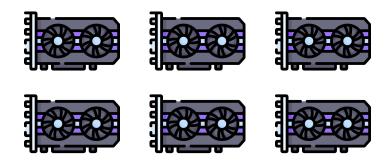# GPHash: An Efficient Hash Index for GPU with Byte-Granularity Persistent Memory

**Menglei Chen**, Yu Hua, Zhangyu Chen, Ming Zhang, Gen Dong

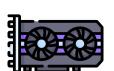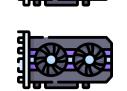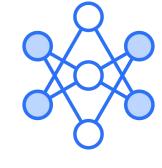*Huazhong University of Science and Technology, China*

# Various Applications Are Powered by GPUs

# Various Applications Are Powered by GPUs

*GPUs*

# Various Applications Are Powered by GPUs

*GPUs*



**Enhance**

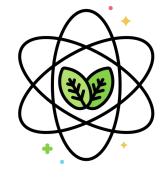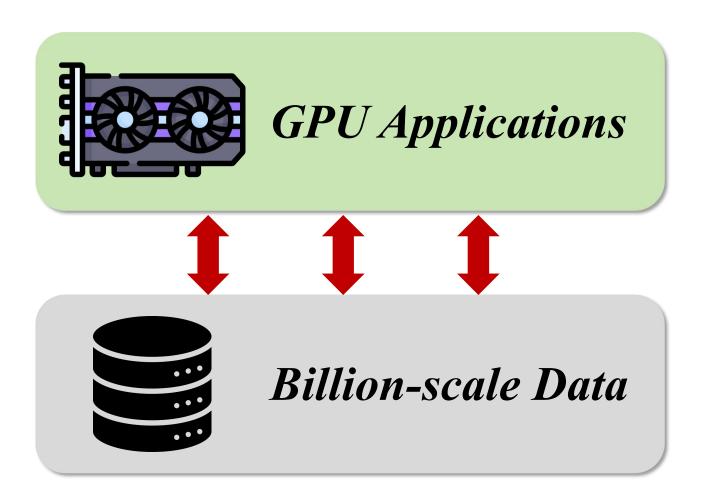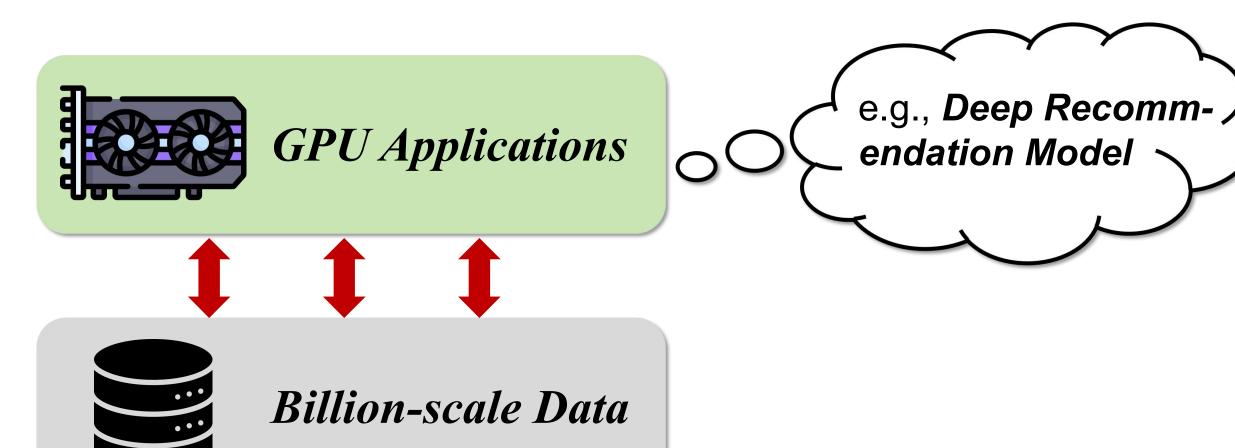# Various Applications Are Powered by GPUs



GPUs

Enhance

Deep Neural Network

Scientific Computing

Autonomous driving

# Unprecedented Data Scale of GPU App.

# Unprecedented Data Scale of GPU App.


GPU Applications

# Unprecedented Data Scale of GPU App.



GPU Applications

Billion-scale Data

# Unprecedented Data Scale of GPU App.

GPU Applications

Billion-scale Data

e.g., *Deep Recomm-endation Model*

# Unprecedented Data Scale of GPU App.



GPU Applications

e.g., *Deep Recomm-endation Model*

Billion-scale Data

e.g., *Billions of Emb-edding Vectors*

# Existing Data Management of GPU App.

# Existing Data Management of GPU App.

*GPU*

*Storage*

# Existing Data Management of GPU App.



GPU

CPU

Storage

PCIe

PCIe

DDR

DRAM

# Existing Data Management of GPU App.



*GPU*

*CPU*

*Storage*

PCIe

PCIe

*GPU Memroy*

DDR

*DRAM*

# Existing Data Management of GPU App.

**GPU**

**CPU**

**Storage**

PCIe

PCIe

**CPU**

**SSD**

DDR

**GPU Memroy**

*Managed by indexes*

**DRAM**

# Existing Data Management of GPU App.



GPU

CPU

Storage

PCIe

Data

DDR

DRAM

GPU Memroy

Managed by indexes

4

# Existing Data Management of GPU App.

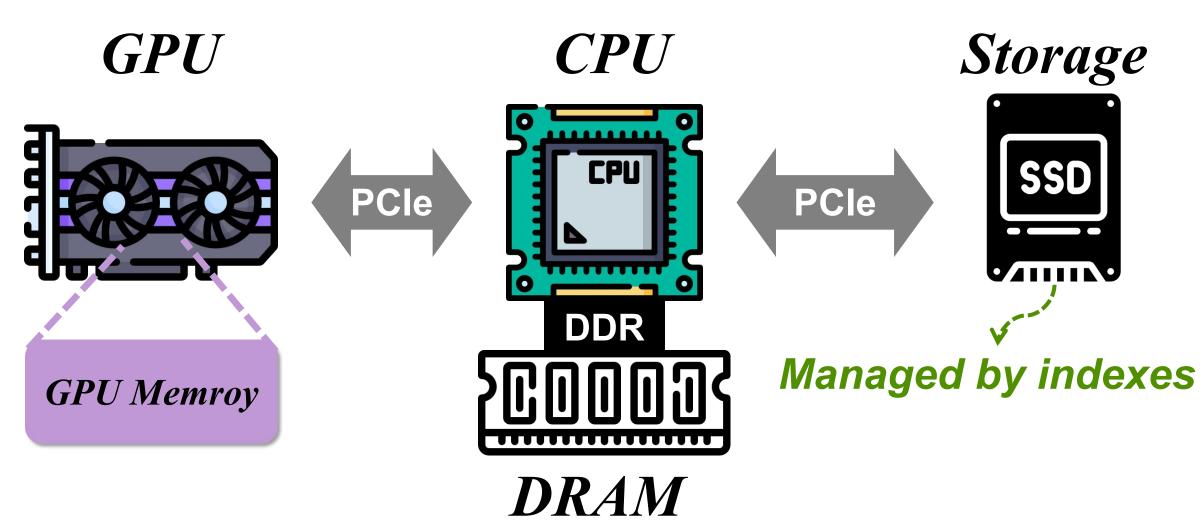# Existing Data Management of GPU App.

✔ **Large capacity and persistence**
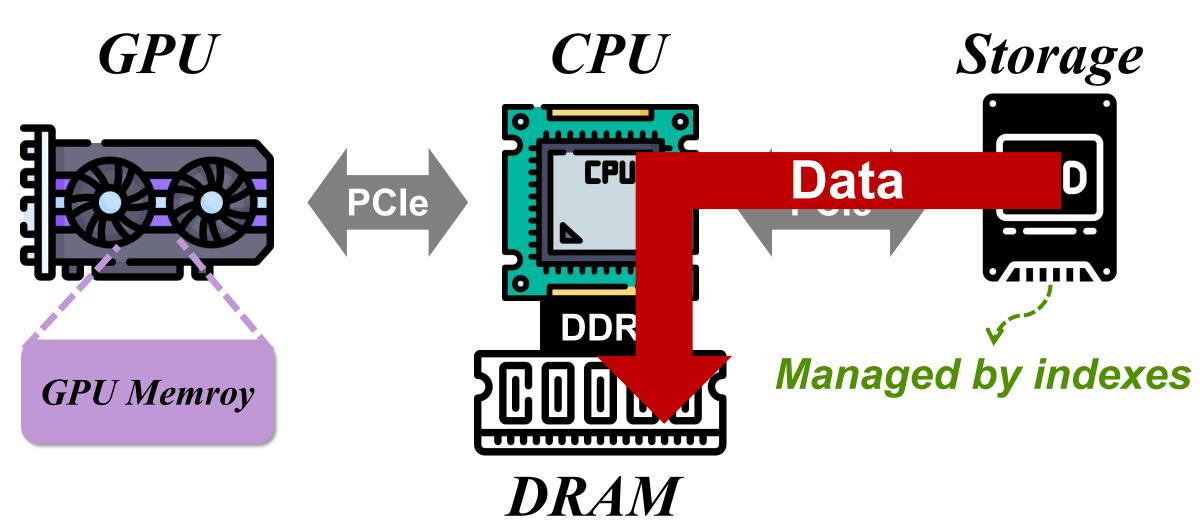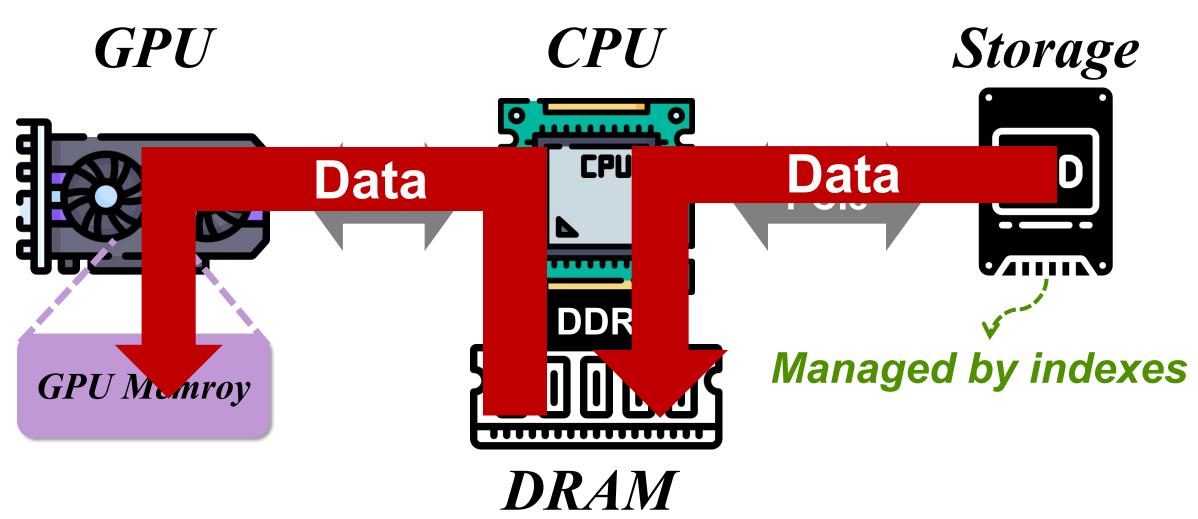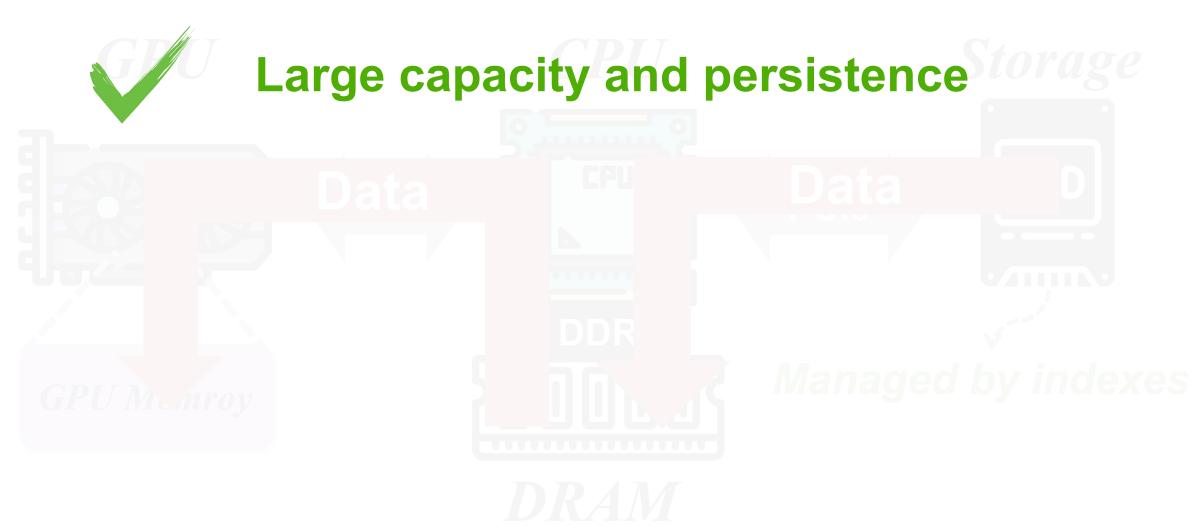
# Existing Data Management of GPU App.

✔ **Large capacity and persistence**

✘ **High overhead for data transfer**

✘ **Extra CPU consumption**

# Direct Data Access from GPU

# Direct Data Access from GPU

*GPU*

*Storage*

# Direct Data Access from GPU

*GPU*

*Storage*

GPU Direct Storage (PCIe)

# Direct Data Access from GPU



GPU

Storage

GPU Direct Storage (PCIe)

**Block-level** *POSIX-liked interfaces*

# Direct Data Access from GPU

# Direct Data Access from GPU

✓ **Large capacity and persistence**

✓ **Cost-efficient data transfer**

# Direct Data Access from GPU

✔ **Large capacity and persistence**

✔ **Cost-efficient data transfer**

✘ **Hard to program data structure**

✘ **Transfer extraneous data**

5

# GPU with Persistent Memory (GPM)

[1] GPM: Leveraging Persistent Memory from a GPU [ASPLOS' 22]
[2] Scoped Buffered Persistency Model for GPUs [ASPLOS' 23]

# GPU with Persistent Memory (GPM)

*GPU*

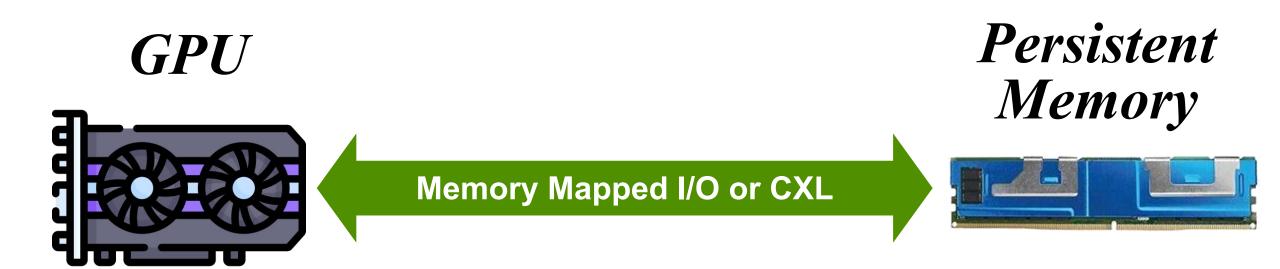*Persistent Memory*

[1] GPM: Leveraging Persistent Memory from a GPU [ASPLOS' 22]
[2] Scoped Buffered Persistency Model for GPUs [ASPLOS' 23]

6

# GPU with Persistent Memory (GPM)



*GPU*

*Persistent Memory*

**Memory Mapped I/O or CXL**

[1] GPM: Leveraging Persistent Memory from a GPU [ASPLOS' 22]
[2] Scoped Buffered Persistency Model for GPUs [ASPLOS' 23]

# GPU with Persistent Memory (GPM)

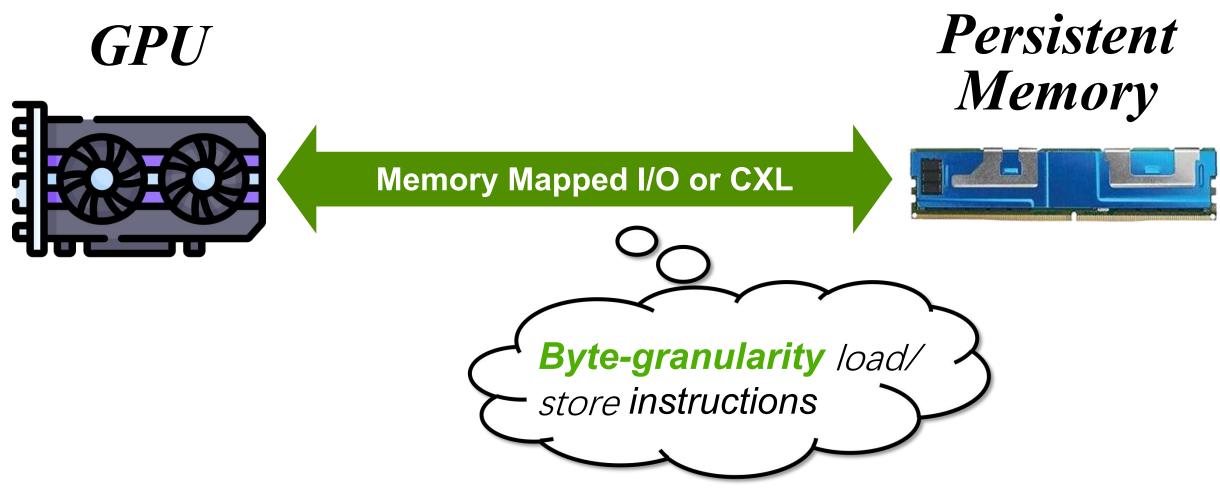**GPU**

**Persistent Memory**

**Memory Mapped I/O or CXL**

*Byte-granularity* load/store *instructions*

[1] GPM: Leveraging Persistent Memory from a GPU [ASPLOS' 22]
[2] Scoped Buffered Persistency Model for GPUs [ASPLOS' 23]

6

# GPU with Persistent Memory (GPM)

GPU

*Persistent Memory*

Memory Mapped I/O or CXL

*Byte-granularity* load/ store instructions

[1] GPM: Leveraging Persistent Memory from a GPU [ASPLOS' 22]
[2] Scoped Buffered Persistency Model for GPUs [ASPLOS' 23]

# GPU with Persistent Memory (GPM)

✓ **Large capacity and persistence**

✓ **Cost-efficient and fine-grained data transfer**

✓ **Easy to program data structure**

[1] GPM: Leveraging Persistent Memory from a GPU [ASPLOS' 22]
[2] Scoped Buffered Persistency Model for GPUs [ASPLOS' 23]

6

# Hash Index

# Hash Index

# Hash Index

**Hash Function**

# Hash Index

**Hash Function**

✓ **Constant-scale point query**

✓ **Good for parallel access**

# Hash Index

**Hash Function**

*GPM Hash Index*

✓ **Constant-scale point query**
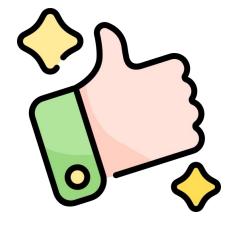
✓ **Good for parallel access**
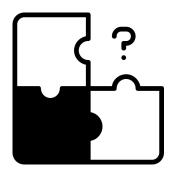
# Hash Index

**Hash Function**

*GPM Hash Index*

✓ **Constant-scale point query**

✓ **Good for parallel access**

*However, it is non-trivial to implement an efficient GPM hash index*

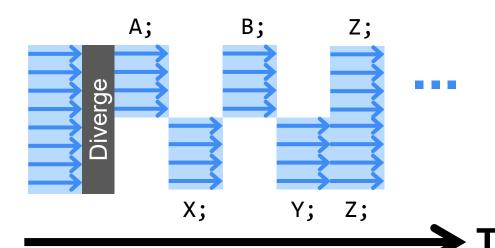# Challenge 1: Agnostic to GPU Execution Manner

# Challenge 1: Agnostic to GPU Execution Manner

*Warp Divergence*

# Challenge 1: Agnostic to GPU Execution Manner

*Warp Divergence*

```
if (thread_id < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;
```

*Warp Divergence*

```
if (thread_id < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;
```

Time

# Challenge 1: Agnostic to GPU Execution Manner

*Warp Divergence*

```
if (thread_id < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;
```

Diverge

**Time**

# Challenge 1: Agnostic to GPU Execution Manner

*Warp Divergence*

```
if (thread_id < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;
```



**Time**

# Challenge 1: Agnostic to GPU Execution Manner

*Warp Divergence*

```
if (thread_id < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;
```

A;                B;              Z;

Diverge

X;                Y;      Z;

Time

# Challenge 1: Agnostic to GPU Execution Manner

*Warp Divergence*

```
if (thread_id < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;
```

A;            B;            Z;
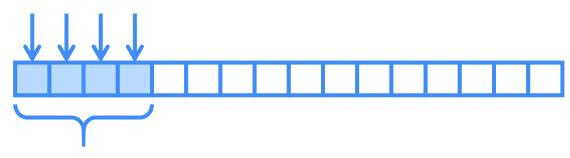
Diverge

X;            Y;   Z;

Time

*Coalesced memory accesses*

# Challenge 1: Agnostic to GPU Execution Manner

*Warp Divergence*

```
if (thread_id < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;
```

A;        B;        Z;

Diverge

X;        Y;   Z;

Time

*Coalesced memory accesses*

# Challenge 1: Agnostic to GPU Execution Manner

# Challenge 1: Agnostic to GPU Execution Manner

*Warp Divergence*

```
if (thread_id < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;
```

*Coalesced memory accesses*

**Memory Block**

Time

# Challenge 1: Agnostic to GPU Execution Manner



*Warp Divergence*

```
if (thread_id < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;
```
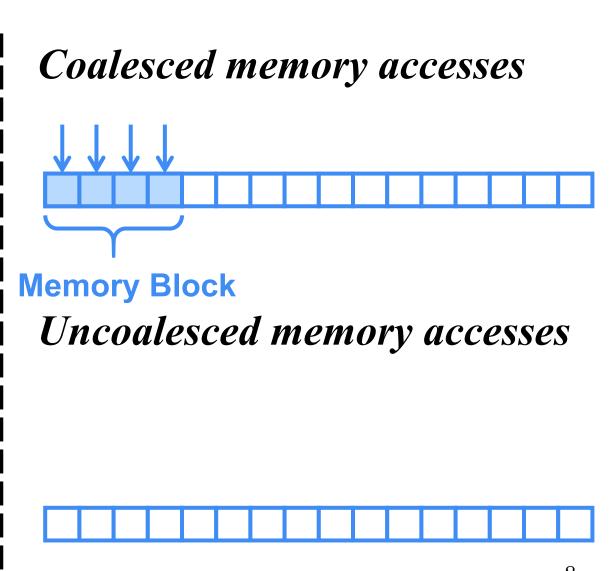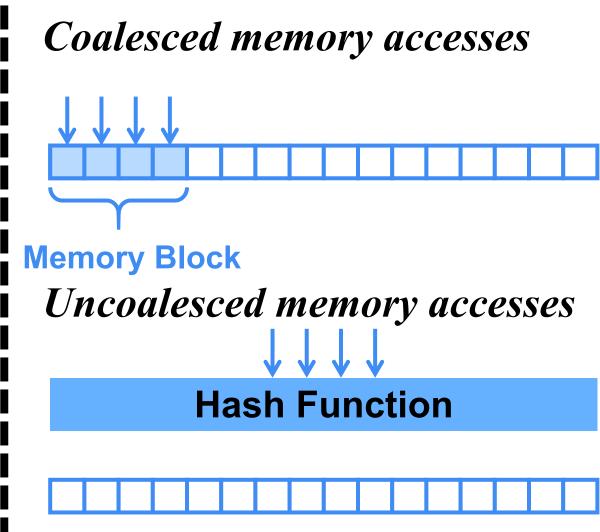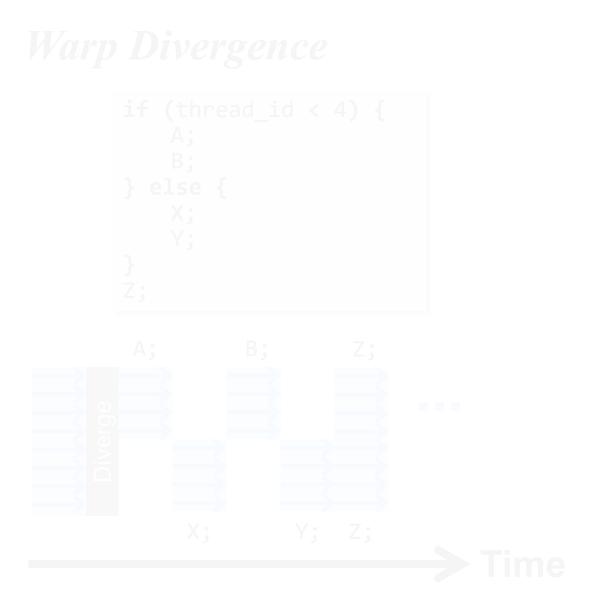
*Coalesced memory accesses*

**Memory Block**

*Uncoalesced memory accesses*

Time

# Challenge 1: Agnostic to GPU Execution Manner



*Warp Divergence*

```
if (thread_id < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;
```

*Coalesced memory accesses*

**Memory Block**

*Uncoalesced memory accesses*

Time

# Challenge 1: Agnostic to GPU Execution Manner

*Warp Divergence*

```
if (thread_id < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;
```

A;          B;          Z;

Diverge

X;          Y;   Z;

Time

*Coalesced memory accesses*

**Memory Block**

*Uncoalesced memory accesses*

**Hash Function**

8

# Challenge 1: Agnostic to GPU Execution Manner



*Warp Divergence*

```
if (thread_id < 4) {
    A;
    B;
} else {
    X;
    Y;
}
Z;
```

*Coalesced memory accesses*

**Memory Block**

*Uncoalesced memory accesses*

**Hash Function**

Time

# Challenge 1: Agnostic to GPU Execution Manner

Severe *warp divergence* and *uncoalesced memory accesses* lead to **Performance Degradation**

# Challenge 2: Ensure Crash Consistency

# Challenge 2: Ensure Crash Consistency

**Without Consistency Guarantee**

# Challenge 2: Ensure Crash Consistency

**Without Consistency Guarantee**



CPU

Memory Bus (8-Byte)

PM

# Challenge 2: **Ensure Crash Consistency**

**Without Consistency Guarantee**

# Challenge 2: Ensure Crash Consistency

**Without Consistency Guarantee**

# Challenge 2: Ensure Crash Consistency

**Without Consistency Guarantee**



CPU

32-Byte Data

Memory Bus (8-Byte)

PM

Partial Update

# Challenge 2: Ensure Crash Consistency

**Without Consistency Guarantee**

# Challenge 2: Ensure Crash Consistency

**Without Consistency Guarantee**



*Data Inconsistency!*

# Challenge 2: Ensure Crash Consistency

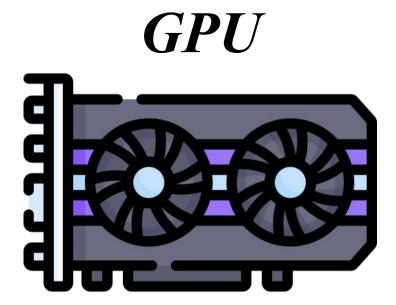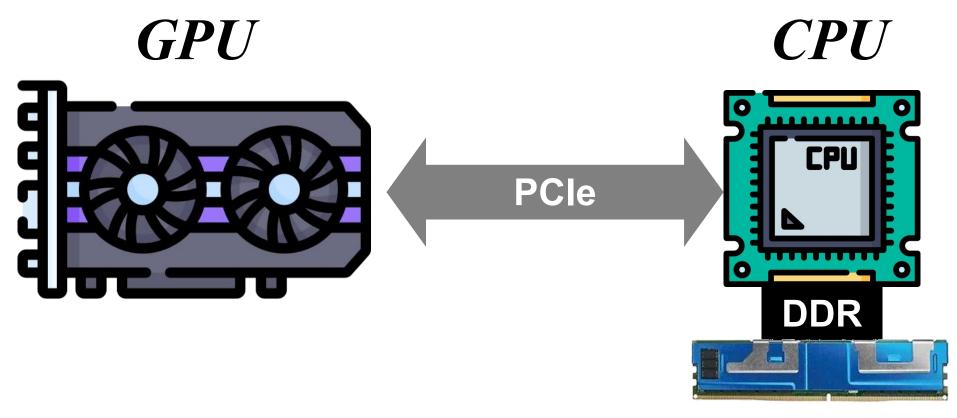**Without Consistency Guarantee**     **With Consistency Guarantee**



**CPU**

**32-Byte Data**

**Memory Bus (8-Byte)**

**PM**

**Partial Update**

**Crash**

*Data Inconsistency!*

# Challenge 2: **Ensure Crash Consistency**

**Without Consistency Guarantee**   **With Consistency Guarantee**



**CPU**

**32-Byte Data**

**Memory Bus (8-Byte)**

**Logging**

**PM**

**Partial Update**

**Crash**

*Data Inconsistency!*

9

# Challenge 2: Ensure Crash Consistency

**Without Consistency Guarantee**

**With Consistency Guarantee**

**CPU**

**32-Byte Data**

**CPU**

**Logging**

**Memory Bus (8-Byte)**

**Memory Bus (8-Byte)**
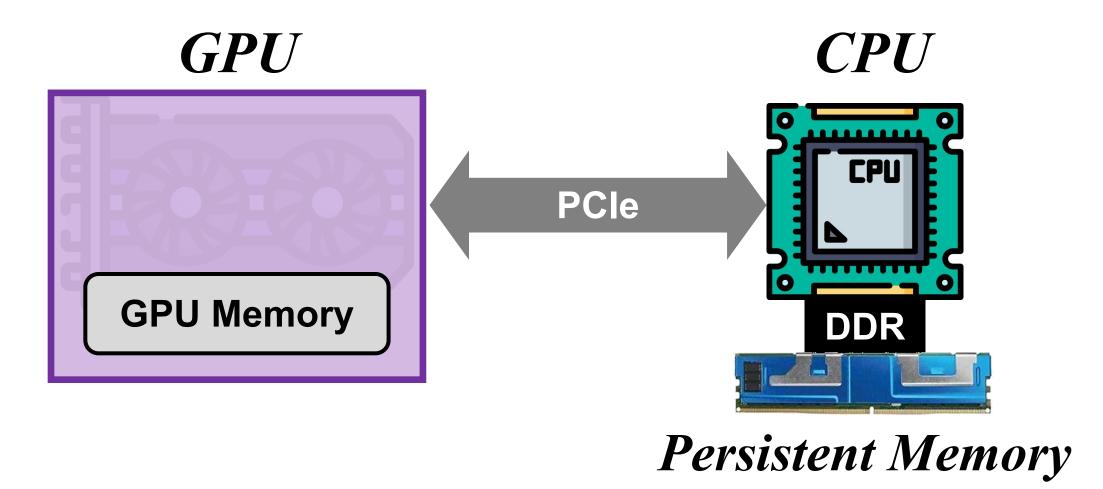
**PM**

**Partial Update**

**Crash**

**PM**

*Data Inconsistency!*

# Challenge 2: Ensure Crash Consistency

**Without Consistency Guarantee**

**With Consistency Guarantee**

CPU

32-Byte Data

Memory Bus (8-Byte)

PM

Partial Update

Crash

*Data Inconsistency!*

Logging

CPU

32-Byte Data

Memory Bus (8-Byte)

PM

# Challenge 2: Ensure Crash Consistency

**Without Consistency Guarantee**

**With Consistency Guarantee**



**Logging**

*Data Inconsistency!*

# Challenge 2: Ensure Crash Consistency

**Without Consistency Guarantee**

**With Consistency Guarantee**



**Logging**

**Data Inconsistency!**

# Challenge 2: Ensure Crash Consistency

**Without Consistency Guarantee**

**With Consistency Guarantee**

**CPU**

**32-Byte Data**

**CPU**

**32-Byte Data**

**Memory Bus (8-Byte)**

**Logging**

**Memory Bus (8-Byte)**

**PM** **Partial Update**

**Crash**

**PM**

**Logs**

**Crash**

*Data Inconsistency!*

# Challenge 2: Ensure Crash Consistency

**Without Consistency Guarantee**

**With Consistency Guarantee**



*Data Inconsistency!* — 2× writes — *Recover*

# Challenge 2: Ensure Crash Consistency

Crash consistency guarantee introduces **High Overhead**

# Challenge 3: Huge Bandwidth Gap

# Challenge 3: Huge Bandwidth Gap

*GPU*

# Challenge 3: Huge Bandwidth Gap

*GPU*

*CPU*

PCIe

DDR

*Persistent Memory*

# Challenge 3: Huge Bandwidth Gap

*GPU*

*CPU*

PCIe

GPU Memory

DDR

*Persistent Memory*

# Challenge 3: Huge Bandwidth Gap



*GPU*

*CPU*

*Threads*

PCIe

GPU Memory

DDR

*Persistent Memory*

# Challenge 3: Huge Bandwidth Gap



*GPU*

*Threads*

GPU Memory

PCIe

*CPU*

CPU

DDR

*Persistent Memory*

# Challenge 3: Huge Bandwidth Gap

*GPU*

*CPU*

*Threads*

...

GPU Memory

PCIe

CPU

DDR

e.g., **900 GB/s** for NVIDIA V100

*Persistent Memory*

# Challenge 3: Huge Bandwidth Gap



**GPU**

Threads

GPU Memory

PCIe

**CPU**

only *tens of GB/s*

DDR

*Persistent Memory*

e.g., *900 GB/s* for NVIDIA V100

# Challenge 3: Huge Bandwidth Gap



*GPU*

*CPU*

Threads

GPU Memory

PCIe

Huge Gap

only *tens of GB/s*

e.g., *900 GB/s* for NVIDIA V100

*Persistent Memory*

# Challenge 3: Huge Bandwidth Gap

Huge bandwidth gap between PM and GPU **Limits the Utilization** of GPU's high parallelism
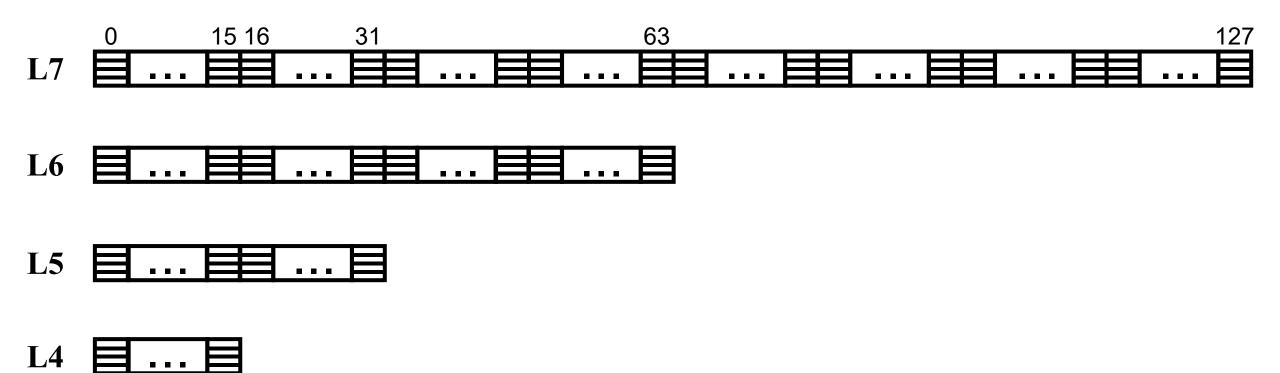
# Our Solution: GPHash

# Our Solution: GPHash

➢ **GPHash:** *An Efficient Hash Index for GPM Systems*

# Our Solution: **GPHash**

➤ **GPHash:** *An Efficient Hash Index for GPM Systems*

# Our Solution: GPHash

➢ **GPHash:** *An Efficient Hash Index for GPM Systems*

# Our Solution: GPHash

> **GPHash:** *An Efficient Hash Index for GPM Systems*



GPU-conscious and PM-friendly hash table

# Our Solution: **GPHash**

➢ **GPHash:** *An Efficient Hash Index for GPM Systems*

# Our Solution: GPHash

> **GPHash:** *An Efficient Hash Index for GPM Systems*



GPU

*Threads*

Lock-free and log-free operations

**Direct Access**

*Buckets*

**GPU Memory**

**Frozen-based bucket cache**

**GPU-conscious and PM-friendly hash table**

11

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table



0   15 16   31   63   127

Slot 0
Slot 1
Slot 2
Slot 3

- *Slot Associativity*

12

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table



- *Inter-level Sharing*

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table



**Warp** §Thread 0 §Thread 1 ⋯ §Thread 15 §Thread 16 §Thread 17 ⋯ §Thread 31

*Activated*

$Hash_1(Key_{31})$

$Hash_2(Key_{31})$

Candidate buckets in one-shot warp access

- *Multiple Hash Locations*
- *One-Shot Warp Access*

# GPU-Conscious and PM-Friendly Hash Table



**Warp** · Thread 0 · Thread 1 ⋯ Thread 15 · Thread 16 · Thread 17 ⋯ Thread 31

$Hash_1(Key_{31})$

$Hash_2(Key_{31})$

*Activated*

L7  0  15 16  31  63  127

L6

**Candidate buckets in one-shot warp access**

L5

• *Multiple Hash Locations*

L4

• *One-Shot Warp Access*

**32x → 1x**

# GPU-Conscious and PM-Friendly Hash Table



Warp  §Thread 0  §Thread 1  ••• §Thread 15  §Thread 16  §Thread 17  ••• §Thread 31

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table

*Pointer-based key placement*

# GPU-Conscious and PM-Friendly Hash Table

*Pointer-based key placement*

**Bucket** | FPs & States | Key Ptr$_0$ | Key Ptr$_1$ | Key Ptr$_2$ | Key Ptr$_3$ | Value Ptrs

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table

**Pointer-based key placement**

Thread 0    Thread 1    Thread 2    Thread 3

Bucket:  FPs & States | Key Ptr$_0$ | Key Ptr$_1$ | Key Ptr$_2$ | Key Ptr$_3$ | Value Ptrs

PM Space:

*Scattered addresses*

**In-place key placement**

Bucket:  FPs & States | Key$_0$ | Key$_1$ | Key$_2$ | Key$_3$ | Value Ptrs

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table

# GPU-Conscious and PM-Friendly Hash Table

✓ **GPU-friendly:** *minimize warp divergence and uncoalesced memory accesses*

✓ **Write-optimized:** *each insertion only involves a constant number of buckets without any data movement*

✓ **Memory-efficient:** *achieve a high load factor that is up to 92%*

# Lock-Free and Log-Free Insertion

# Lock-Free and Log-Free Insertion

**Warp**  Thread 0    Thread 1    ...    Thread 31

# Lock-Free and Log-Free Insertion

# Lock-Free and Log-Free Insertion

# Lock-Free and Log-Free Insertion

# Lock-Free and Log-Free Insertion

# Lock-Free and Log-Free Insertion



**Warp**

**Buckets**

❶ Check if the key exists

❷ Find the empty slot belonging to the less-loaded bucket

**Target Slot**

| FP & State: EMPTY_KEY | Key | Value Ptr |

EMPTY: 0xFFFFFFFFFFFFFFFF

❸ Set the State to INSERTING using AtomicCAS operation

# Lock-Free and Log-Free Insertion



**Warp**

Thread 0   Thread 1   ...   Thread 31

❶ Check if the key exists

**Buckets**

...   ...   ...

❷ Find the empty slot belonging to the less-loaded bucket

**Target Slot**

| FP & State: INSERT | Key | Value Ptr |

INSERT: 0xFFFFFFFFFFFFFFF<u>E</u>

**CAS Succeeds**   ❸ Set the State to INSERTING using AtomicCAS operation

# Lock-Free and Log-Free Insertion



**Warp**: Thread 0, Thread 1, ..., Thread 31

❶ Check if the key exists

**Buckets**

❷ Find the empty slot belonging to the less-loaded bucket

**CAS Fails**

**Target Slot**: FP & State: INSERT | Key | Value Ptr

INSERT: 0xFFFFFFFFFFFFFF**E**

**CAS Succeeds**  ❸ Set the State to INSERTING using AtomicCAS operation

# Lock-Free and Log-Free Insertion



Warp

Thread 0   Thread 1   ...   Thread 31

❶ Check if the key exists

Buckets

...   ...   ...

❷ Find the empty slot belonging to the less-loaded bucket

Target Slot   FP & State: INSERT   Key   Value Ptr

INSERT: 0xFFFFFFFFFFFFFFFE

# Lock-Free and Log-Free Insertion

# Lock-Free and Log-Free Insertion



**Warp**

❶ Check if the key exists

**Buckets**

❷ Find the empty slot belonging to the less-loaded bucket

In-place writing

**Target Slot**

| FP & State: INSERT | Key | | Value Ptr |

INSERT: 0xFFFFFFFFFFFFFFF**E**

❹ Write the key and value ptr

# Lock-Free and Log-Free Insertion



Warp

Thread 0    Thread 1    ...    Thread 31

❶ Check if the key exists

Buckets

...    ...    ...

❷ Find the empty slot belonging to the less-loaded bucket

In-place writing

Target Slot | FP & State: INSERT | Key | | Value Ptr

INSERT: 0xFFFFFFFFFFFFFF**E**

❺ Set the FP to the hash value of the key

# Lock-Free and Log-Free Insertion

**Warp**



❶ Check if the key exists

**Buckets**

❷ Find the empty slot belonging to the less-loaded bucket

**In-place writing**

**Target Slot** | FP & State: INSERT | Key | | Value Ptr

INSERT: 0xFFFFFFFFFFFFFFF**E**

❺ Set the FP to the hash value of the key

**Target Slot** | FP & State: Hash Value | Key | Value Ptr

# Lock-Free and Log-Free Insertion

Warp



❶ Check if the key exists

Buckets

❷ Find the empty slot belonging to the less-loaded bucket

In-place writing

Target Slot | FP & State: INSERT | Key **Crash** | Value Ptr

INSERT: 0xFFFFFFFFFFFFFF**E**

❺ Set the FP to the hash value of the key

Target Slot | FP & State: Hash Value | Key | Value Ptr

# Lock-Free and Log-Free Insertion

**Warp**



Thread 0    Thread 1    ...    Thread 31

❶ Check if the key exists

**Buckets**

*Check slot state*

...belonging to the less-loaded bucket

**In-place writing**

**Target Slot** | FP & State: INSERT | Key **Crash** | Value Ptr

INSERT: 0xFFFFFFFFFFFFFFF**E**

❺ Set the FP to the hash value of the key

**Target Slot** | FP & State: Hash Value | Key | Value Ptr

# Lock-Free and Log-Free Insertion

**Warp**

Thread 0    Thread 1    ...    Thread 31

❶ Check if the key exists

**Buckets**

❷ Find the empty slot belonging to the less-loaded bucket

**Recover**

**Target Slot** | FP & State: EMPTY_KEY | Key | Value Ptr

EMPTY: 0xFFFFFFFFFFFFFFF**F**

❺ Set the FP to the hash value of the key

**Target Slot** | FP & State: Hash Value | Key | Value Ptr

# Frozen-Based Bucket Cache

# Frozen-Based Bucket Cache

**GPU**

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**PM**

# Frozen-Based Bucket Cache

**GPU**

**PM**

GPHash Buckets

14

# Frozen-Based Bucket Cache

**GPU**

**Cache**

**PM**

**GPHash Buckets**

14

# Frozen-Based Bucket Cache



GPU

Cache

PM

GPHash Buckets

Hot Buckets

# Frozen-Based Bucket Cache

**GPU**

- - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - - -

**PM**

# Frozen-Based Bucket Cache

**GPU**

*BktCache*

**PM**

*GPHash*

# Frozen-Based Bucket Cache

# Frozen-Based Bucket Cache

# Frozen-Based Bucket Cache

Thread 0    Thread 1    ...    Thread 31

**GPU**

*BktCache*

...

**PM**

*GPHash*

...    ...

...

# Frozen-Based Bucket Cache



14

# Frozen-Based Bucket Cache

# Frozen-Based Bucket Cache



14

# Frozen-Based Bucket Cache



14

# Frozen-Based Bucket Cache



**Cache hit**

**Cache miss**

Thread 0  Thread 1  ...  Thread 31

*BktCache*

GPU

*GPHash*

PM

*Concurrent Loading*

❶ Invalidate the old cached bucket

14

# Frozen-Based Bucket Cache



Thread 0   Thread 1   ...   Thread 31

➡ **Cache hit**

➡ **Cache miss**

*BktCache*

**GPU**

*GPHash*

**PM**

*Concurrent Loading*

❶ Invalidate the old cached bucket

14

# Frozen-Based Bucket Cache



14

# Frozen-Based Bucket Cache



§ Thread 0    § Thread 1    ...    § Thread 31    --→ **Cache hit**

--→ **Cache miss**

*BktCache*

**GPU**

*GPHash*

**PM**

*Concurrent Loading*

❶ Invalidate the old cached bucket

❷ Fetch the new bucket

❸ Validate the new cached bucket

14

# More Details

# More Details

- ➢ Warp-cooperative Exec-
  ution Manner
- ➢ More Index Operations
- ➢ Bucket Caching Granularity
- ➢ ...

# More Details

- ➤ Warp-cooperative Exec-
  ution Manner

- ➤ More Index Operations

- ➤ Bucket Caching Granularity

- ➤ ...

# More Details

➢ Warp-cooperative Execution Manner

➢ More Index Operations

➢ Bucket Caching Granularity

➢ ...

READ PAPER

# Evaluation

# Evaluation

> **Platform**

- 1 V100 GPU + 768 GB Intel Optane DC PM (6 × 128 GB)

# Evaluation

➤ **Platform**

- 1 V100 GPU + 768 GB Intel Optane DC PM (6 × 128 GB)

➤ **Comparisons**

- CPU-assisted approaches[1]: *Clevel[ATC'20]* , *Dash[VLDB'20], and SEPH[OSDI'23]*

- GPM-enabled approaches: *Clevel-GPM and SlabHash[IPDPS'18]-GPM*

[1] For fair evaluation, we use PM to store the data and leverage PM hash indexes to manage the data in PM
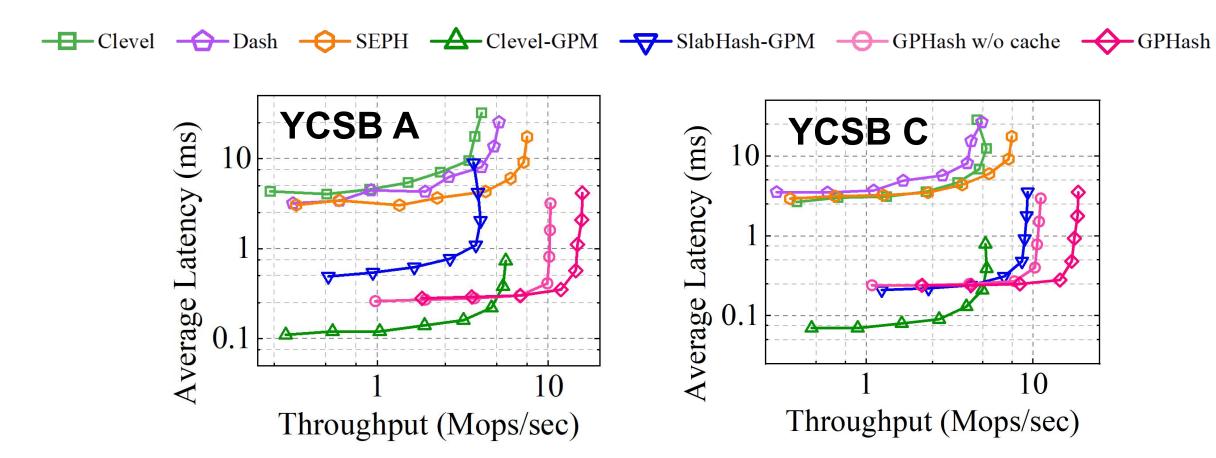
# Evaluation

➢ **Platform**
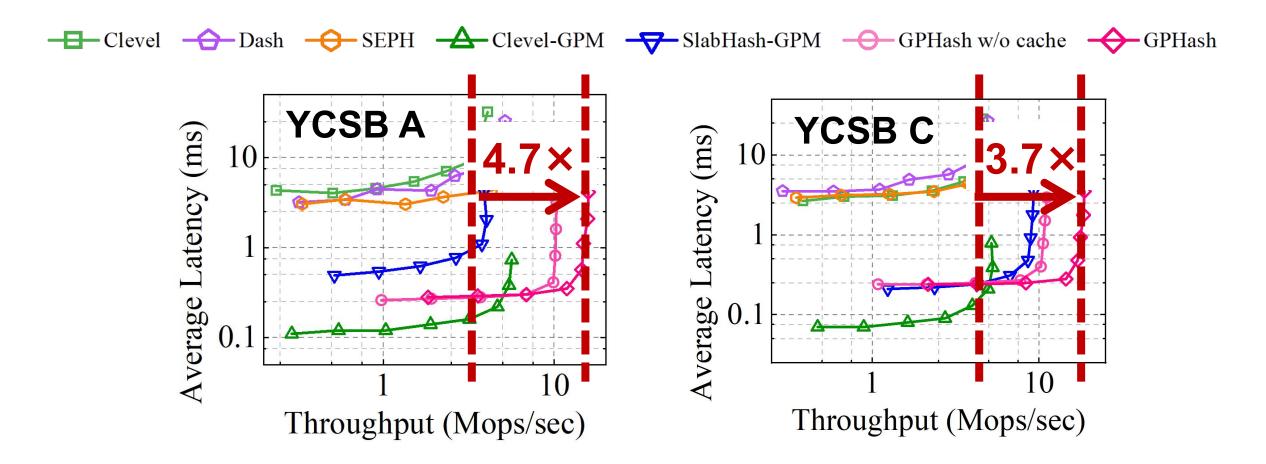
- 1 V100 GPU + 768 GB Intel Optane DC PM (6 × 128 GB)

➢ **Comparisons**

- CPU-assisted approaches[1]: *Clevel[ATC'20]* , *Dash[VLDB'20], and SEPH[OSDI'23]*

- GPM-enabled approaches: *Clevel-GPM and SlabHash[IPDPS'18]-GPM*

➢ **Workloads**

- YCSB workloads: *8-byte and 32-byte keys, 128-byte values*

- Real-world workloads: *DLRM and PageRank*

[1] For fair evaluation, we use PM to store the data and leverage PM hash indexes to manage the data in PM

# End-to-End Performance

# End-to-End Performance

# End-to-End Performance

# End-to-End Performance



GPHash improves the throughput by **1.9~6.3×**

# Latency Breakdown

# Latency Breakdown



CPU-assisted approaches suffer from *high transmission cost*

# Latency Breakdown



Naive GPM-enabled approaches suffer from *severe warp divergence* and *high-overhead consistency guarantee*

# Latency Breakdown



GPHash fully leverages the *high parallelism of GPU* and provides a *low-overhead consistency guarantee*

# Conclusion

# Conclusion

➢ Inefficient GPU data management on large-scale data

- *High overhead for data transfer*
- *Extra CPU consumption*

# Conclusion

➢ Inefficient GPU data management on large-scale data

- *High overhead for data transfer*
- *Extra CPU consumption*

➢ **GPHash:** an efficient GPM-enabled hash index

- *GPU-conscious and PM-friendly hash table*
- *Lock-free and log-free operations*
- *Frozen-based bucket cache*

# Conclusion

➢ Inefficient GPU data management on large-scale data

- *High overhead for data transfer*
- *Extra CPU consumption*

➢ **GPHash:** an efficient GPM-enabled hash index

- *GPU-conscious and PM-friendly hash table*
- *Lock-free and log-free operations*
- *Frozen-based bucket cache*

# Conclusion

https://github.com/LighT-chenml/GPHash

chenml@hust.edu.cn

➢ Inefficient GPU data management on large-scale data

- *High overhead for data transfer*

- *Extra CPU consumption*

➢ **GPHash:** an efficient GPM-enabled hash index

- *GPU-conscious and PM-friendly hash table*

- *Lock-free and log-free operations*

- *Frozen-based bucket cache*

## *Thank you! Q&A*