



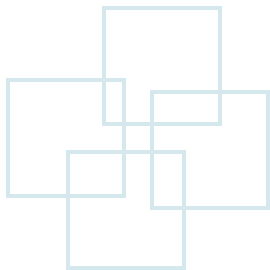
Scheduling-Aware Prefetching: Enabling the PCIe SSD to Extend the Global Memory of GPU Device

Tse-Yuan Wang^{1,2}, **Chun-Feng Wu^{1,2}**, **Che-Wei Tsao^{1,2}**, **Yuan-Hao Chang²**,
and **Tei-Wei Kuo^{1,3}**

¹Department of Computer Science and Information Engineering, National Taiwan University,

²Institute of Information Science, Academia Sinica

³Hong Kong Institute for Advanced Study, Department of Computer Science, City University of Hong Kong, Kowloon, Hong Kong



臺灣大學

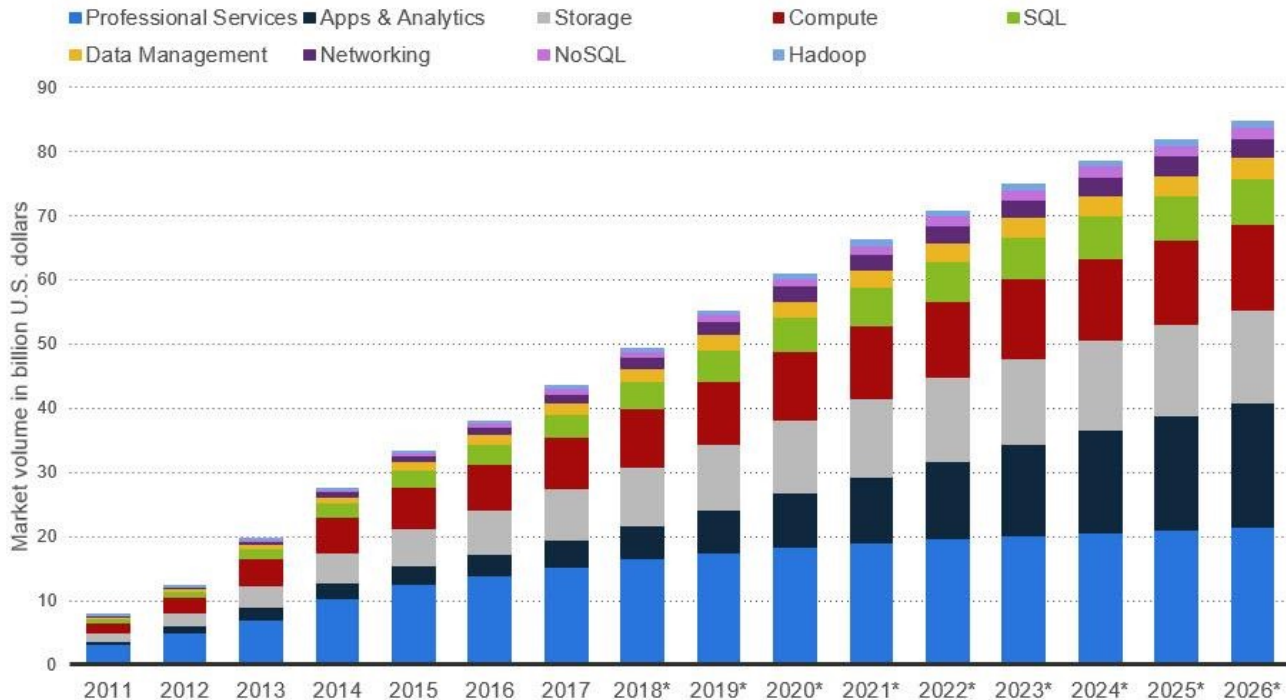


Outline

- Background and Motivation
- Scheduling-Aware Prefetching
- Performance Evaluation
- Conclusion



Big Data Market Forecast Worldwide from 2011 to 2026





Observation

- GPUs are widely used to accelerate these data-intensive big data applications, such as computer graphic and deep learning
- These applications also need a huge DRAM size of GPU device to as the temporary storage for data models
- The growth rate of the required memory size of applications is higher than the process technology of DRAM
 - High leakage power with huge DRAM size

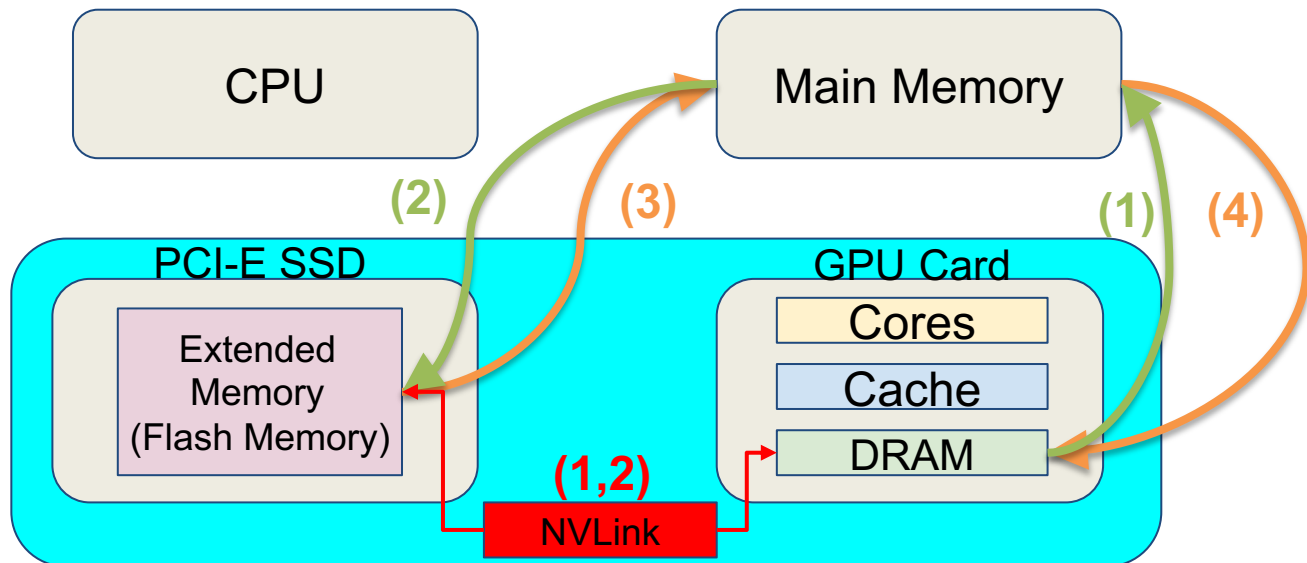


Observation

- NVIDIA provides Unified Memory Architecture (memories of Host and GPU inside the same address space)
 - Reducing the complexity of writing the CUDA program (programmers do not need to consider the data movement between Host and GPU device)
 - Ultra-scale memory drive with high access latency is not considered
- The new interface “RDMA” between the PCIe device and GPU device is provided to reduce the loading of host
- GPU is data parallelism, not task parallelism, and the segments of the program are decided to perform by the internal warp schedulers. The programmer and host are difficult to catch or create the data access locality.



The Complex Data Movement among Host, GPU Device, and PCIe SSD



Memory Expansion



Motivation

➤ Performance degradation

- Although GPUs can directly access PCIe SSD without moving the data to host DRAM beforehand, it leads to serious performance degradation due to the large performance gap between the PCIe SSD and the off-chip DRAM in the GPUs
- How to pre-fetch the data between DRAM and Flash Memory?

Type	DRAM	Flash SLC
Read time	60 ns	50 us
Write time	60 ns	550 us

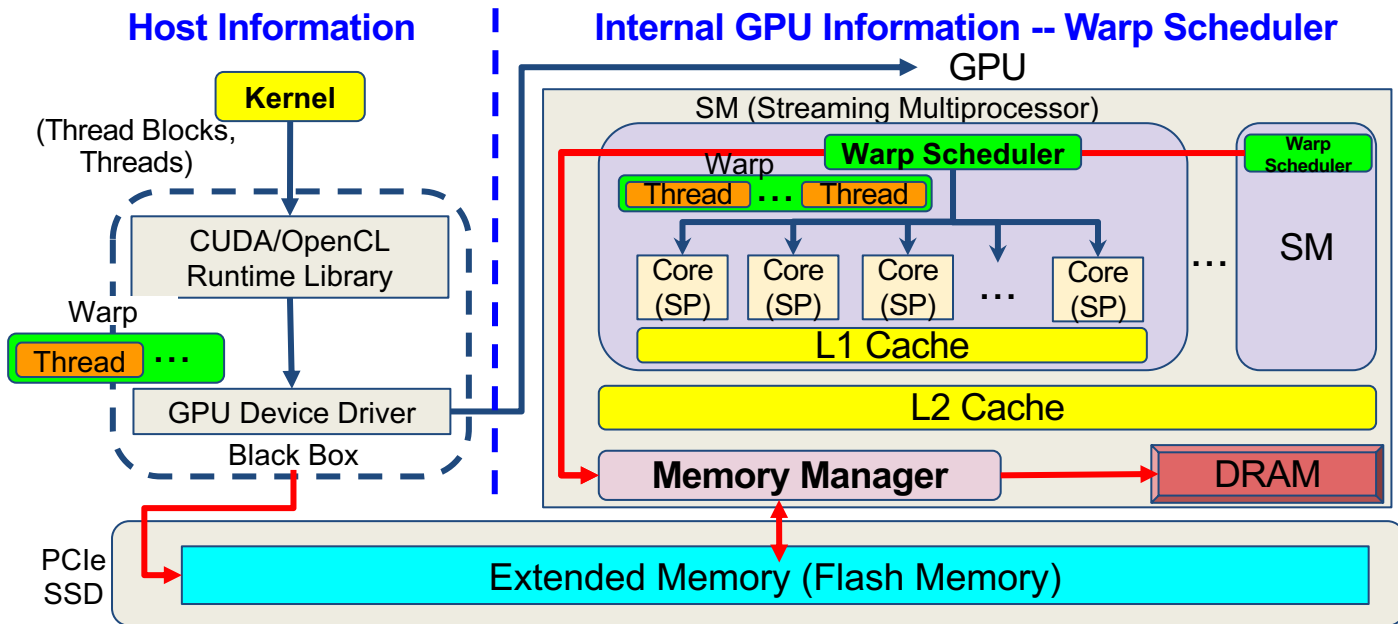


Outline

- Background and Motivation
- Scheduling-Aware Prefetching
- Performance Evaluation
- Conclusion

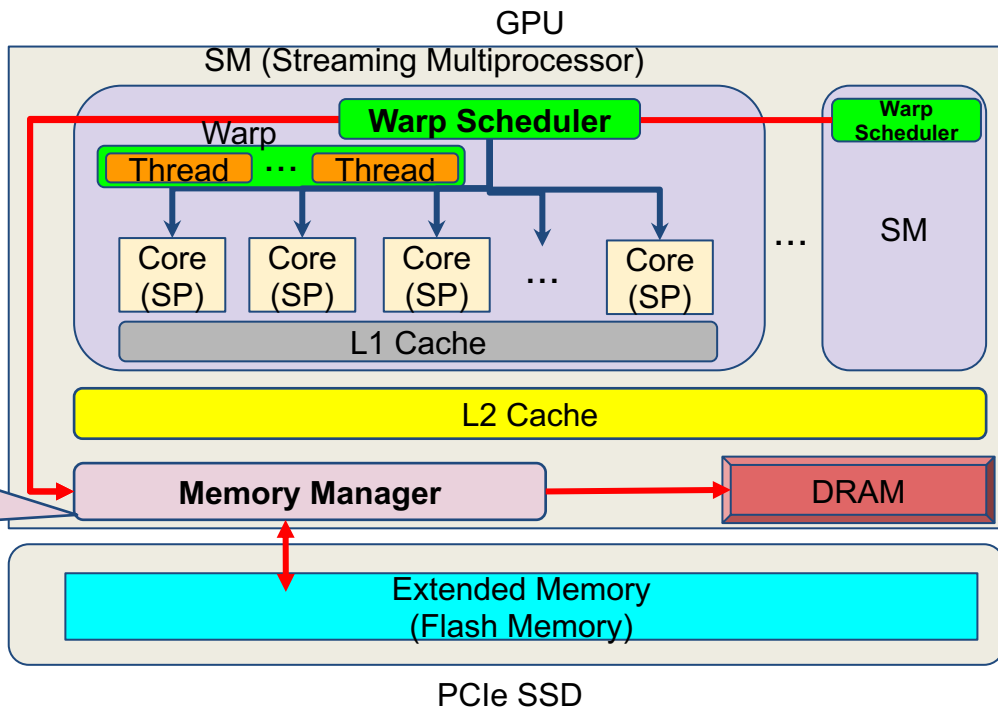
System Architecture

- Without modifying the host design (e.g., GPU device driver) and the programming way of the programmer
 - A new component **Memory Manager** is needed to coordinate between DRAM and SSD with considering the warp schedulers





Internal GPU Information -- Warp Scheduling





Internal GPU Information -- Warp Scheduling (cont.)

- The instructions are issued and executed per **warp**, and **memory operations** are also issued per warp
 - When executing a memory instruction, **each thread in a warp** provides a memory address which is loading or storing
- Memory Behavior Tracking
 - Warp = **32 threads**
 - The required memory addresses of threads in a warp are **continuous**
 - The memory request unit at once is based on the **warp**



Internal GPU Information -- Warp Scheduling (cont.)

- When op_code of the warp in warp scheduler is **LOAD** or **STORE**, the required memory addresses will notify **Memory Manager**
- When data of the accessed memory addresses are **missing in DRAM**, Memory Manager uses warp scheduler information to **pre-fetch** data to DRAM

```

scheduler : warp_id = 39, type = 0, addr = 2147545984, size = 128, write = 0
scheduler : warp_id = 40, type = 0, addr = 2147560448, size = 128, write = 0
scheduler : warp_id = 41, type = 0, addr = 2147560576, size = 128, write = 0
scheduler : warp_id = 42, type = 0, addr = 2147560704, size = 128, write = 0
exec : warp_id = 39, type = 0, addr = 2147545984, size = 128, write = 0
scheduler : warp_id = 43, type = 0, addr = 2147560832, size = 128, write = 0
exec : warp_id = 40, type = 0, addr = 2147560448, size = 128, write = 0
scheduler : warp_id = 44, type = 0, addr = 2147560960, size = 128, write = 0
exec : warp_id = 41, type = 0, addr = 2147560576, size = 128, write = 0
scheduler : warp_id = 45, type = 0, addr = 2147561088, size = 128, write = 0
scheduler : warp_id = 9, type = 0, addr = 2147501184, size = 128, write = 0
exec : warp_id = 42, type = 0, addr = 2147560704, size = 128, write = 0
scheduler : warp_id = 46, type = 0, addr = 2147561216, size = 128, write = 0
scheduler : warp_id = 13, type = 0, addr = 2147501696, size = 128, write = 0
exec : warp_id = 43, type = 0, addr = 2147560832, size = 128, write = 0
scheduler : warp_id = 47, type = 0, addr = 2147561344, size = 128, write = 0
scheduler : warp_id = 15, type = 0, addr = 2147501952, size = 128, write = 0
    
```

miss !

prefetch !

prefetch !!

miss !!



Outline

- Background and Motivation
- Scheduling-Aware Prefetching
- Performance Evaluation
- Conclusion



Experiment Setup

Evaluated Unified Memory

Type	Access Latency
DRAM Random Access Latency (Read/Write)	60 ns
Flash Memory Random Access Latency (Read)	50,000 ns
Flash Memory Random Access Latency (Write)	550,000 ns
Flash Memory Serial Access Time (ns/Bytes)	5 ns

Benchmarks

ISPASS2009-benchmarks – AES, BFS, CP, LPS, MUM, and NN

Memory Trace and Warp Scheduler Information are collected by the modified GPGPU-Sim

Compared Strategies: LRU and FIFO

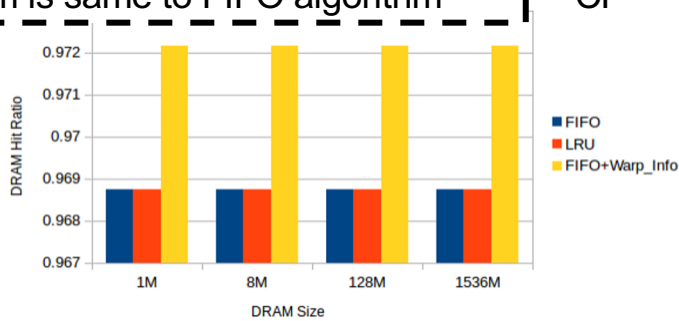
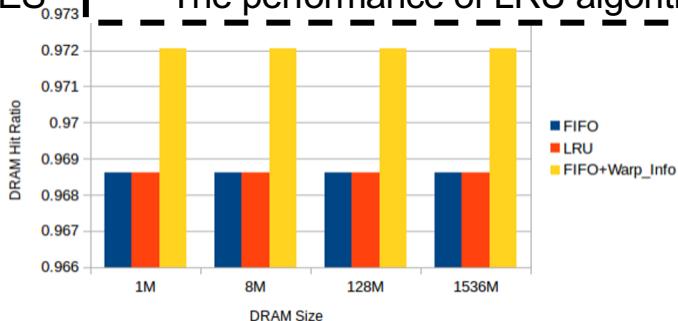


DRAM Hit Ratio on Different DRAM Size with Different Policy

AES

The performance of LRU algorithm is same to FIFO algorithm

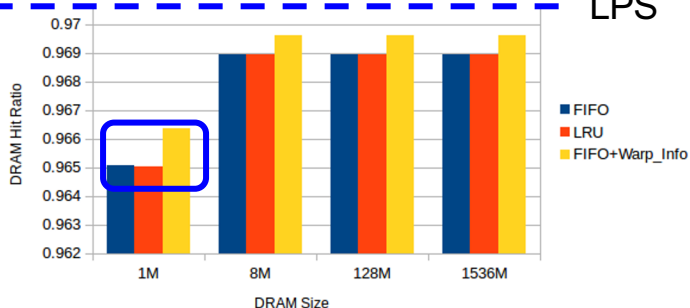
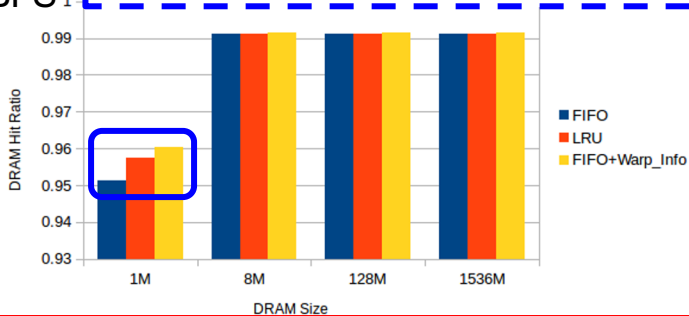
CP



BFS

The performance of LRU algorithm is higher/lower than FIFO algorithm

LPS



LRU algorithm is not suitable to manage the GPU Memory



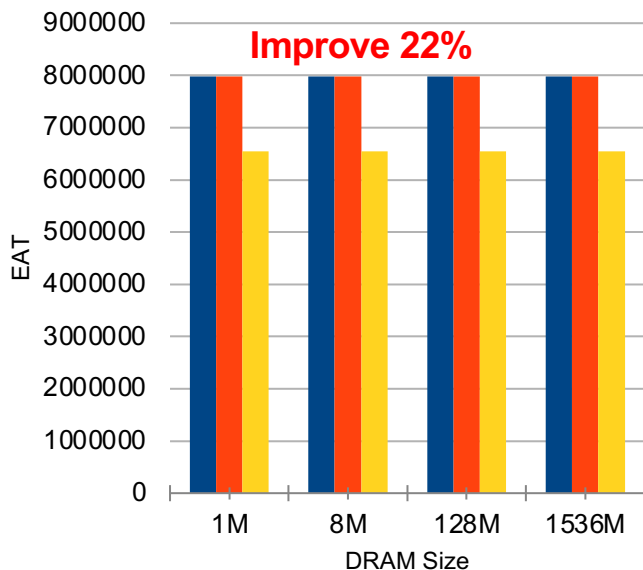
Effective Access Time

FIFO

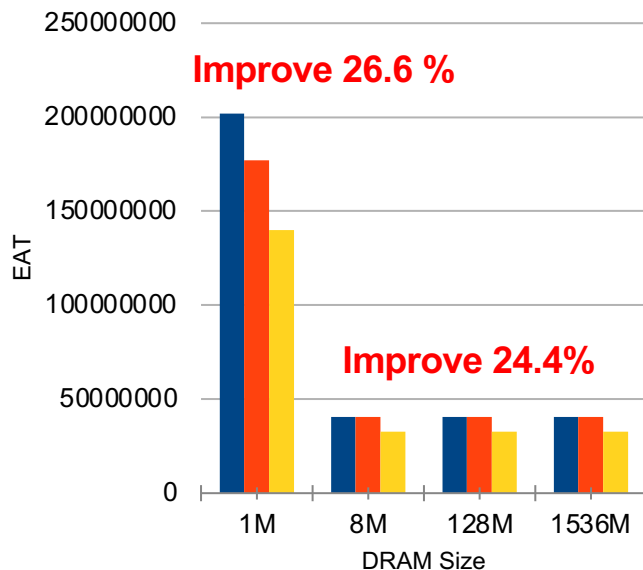
LRU

FIFO+Warp_Info

AES



BFS

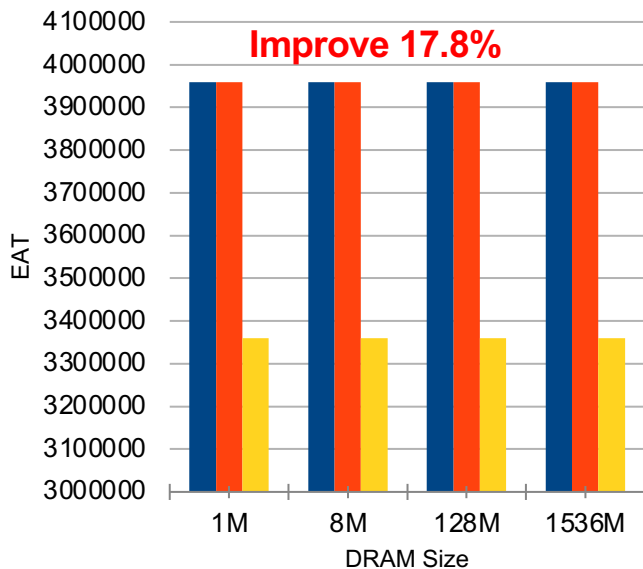




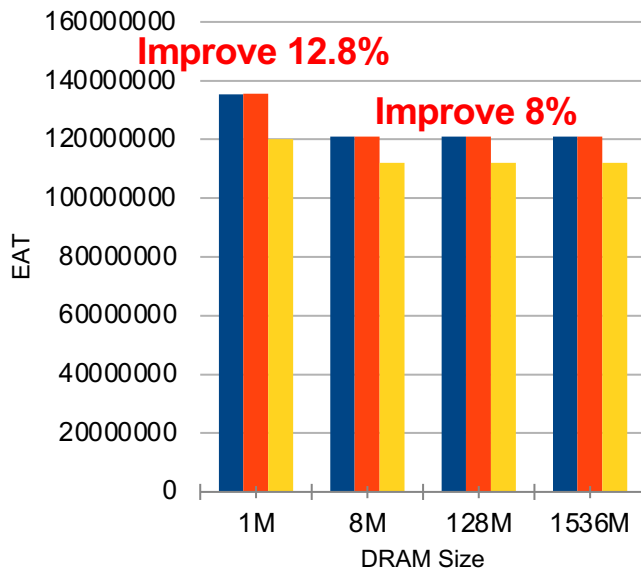
Effective Access Time (cont.)

- FIFO
- LRU
- FIFO+Warp_Info

CP



LPS

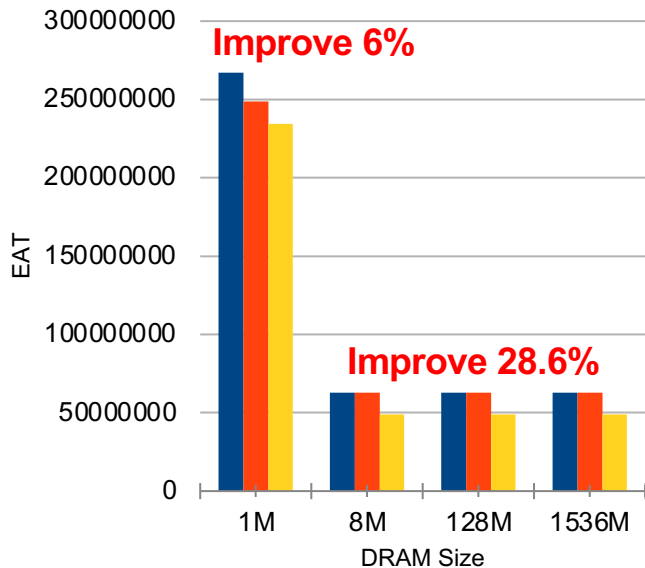




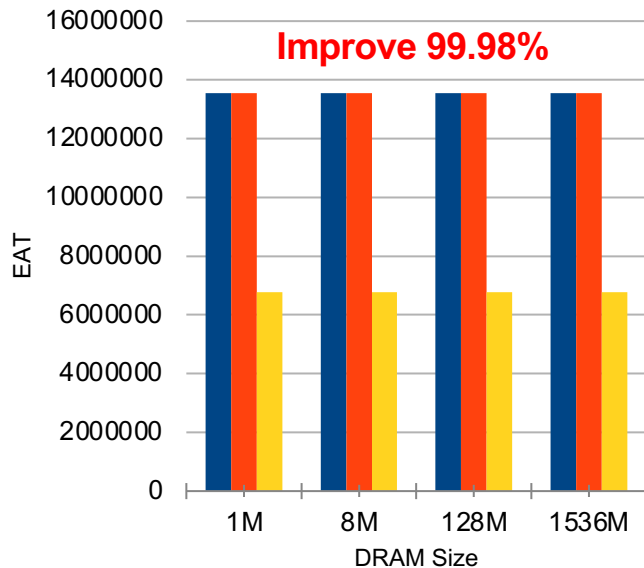
Effective Access Time (cont.)

- FIFO
- LRU
- FIFO+Warp_Info

MUM



NN



香港大學



CityU

Outline

- Background and Motivation
- Scheduling-Aware Prefetching
- Performance Evaluation
- Conclusion



Conclusion

- We proposed a scheduling-aware prefetching design in Unified Memory for GPU Device to exploit the process information provided by the warp scheduler so as to predict memory access patterns and perform accurate data prefetching
- The experiment results show that scheduling-aware prefetching design can efficiently improve the DRAM cache hit ratio
- The proposed scheduling-aware prefetching design can improve up to **99%** of effective access time compared to LRU algorithm.



Thanks for your attention

Q&A

Contact: Tse-Yuan Wang – tseyuan20@iis.sinica.edu.tw

Thank you