AGILE: Lightweight and Efficient Asynchronous GPU-SSD Integration

Zhuoping Yang*, Jinming Zhuang*, Xingzhen Chen*, Alex K. Jones†, and Peipei Zhou*

Brown University*; Syracuse University†;

zhuoping_yang@brown.edu peipei_zhou@brown.edu

https://peipeizhou-eecs.github.io/

https://github.com/arc-research-lab/AGILE





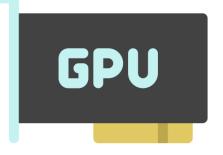














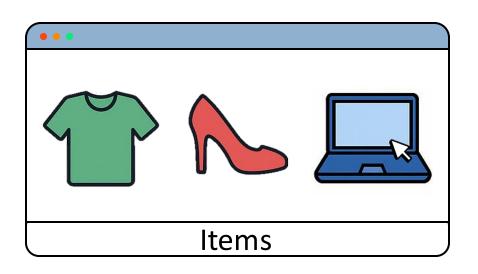
GPGPU: The engine behind modern applications.

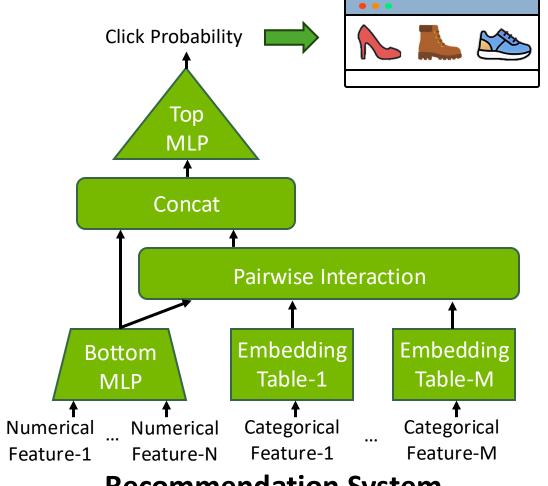






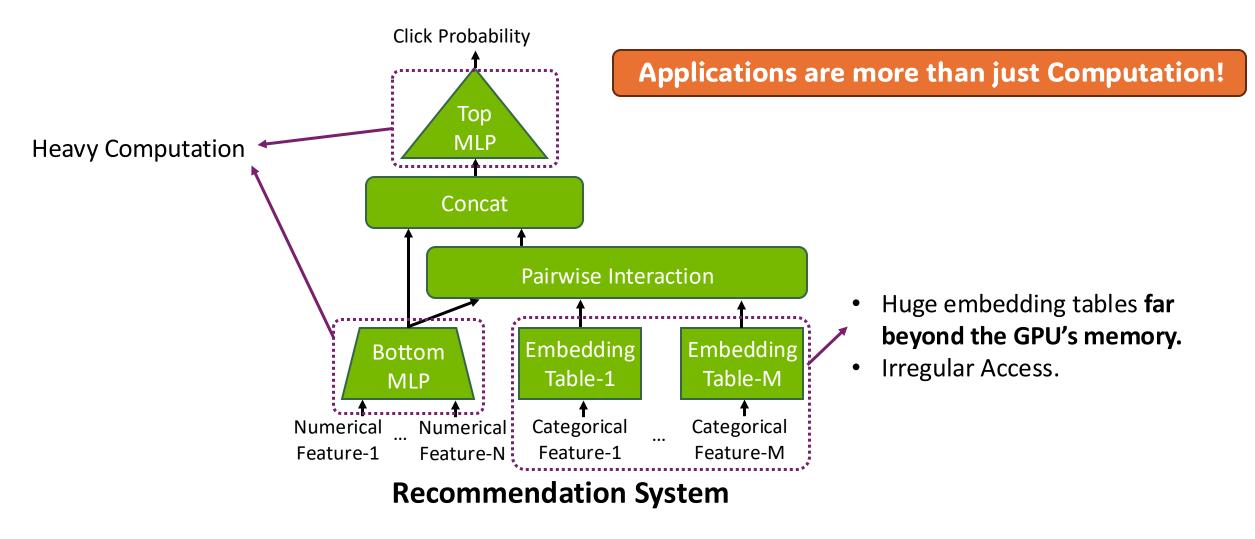


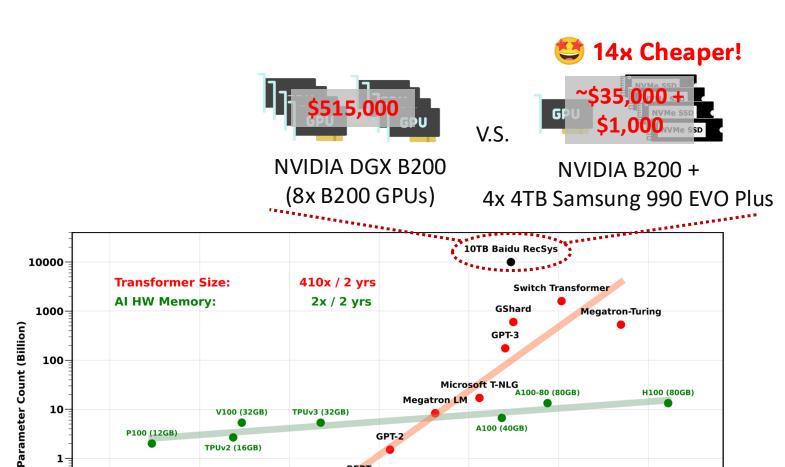




Recommendation System







Microsoft T-NLG

2020

A100 (40GB)

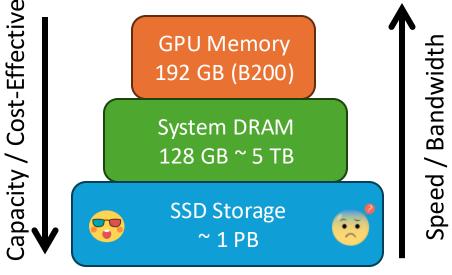
2021

Megatron LM

H100 (80GB)

2022





Hierarchical memory in modern computing systems

GPUs are increasingly constrained by the memory wall

2019 YEAR

Al and Memory Wall^[1]

GPT-2

TPUv3 (32GB)

Transformer

2018

V100 (32GB)

TPUv2 (16GB)

P100 (12GB)

Inception V4 ResNext101

DenseNet

2017

10-

0.1-

0.01

ResNet50

2016

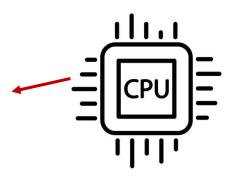
[1] Gholami, Amir, et al. "Ai and memory wall." *IEEE Micro* 44.3 (2024): 33-39.

BaM^[1] (A GPU-Centric Storage Access Model)



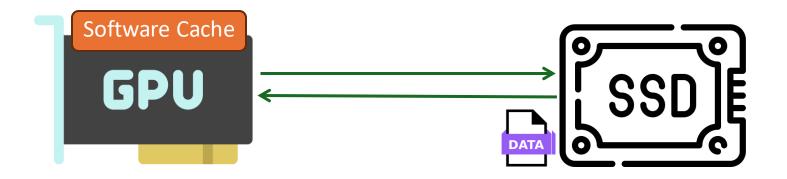
CPU is bypassed completely in both the control plane and the data plane.

No extra copies or synchronization overheads





Can we further improve the GPU-centric storage access model?



Programming Models for GPU-Centric I/O

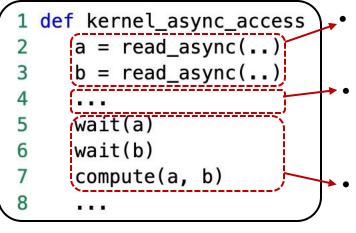


Async I/O enables overlapping the slow transfer with other tasks.

```
1 def kernel block access
    a = read_block(..)
    b = read_block(..)
    compute(a, b)
```

Synchronous I/O Model (BaM)

- Return after the slow transfer is finished.
- GPU thread must wait for the slow transfer.



Asynchronous I/O Model

Return quickly after the request is issued.

Switch to other tasks during the slow transfer.

Continue computation once the requested data is ready.



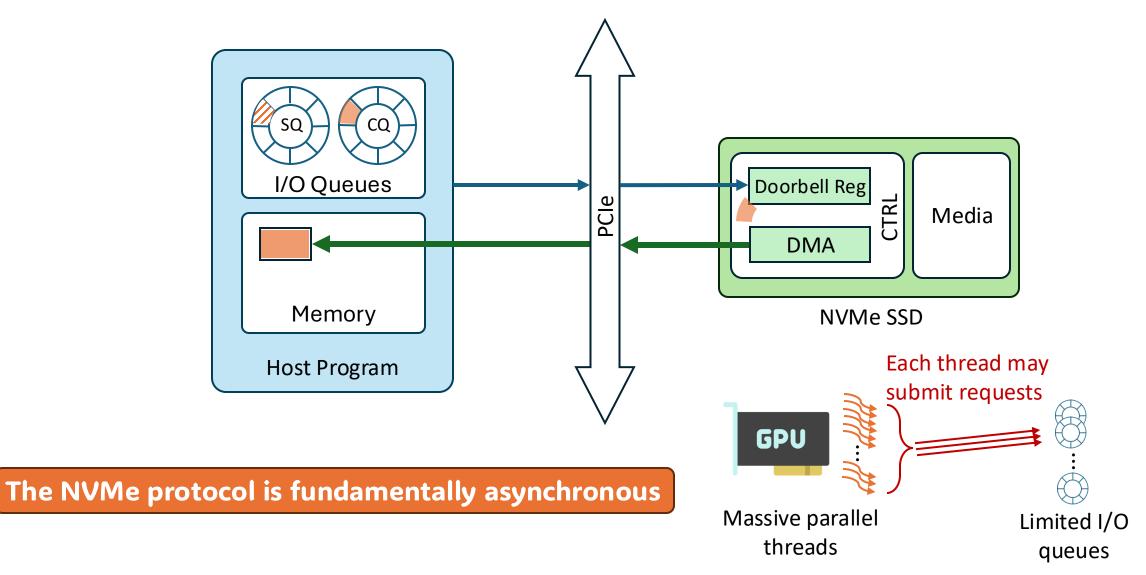
How can we support an async GPU-centric storage access model efficiently?

Use AGILE! (this work)



Background of NVMe Protocol





Deadlock Risks in Async GPU-Centric NVMe I/O



Async GPU-centric NVMe I/O can easily lead to a deadlock.

> A deadlock example in GPU-centric NVMe read:

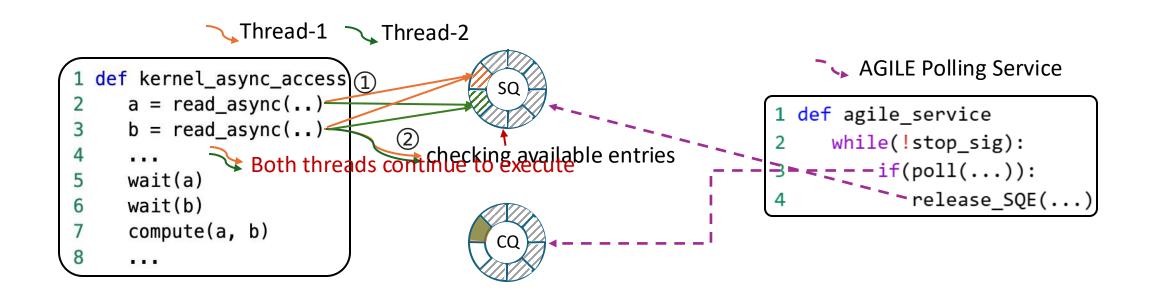
The root cause of this deadlock is allowing user threads to **hold multiple** locks that can block other user threads. (I/O queues, software cache lines)

AGILE Polling Service



To avoid deadlock from asynchronous NVMe I/O:

A warp-centric AGILE polling service running in the background



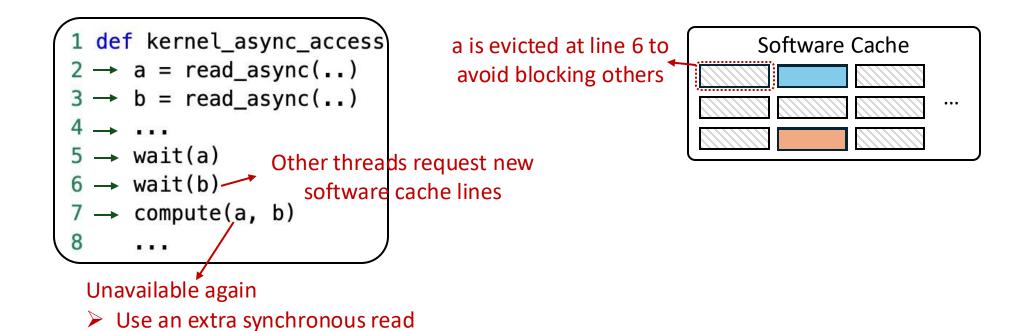
AGILE Software Cache

when accessing 'a'

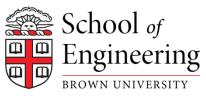


To eliminate deadlock from the software cache:

- AGILE does not allow user threads to avoid cache line eviction.
- ➤ Avoid Deadlocks at the cost of additional I/Os.



Integrate User-managed Buffers into AGILE



By default, all accesses are proxied by AGILE software cache.

- Requested data may be unavailable again to eliminate deadlocks.
- To address this problem, AGILE allows user-managed buffers.

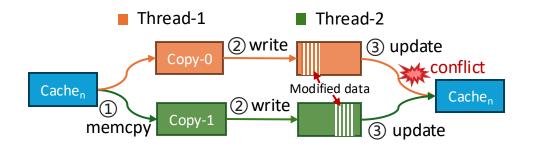
```
1 def kernel_async_access
2    a = read_async(..)
3    b = read_async(..)
4    ...
5    wait(a)
6    wait(b)
7    compute(a, b)
8    ...
```

Access via AGILE software cache (maybe require additional reads)

Integrate User-managed Buffers into AGILE



When multiple threads request the same data, data conflicts may occur.



AGILE provides a compile-time optional mechanism (share-table) to let user threads access the same user-managed buffer.

When the share-table is enabled at compile time:

- user-managed buffer (L1; managed by user; registered to share-table)
- Software cache (L2; managed by AGILE)
- SSDs (L3; managed by AGILE)

AGILE's API



Customizing software cache policy

```
class GPUCache:public GPUCacheBase < GPUCache > { . . . };
#define AGILE_CTRL AgileCtrl<GPUCache, ShareTable>
__global__
void kernel(AGILE_CTRL * ctrl, void * data){
                         → Debug option for
 AgileLockChain chain;
                           detecting deadlock
 // Method-1: AGILE prefetch
  ctrl->prefetch(dev_idx, blk_idx, chain);
  // Method-2: AGILE async_issue
  AgileBufPtr buf(data + tid * ctrl->line_size);
  ctrl->asyncRead(dev_idx, blk_idx, buf, chain);
  buf.wait();
  ctrl->asyncWrite(dev_idx, blk_idx, buf, chain);
 // Method-3: AGILE array-like synchronous API
  auto agileArr = ctrl->getArrayWrap<int>(chain);
  int val = agileArr[dev_idx][idx];
```

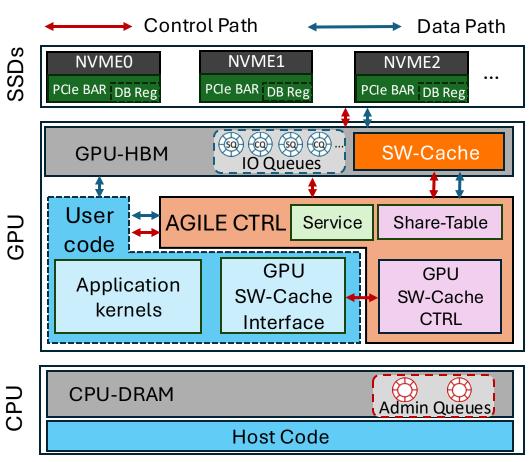
```
int main(int argc, char** argv){
   // GPU Configurations
    AGILE_HOST host(...);
    // Policy Configurations
    SHARE_TABLE_IMPL s_table(...);
   GPU_CACHE_IMPL g_cache(...);
    host.setGPUCache(g_cache);
   host.setShareTable(s_table);
   // Add and open target SSDs in the program
   host.addNvmeDev("/dev/AGILE-xxx", ...);
                                                  Init AGILE &
    host.addNvmeDev("/dev/AGILE-xxx", ...);
   host.initNvme();
                                                  Start AGILE
   // Initialize AGILE controller
                                                  service
   host.initializeAgile(...);
    // CUDA kernel parallelism configurations
    host.configKernelParallelism(...);
   host.queryOccupancy(kernel);
    // Start the lightweight AGILE service
   host.startAgile();
   // Execute the CUDA kernel
                                               Start user kernel
   host.runKernel(kernel, args...);
   // Stop AGILE service
   host.stopAgile();
                                               Stop AGILE
   // Close the opened SSDs
                                                service & close
   host.closeNvme();
                                                SSDs
```

Summary of AGILE



Feature Highlights:

- 1. AGILE is the first async GPU-centric I/O model, allowing GPU threads to access SSDs asynchronously.
- 2. AGILE eliminates the deadlock risks in asynchronous NVMe I/O via AGILE polling service.
- 3. AGILE allows users to customize software cache policies.
- 4. AGILE allows user-maintained buffers to be integrated into the software cache hierarchy.



Experimental Setup



- Dell R750 Server
 - Ubuntu 20.04 (Linux 5.4.0-200-generic)
- Nvidia RTX 5000 Ada GPU (PCle Gen4 x16)
- NVMe SSDs
 - Dell Ent NVMe AGN MU AIC 1.6TB SSD (PCIe Gen4 x4)
 - Two Samsung 990 PRO 1TB SSDs (PCle Gen4 x4)
- Software Cache Policy
 - Clock replacement algorithm^[1] (keep the same with BaM)



- 1. Compute and communication overlap.
 - 1 thread block issues 64 NVMe commands, and uses fetched data for compute.
- >AGILE can hide slow communication with computation effectively.

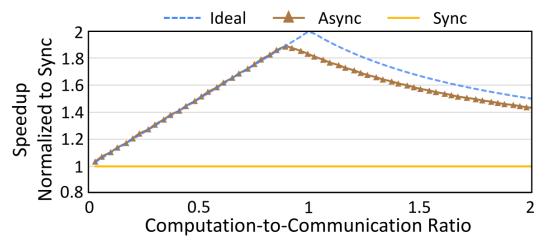


Figure 4: Speedup comparison of asynchronous I/O over synchronous I/O on workloads with different Computation-to-Communication Ratio (CTC).



- 2. Scalability in 4KB random read/write
 - Access different SSDs in an interleaving fashion.
- ➤ AGILE shows good scalability.

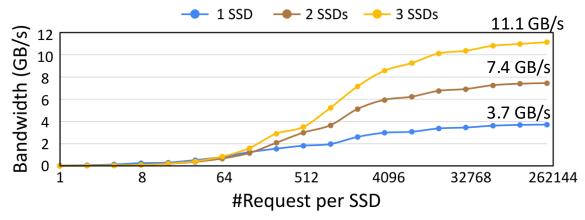


Figure 5: AGILE 4KB random read on multiple SSDs

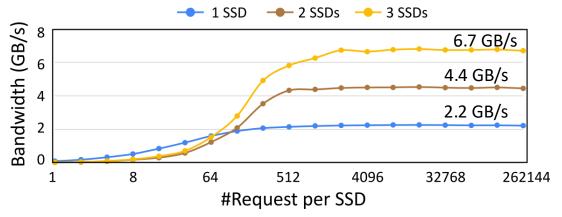


Figure 6: AGILE 4KB random write on multiple SSDs



3. Comparison with BaM on DLRM inference MicroBenchmark.

- The DLRM model is adopted from [1] with Criteo 1TB Click Logs dataset[2].
- cuBLAS is used for all matrix multiplications.
- BaM and AGILE are used for fetching data.
- AGILE is used in both synchronous mode and asynchronous mode

Default parameters:

- Batch Size: 2048
- #I/O queues: 128
- Software cache size: 2 GB

■ Sweep key parameter: Batch Size & the Number of I/O Queues.

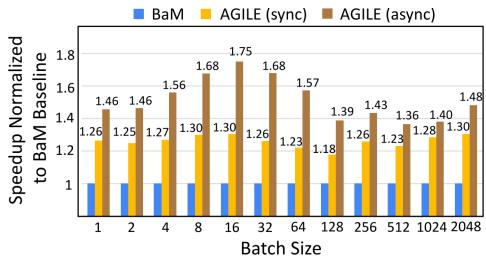


Figure 8: Speedup comparison of AGILE (async and sync modes) and BaM across varying batch sizes in DLRM inference.

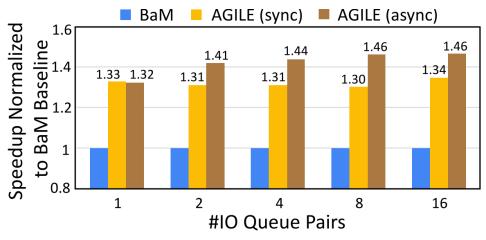


Figure 9: Speedup comparison of AGILE (async and sync modes) and BaM under varying numbers of I/O queue pairs in DLRM inference.

^[1] Naumov, Maxim, et al. "Deep learning recommendation model for personalization and recommendation systems." arXiv preprint arXiv:1906.00091 (2019).



3. Comparison with BaM on DLRM inference MicroBenchmark.

- The DLRM model is adopted from [1] with Criteo 1TB Click Logs dataset[2].
- cuBLAS is used for all matrix multiplications.
- BaM and AGILE are used for fetching data.
- AGILE is used in both synchronous mode and asynchronous mode
- Sweep key parameter: Software Cache Size

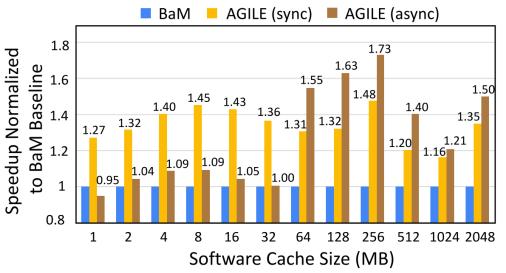


Figure 10: Speedup comparison of AGILE (async and sync modes) and BaM under varying software cache sizes in DLRM inference.

Default parameters:

- Batch Size: 2048
- #I/O queues: 128
- Software cache size: 2 GB

THANK YOU? QUESTIONS?

zhuoping_yang@brown.edu peipei_zhou@brown.edu

https://peipeizhou-eecs.github.io/

https://github.com/arc-research-lab/AGILE





