
**(MariusGNN, Ginex)**

- Unnecessary I/O from loading full neighborhoods into memory
- Lower computation power than GPU
- High data movement overhead between memory and SSD

### GPU-based
**(NextDoor, gSampler)**

- Constrained GPU memory
- Sampling competes with other tasks for GPU resources
- High computation cost

### SSD-based
**(In-situ SmartSSD, FlashGNN)**

- Difficult to adopt widely across different hardware configurations
- Limited bandwidth compared to main memory

**My proposal: Leverage modern storage APIs and high-bandwidth SSDs to perform sampling on larger-than-memory graphs**

## CPU-Based Sampling System Accelerated by io_uring

1. **How to** minimize data movement between SSD and CPU?
2. **How to** integrate io_uring to GNN Sampling?
3. **How to** implement multi-threading to maximize CPU utilization?
4. **How to** take the advantage of asynchronous I/O?

**Main memory**

node 1 nrb offset range

offset index

| 0 | 0 | 5 | 9 | 11 | 12 | 14 | 16 | 16 | 17 | ... |
| 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | ... |

target index

| 1 | 32 | 22 | ... | 80 | 33 | ... | 17 | 5 | ... |

batch 1 — batch 2 — batch n

Offset Storage
Partition 1 ...... Partition n

Neighbor Storage
Partition 1 ...... Partition n

Target Storage
Partition 1 ...... Partition n

**CPU**

SQ — tail 2 1 0 — head — Thread 1 — CQ 6 3 2 — tail — head

...... SQ CQ Thread n

Polling

**Kernel**

**Disk**

Edge file

node 1 nrb     node 2 nrb

neighbors | 2 | 3 | 6 | 7 | 11 | 6 | 10 | 14 | 12 | 1 | 4 | 9 | 5 | 10 | 10 | ... |

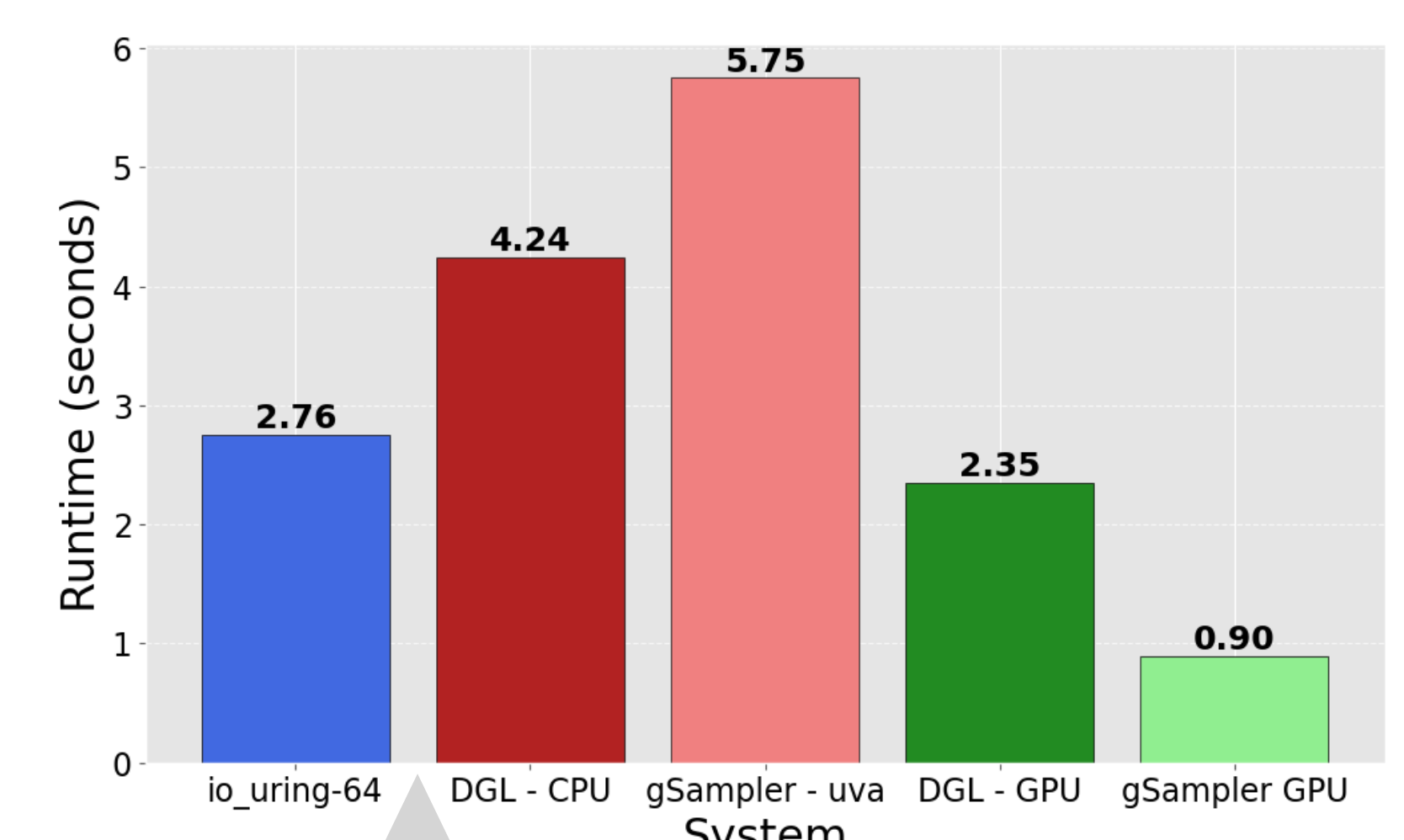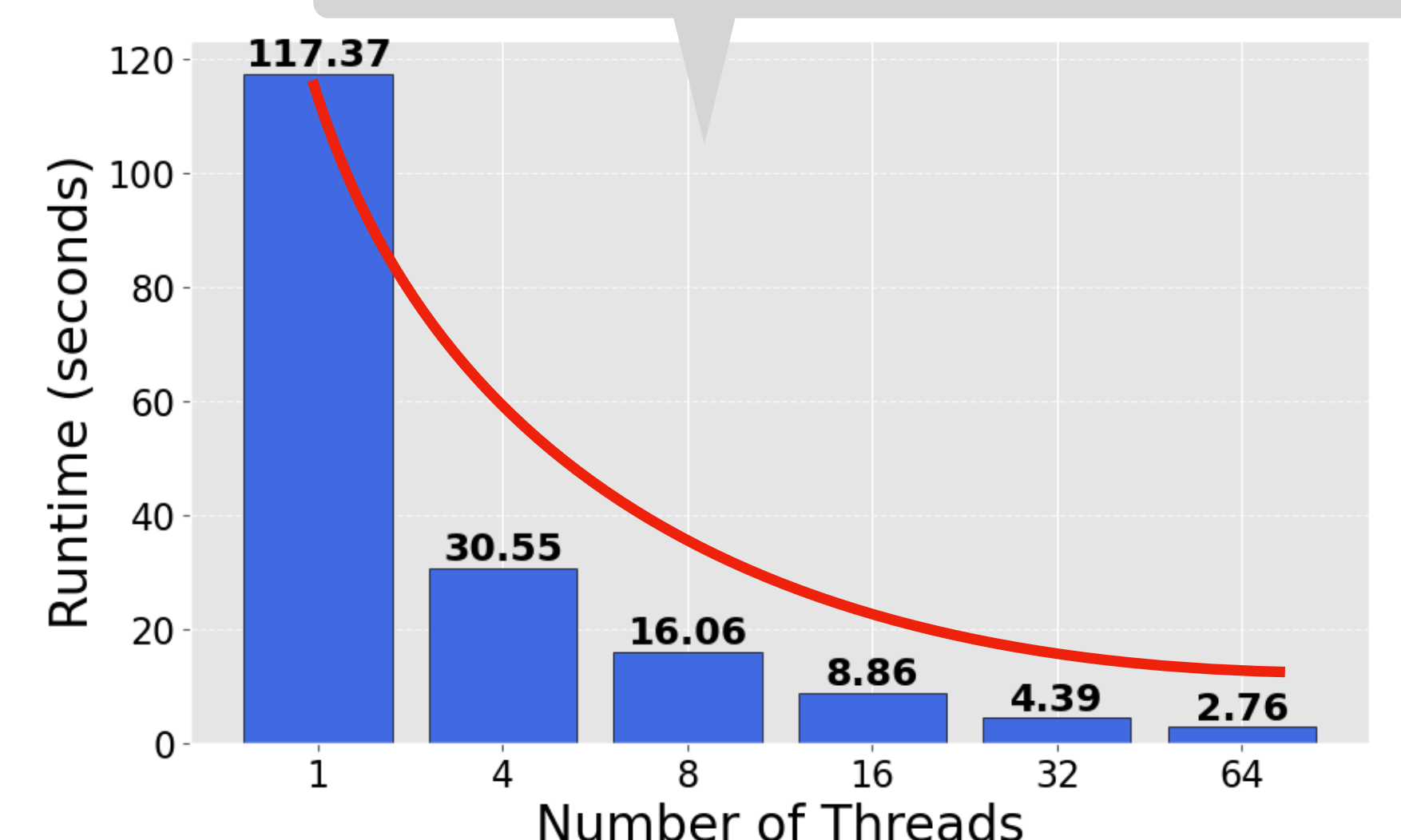index | 0 | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | ... |

1. Create two indexes to efficiently select neighbor offsets, allowing direct access to sampled neighbors without loading all neighbors from disk

2. Batch I/O requests to read neighbors

3. Equally distribute batches to threads. Each thread works in parallel with their own io_uring buffers without interrupting each other

4. While CQ is polling completions, we prepare and load the next I/O batch into the SQ, submitting it once the previous batch is done

**performance scales almost linearly with thread count**

Runtime (seconds) vs Number of Threads

| 1 | 117.37 |
| 4 | 30.55 |
| 8 | 16.06 |
| 16 | 8.86 |
| 32 | 4.39 |
| 64 | 2.76 |

Runtime (seconds) vs System

| io_uring-64 | 2.76 |
| DGL - CPU | 4.24 |
| gSampler - uva | 5.75 |
| DGL - GPU | 2.35 |
| gSampler GPU | 0.90 |

**Sampling billion-edge graph**
- 1.5x faster than baseline
- Comparable with GPU-based
- Can process larger-than-memory graphs