# Common Counters:
# Compressed Encryption Counters for Secure GPU Memory

**Seonjin Na,** Sunho Lee, Yeonjae Kim, Jongse Park, Jaehyuk Huh
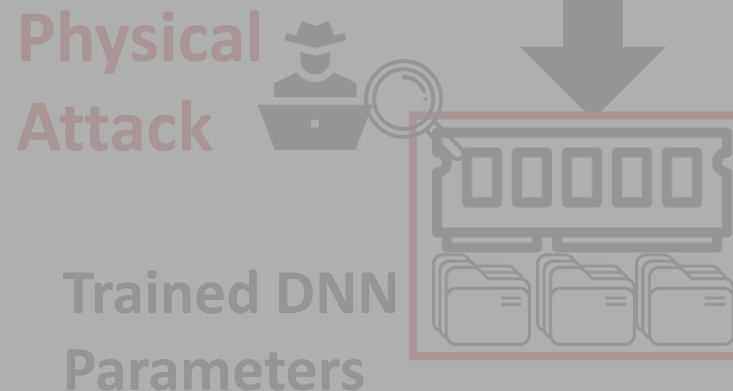
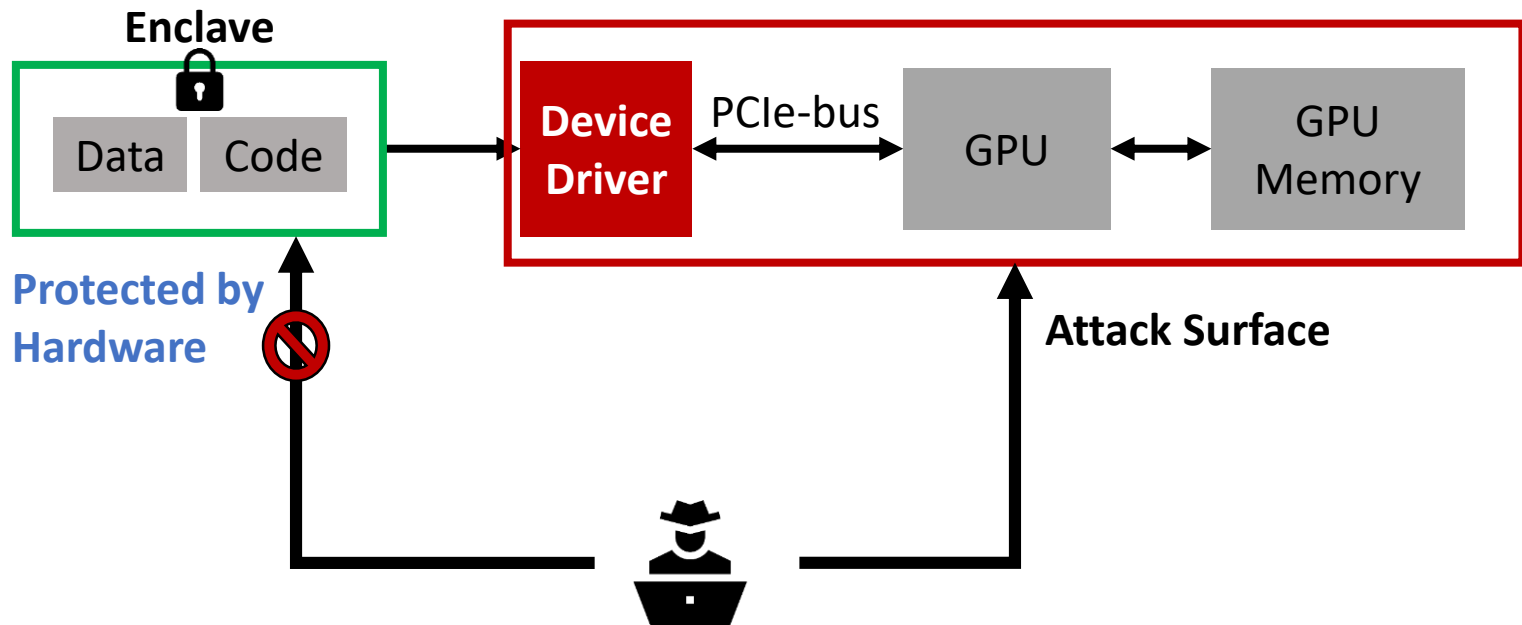**KAIST, School of Computing**

# Need for Secure GPU Computing



**Privileged SW Attack**

Untrusted Device Driver

Training Data & DNN Model

**We need to consider Secure GPU computing !**
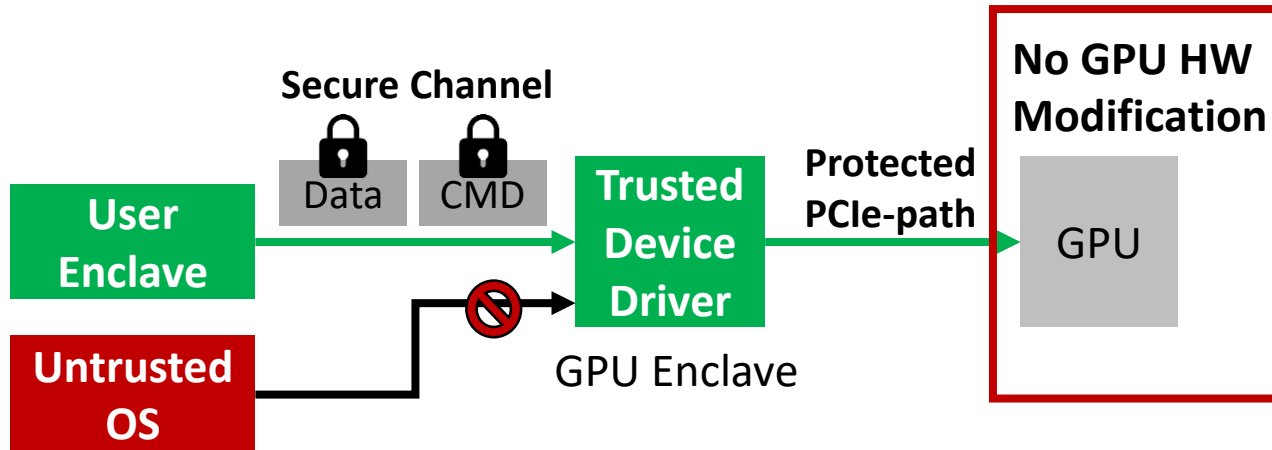
**Physical Attack**

Trained DNN Parameters

# Trusted GPU Computing

- Trusted Execution Environment (TEE)
  - Intel SGX, ARM TrustZone

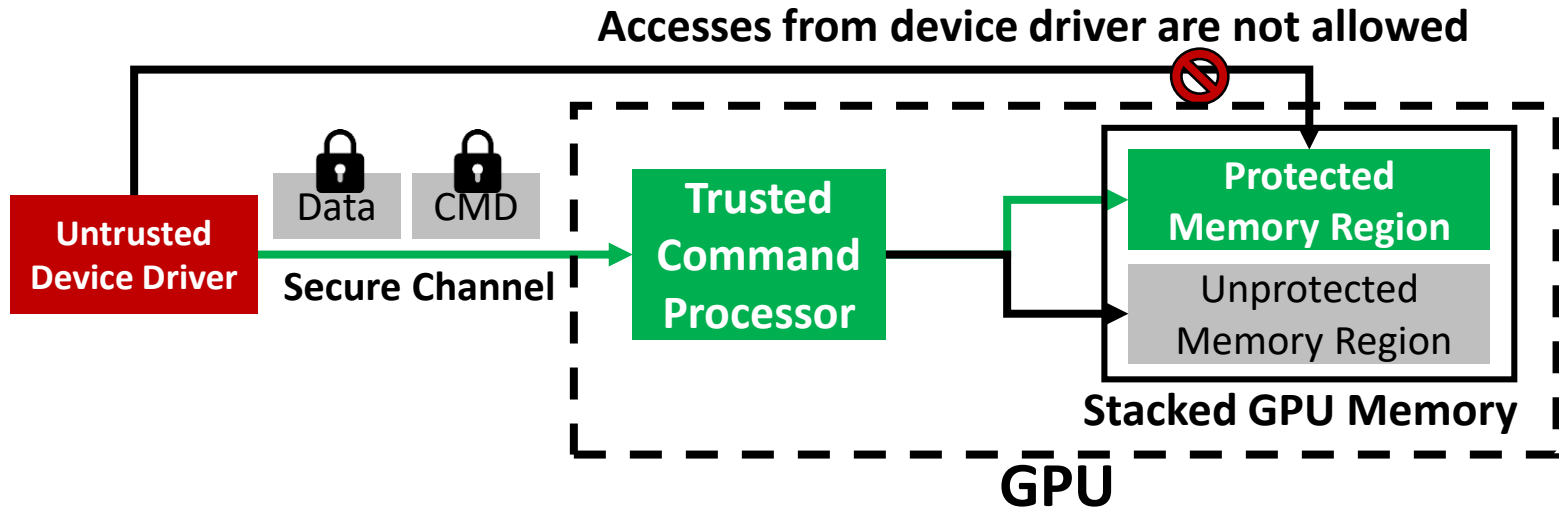- Existing TEEs **does not provide TEE on GPUs**

# Prior Work : HIX

- **HIX [ASPLOS '19]: Securing I/O Path from CPU to GPU**
  - All device I/O accesses to GPU are controlled by **trusted device driver**
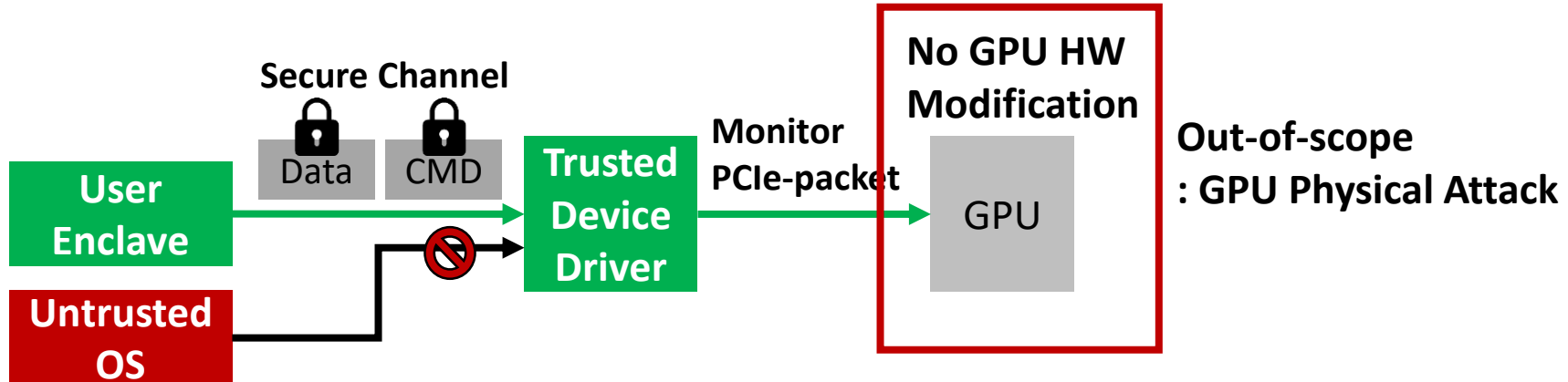
# Prior Work : Graviton

- **Graviton [OSDI '18]: Trusted GPU by changing GPU HW**
  - **Trusted Command Processor** handles critical GPU operations instead of driver



Accesses from device driver are not allowed

Untrusted Device Driver — Secure Channel → Trusted Command Processor

Data | CMD

Protected Memory Region
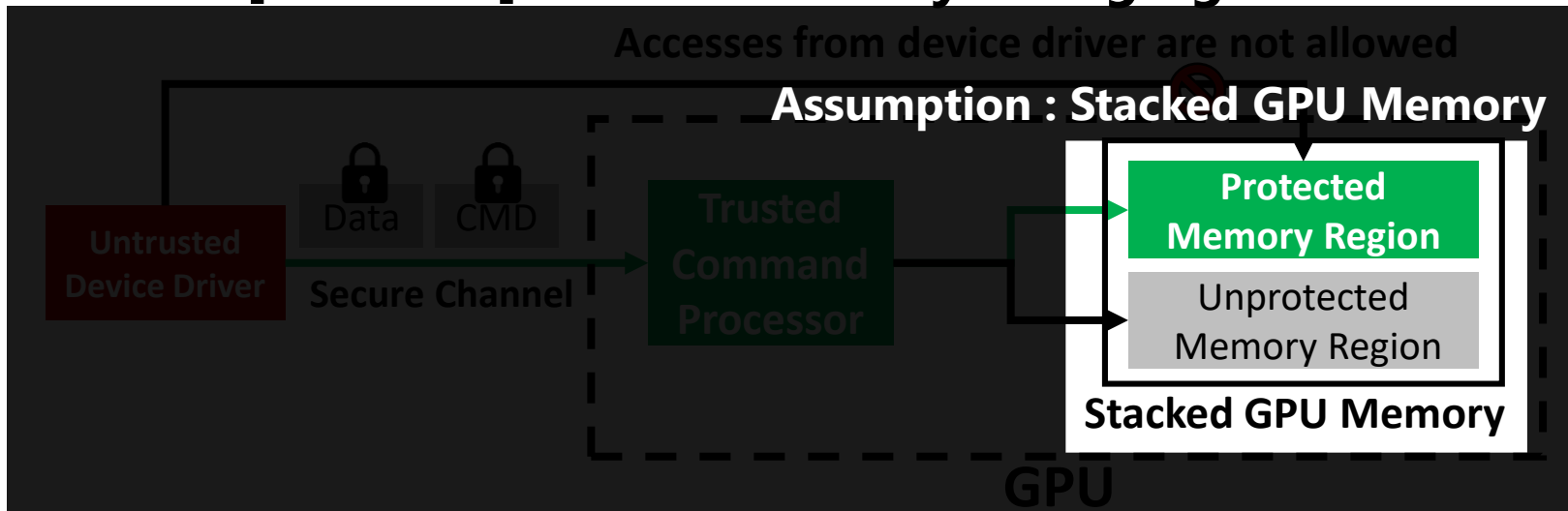Unprotected Memory Region
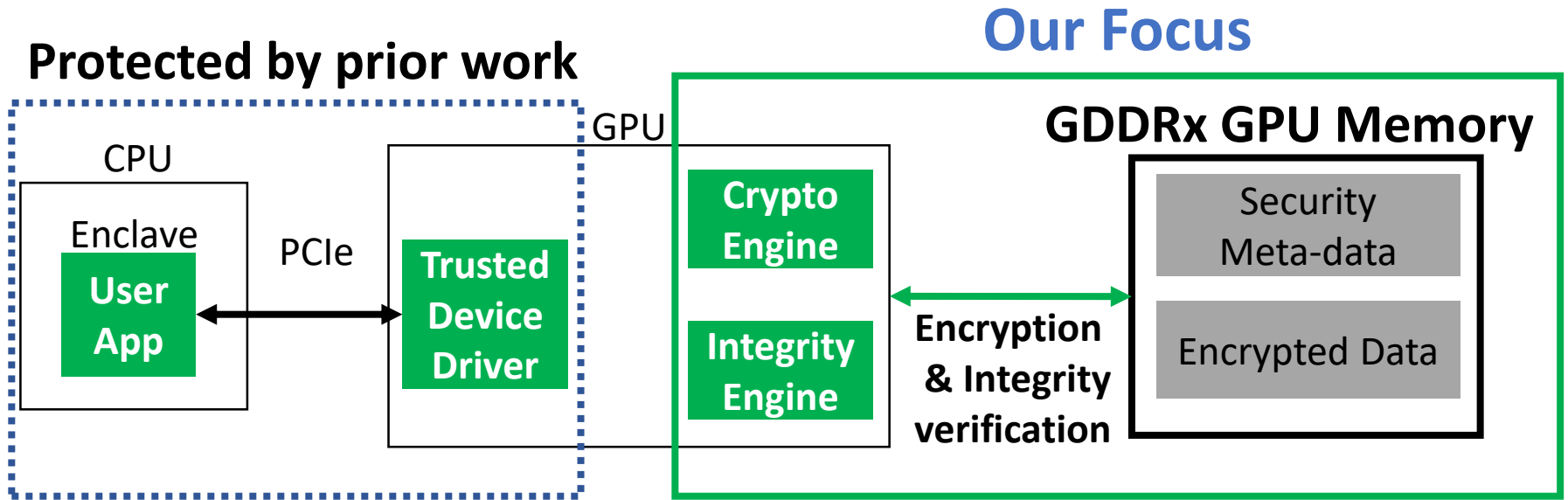Stacked GPU Memory

GPU

# Limitations of Prior Work

- ## HIX [ASPLOS '19]: Securing I/O Path from CPU to GPU



- ## Graviton [OSDI '18]: Trusted GPU by changing GPU HW

# Goal: Secure GPU Memory

**Protected by prior work**

**Our Focus**

**GPU**

CPU

Enclave

**User App**

PCIe

**Trusted Device Driver**

**Crypto Engine**

**Integrity Engine**

**Encryption & Integrity verification**

**GDDRx GPU Memory**

Security Meta-data

Encrypted Data

- **Main Contributions**
  - Provide secure GPU memory **with low performance overheads**
  - Exploit **unique memory update behavior** of common GPU applications
  - Reduce the average performance overhead to **2.9 %**
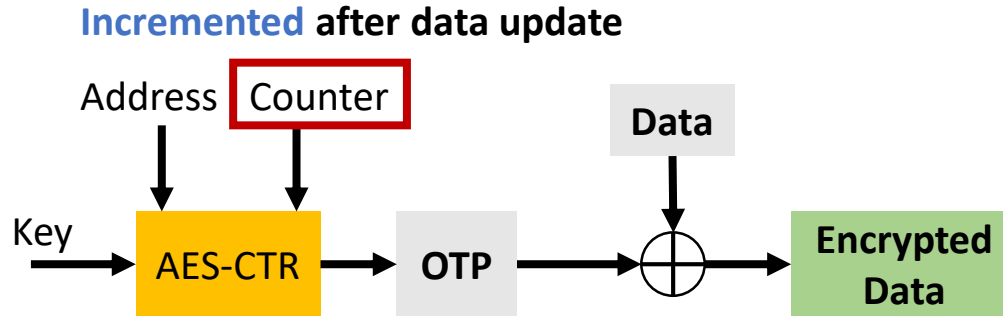
# Threat Model & Assumptions

- Threat Model
    - Attackers can fully control **operating system/hypervisor**
    - Attackers can **physically access the whole system**

- Trusted Computing Base (TCB)
    - GPU processor & GPU software running on the GPU
    - CPU chip & user application in an CPU Enclave

- Out of Scope
    - Denial of Service(DoS) attacks
    - Side-channel attacks

# Outline

- Introduction

- **Background & Motivation**

- Common Counter
  - Main Idea
  - Additional Metadata
  - Common Counter Mechanism

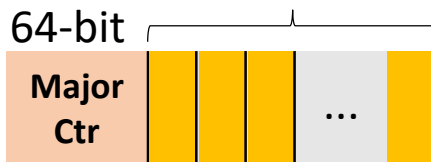- Evaluation

# Background : Securing Memory

- Memory Encryption
  - Counter mode encryption

**Incremented** after data update

Address | Counter

Key → AES-CTR → OTP → ⊕ → **Encrypted Data**

**Data**

  - Split Counter scheme

## Counter = Major | Minor

Minor Counters(7-bit for each)

64-bit

| **Major Ctr** | | | | ... | |

Cache block :128B
→ **128 minor counters**

# Background : Securing Memory

- Memory Integrity Verification
  - Message Authentication Code (MAC)

| Data | → | **Cryptographic Hash** | → | **Data MAC** |

  - Counter Integrity Tree

**Replay Attack**

{Data1,MAC1, Ctr1}

**Processor** ◄- - - **DRAM**

**Sniff**

{Data0, MAC0, Ctr0}
**(older copy)**

**On-chip secured**

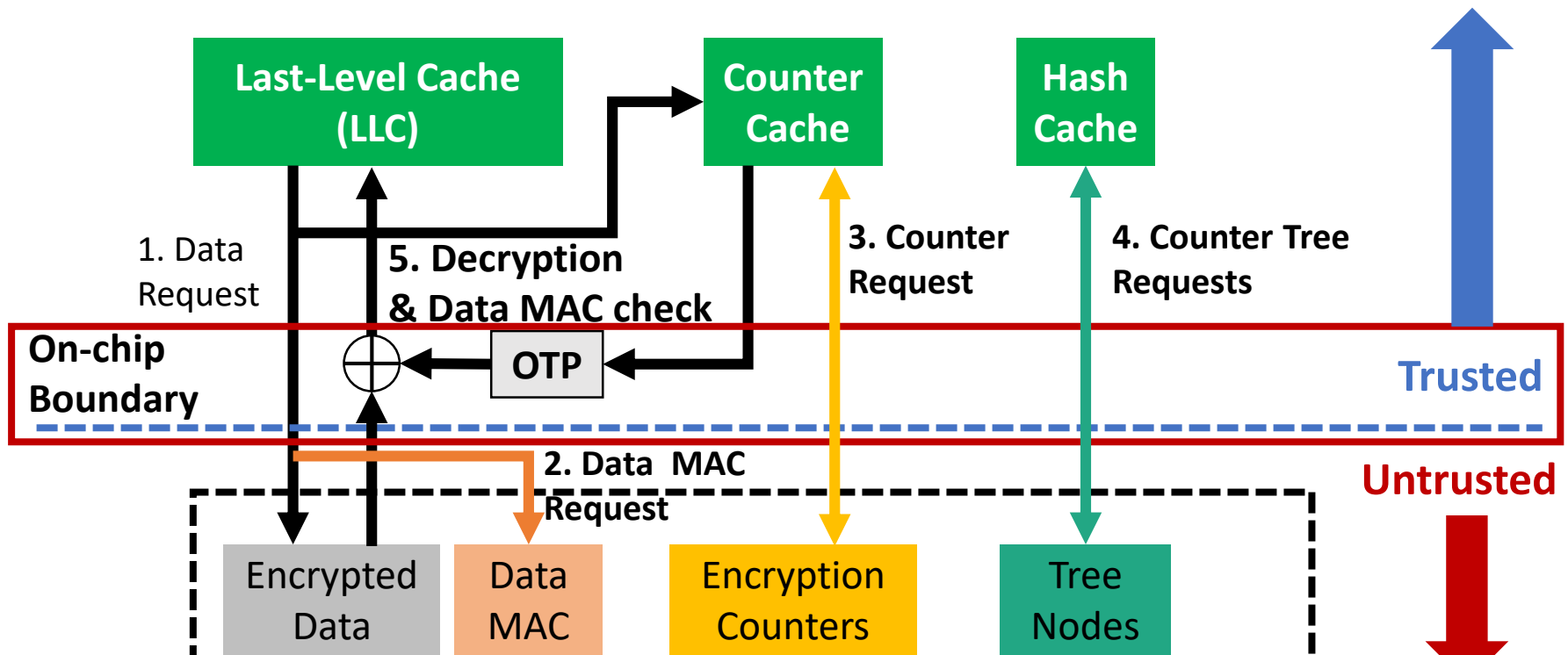**Counter Integrity Tree**

**Data Blocks**   **Counters**

Our baseline : **SC-128**
**128-ary** (Split Counter + Counter integrity tree)
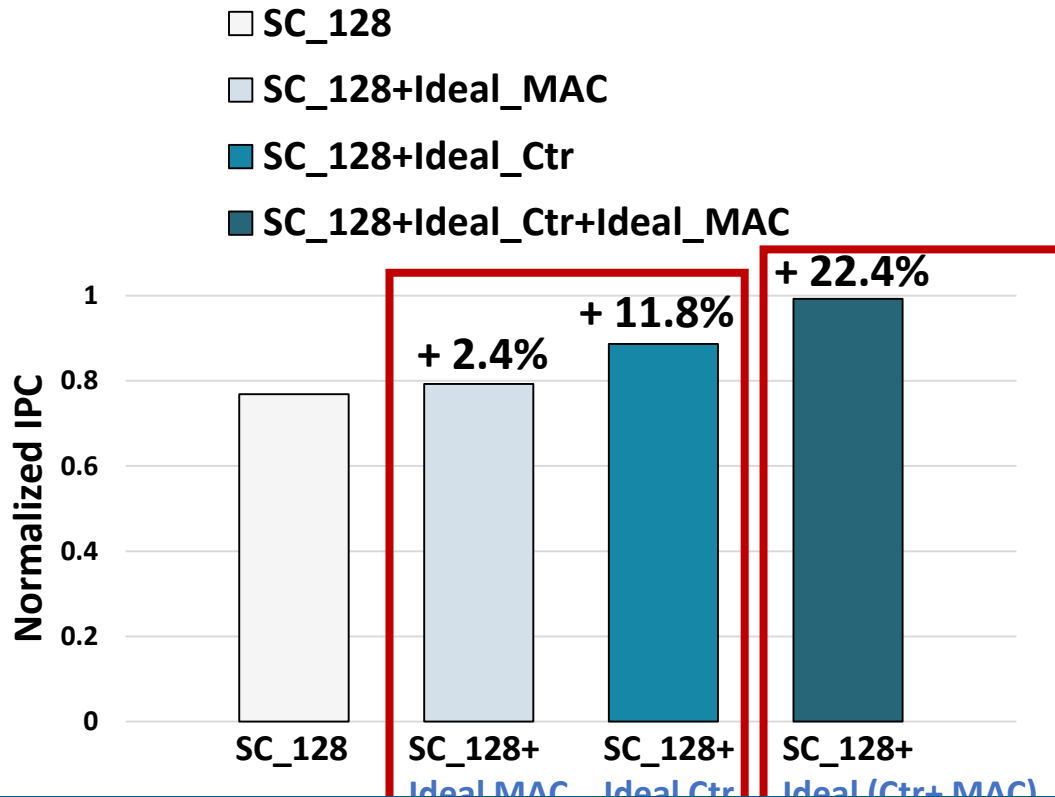
# Problem : Performance Overhead

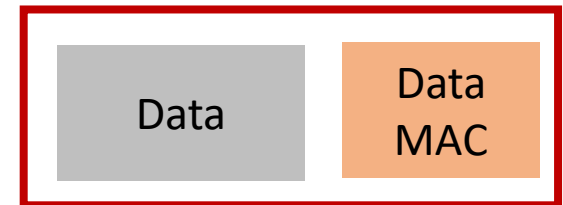- Secure memory require **additional meta-data requests**



Secure memory **adds decryption latency** and **increases memory bandwidth**

# Performance Breakdown Analysis

- GPU memory protection overhead result for GPU benchmark suites

□ **SC_128**
□ **SC_128+Ideal_MAC**
■ **SC_128+Ideal_Ctr**
■ **SC_128+Ideal_Ctr+Ideal_MAC**

**For data MAC overhead**

| Data | Data MAC |
|------|----------|

With ECC memory,
Data & MAC can be provided by
**1 memory access** using **Synergy[1]**

**Normalized IPC**

Chart values:
- SC_128
- SC_128+ Ideal MAC: **+ 2.4%**
- SC_128+ Ideal Ctr: **+ 11.8%**
- SC_128+ Ideal (Ctr+ MAC): **+ 22.4%**

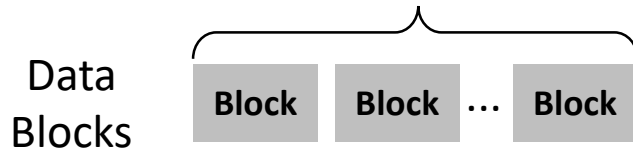**Counter mode encryption is one of the key bottlenecks**

[1] :SYNERGY: Rethinking Secure-Memory Design for Error-Correcting Memories, HPCA'18
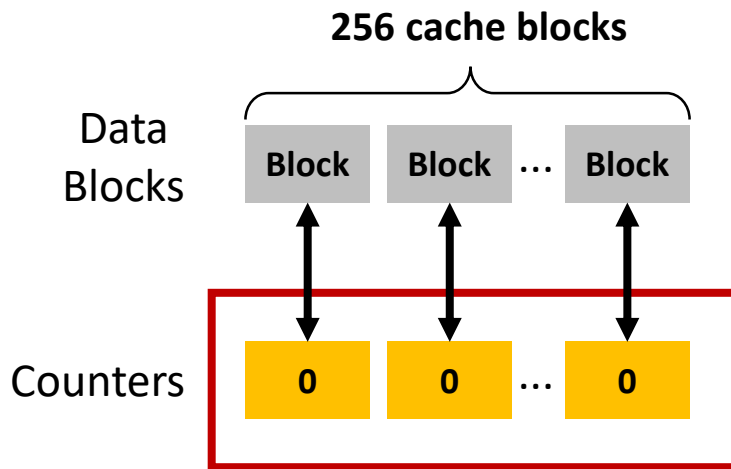
# Uniformly Updated Segments

- Memory segment: Contiguous memory region

  Example granularity: 32KB

  32KB/ 128B = **256 Data cache blocks**

  Data Blocks  | **Block** | **Block** | ... | **Block** |

- Uniformly updated segment: **Read-only** + **uniformly written**

**256 cache blocks**

Data Blocks | **Block** | **Block** | ... | **Block** |

Counters | 0 | 0 | ... | 0 |

**1. Read-only**

**256 cache blocks**

Data Blocks | **Block** | **Block** | ... | **Block** |

Counters | 2 | 2 | ... | 2 |

**2. Uniformly written**

14

# Observation : GPU SW Write Patterns

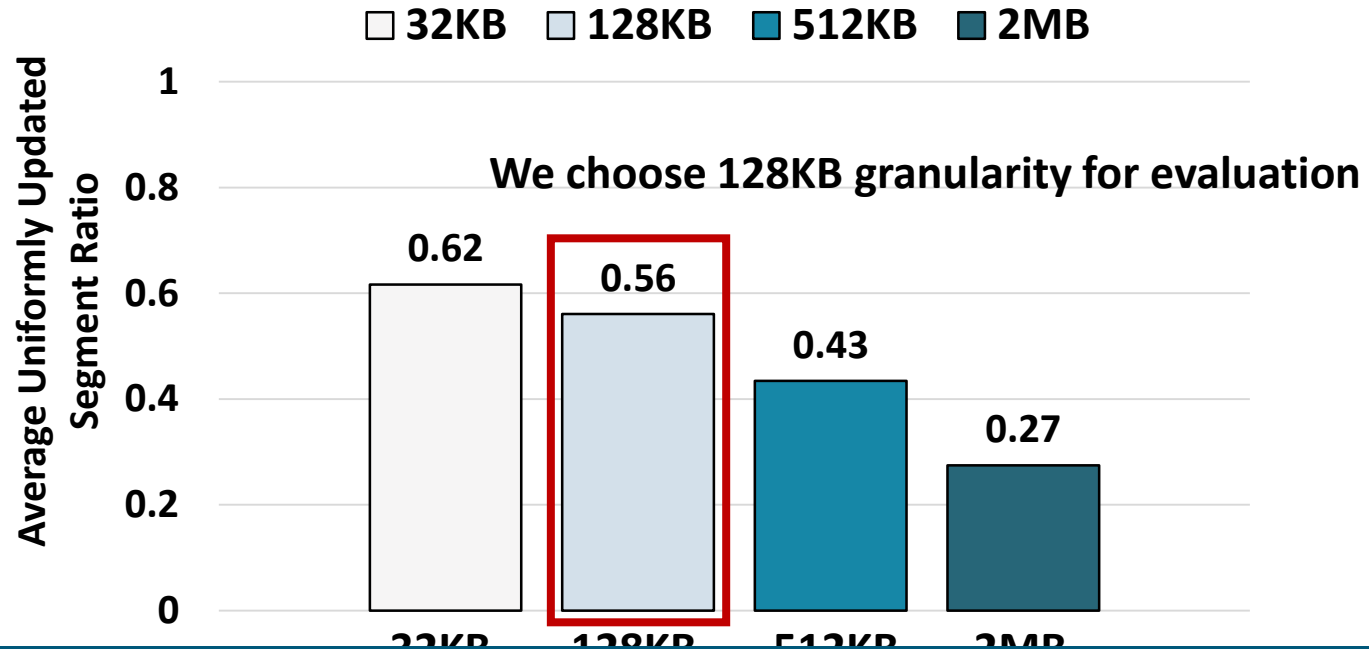- Analyze memory read/write behavior by using **NVBit [MICRO '19]**

**Result of GPU Benchmark Suite**

□ 32KB    □ 128KB    ■ 512KB    ■ 2MB

We choose 128KB granularity for evaluation

Average Uniformly Updated Segment Ratio

- 32KB: 0.62
- 128KB: 0.56
- 512KB: 0.43
- 2MB: 0.27

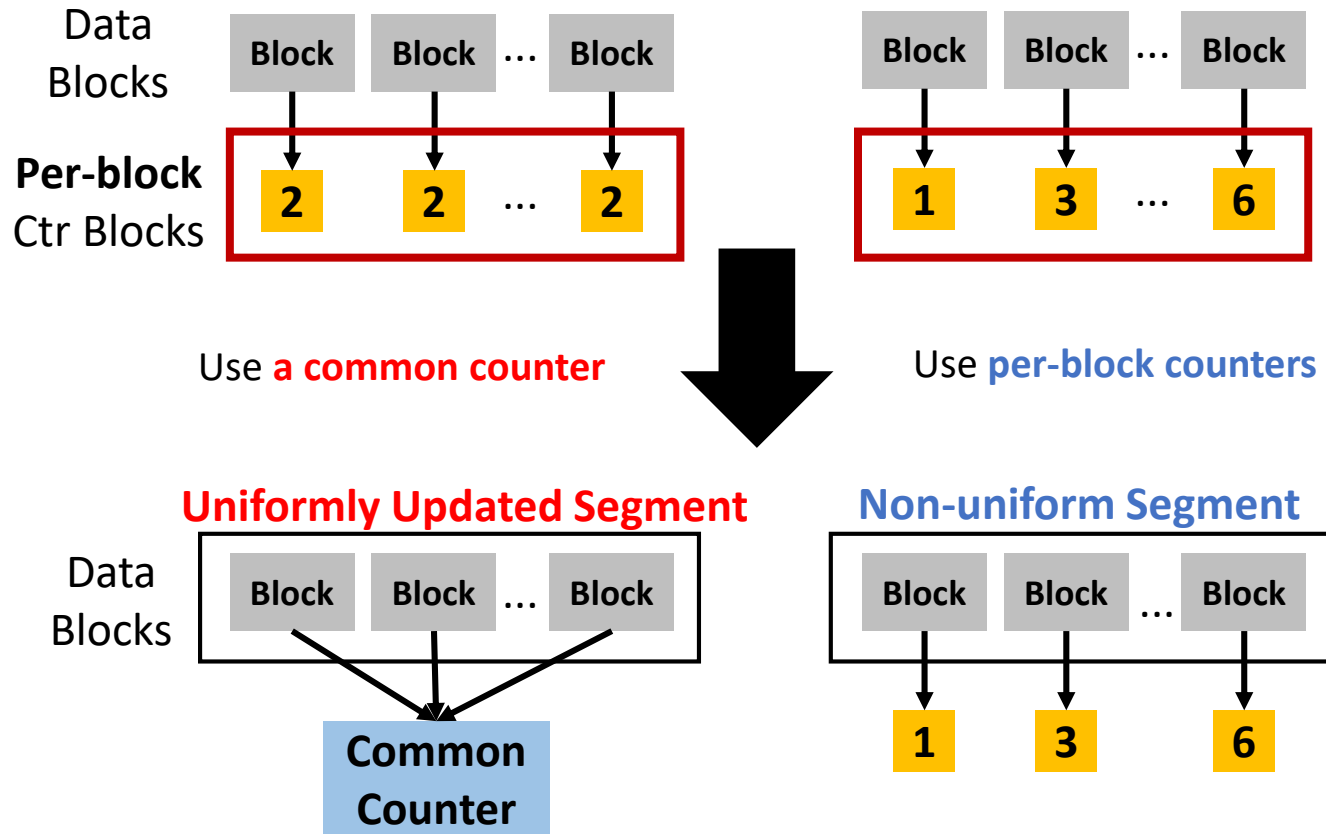Observation 1: **GPU programs tend to uniformly update memory**
Observation 2: **The number of distinct counter values is small**

# Outline

- Introduction

- Background & Motivation

- **Common Counter**
  - **Main Idea**
  - **Additional Metadata**
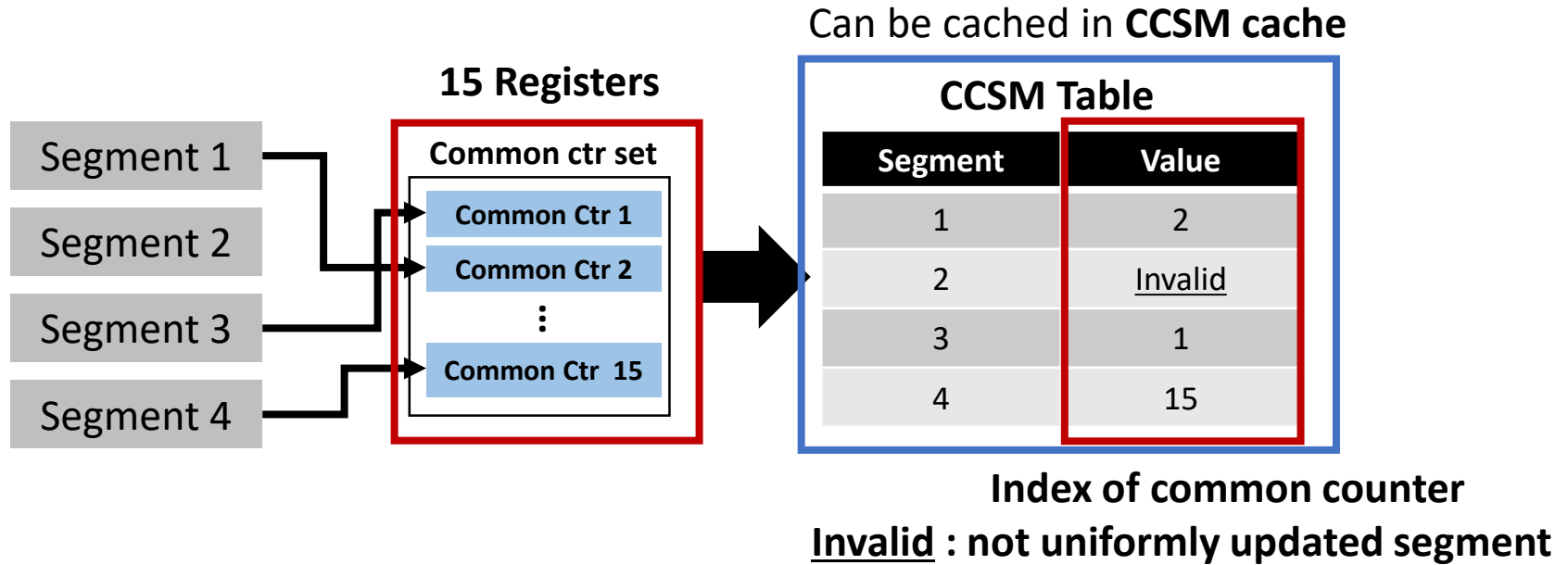  - **Common Counter Mechanism**

- Evaluation

# Common Counter : Main Idea

- Use **coarse-grained counters** for uniformly updated segments

Data Blocks

| Block | Block | ... | Block |

**Per-block** Ctr Blocks

| **2** | **2** | ... | **2** |

| Block | Block | ... | Block |

| **1** | **3** | ... | **6** |

Use **a common counter**

Use **per-block counters**

**Uniformly Updated Segment**

**Non-uniform Segment**

Data Blocks

| Block | Block | ... | Block |

**Common Counter**
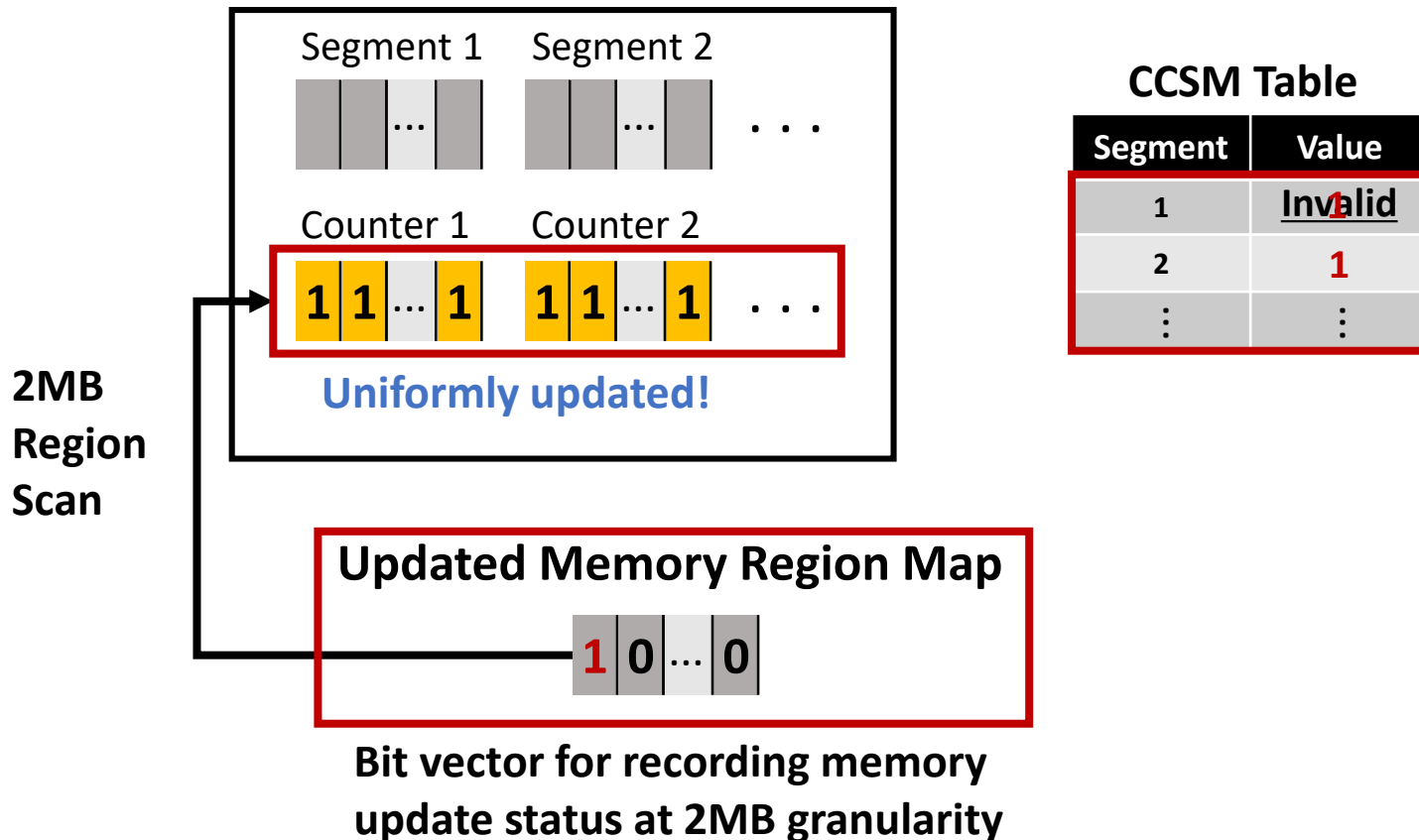
| Block | Block | ... | Block |

| **1** | **3** | **6** |

# Finding Uniformly Updated Segments

- Common Counter Status Map (CCSM)
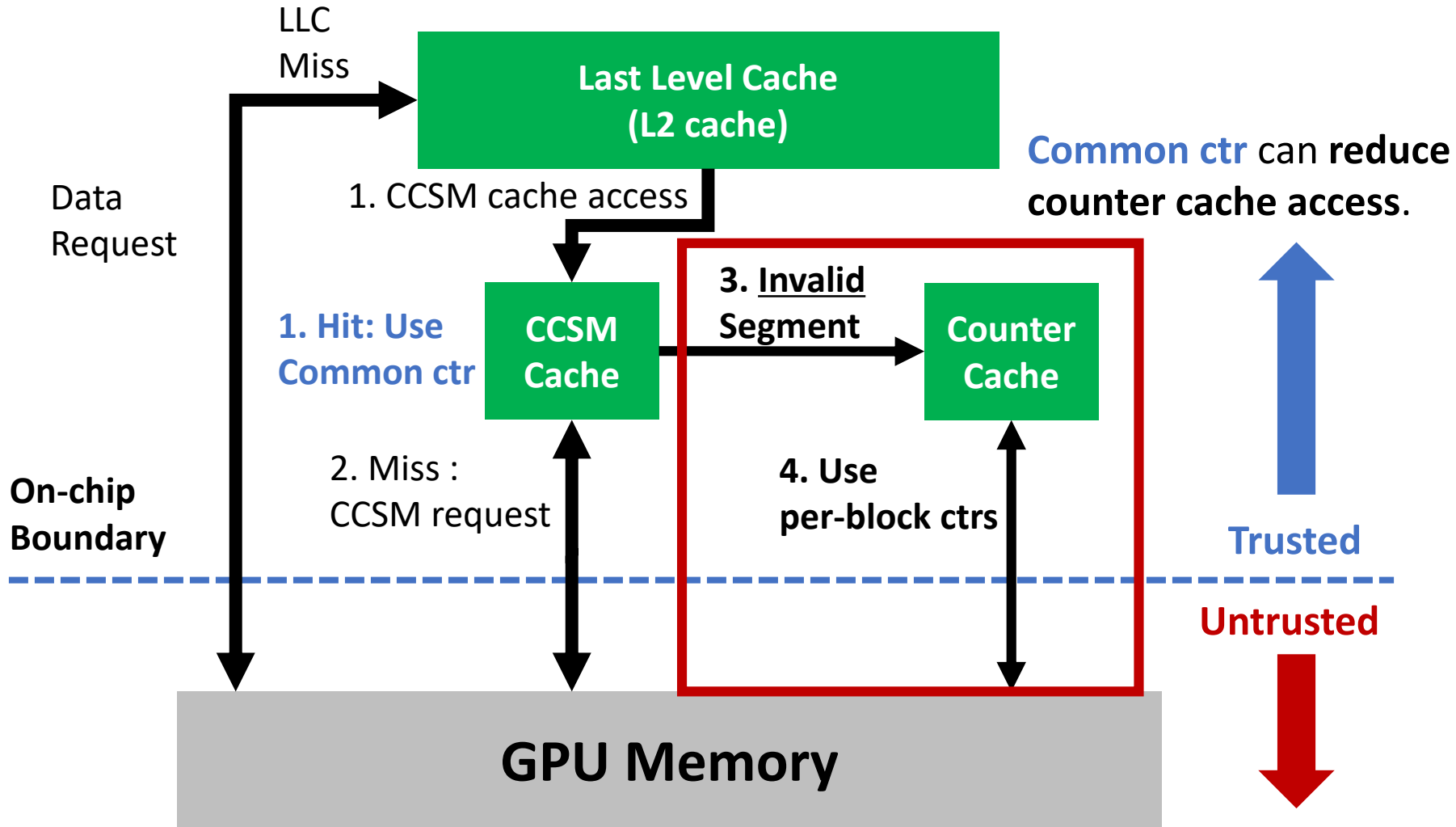  - Check whether a memory segment uses a common counter or not

Can be cached in **CCSM cache**

**15 Registers**

**CCSM Table**

| Segment 1 | | Common ctr set |
|---|---|---|

| Common Ctr 1 |
|---|
| Common Ctr 2 |
| ⋮ |
| Common Ctr 15 |

| Segment | Value |
|---|---|
| 1 | 2 |
| 2 | Invalid |
| 3 | 1 |
| 4 | 15 |

**Index of common counter**
**Invalid : not uniformly updated segment**

# Updating CCSM Table

- Initialized at application launch
- Scanning Procedure
  - When? After a kernel is completed



**CCSM Table**

| Segment | Value |
|---------|---------|
| 1 | Invalid 1 |
| 2 | 1 |
| ⋮ | ⋮ |

**2MB Region Scan**

Segment 1    Segment 2

Counter 1    Counter 2

**1 1 … 1    1 1 … 1 . . .**

**Uniformly updated!**

**Updated Memory Region Map**

**1 0 … 0**

**Bit vector for recording memory update status at 2MB granularity**

# LLC Miss Handling with Common Counters

Last Level Cache
(L2 cache)

LLC Miss

Data Request

1. CCSM cache access

**Common ctr** can **reduce counter cache access**.

**1. Hit: Use Common ctr**

CCSM Cache

**3. Invalid Segment**

Counter Cache

**2. Miss : CCSM request**

**4. Use per-block ctrs**

**On-chip Boundary**

**Trusted**

**Untrusted**

**GPU Memory**

# GPU Execution with Common Counter

# Outline

- Introduction

- Background & Motivation

- Common Counter
  - Main Idea
  - Additional Metadata
  - Common Counter Mechanism
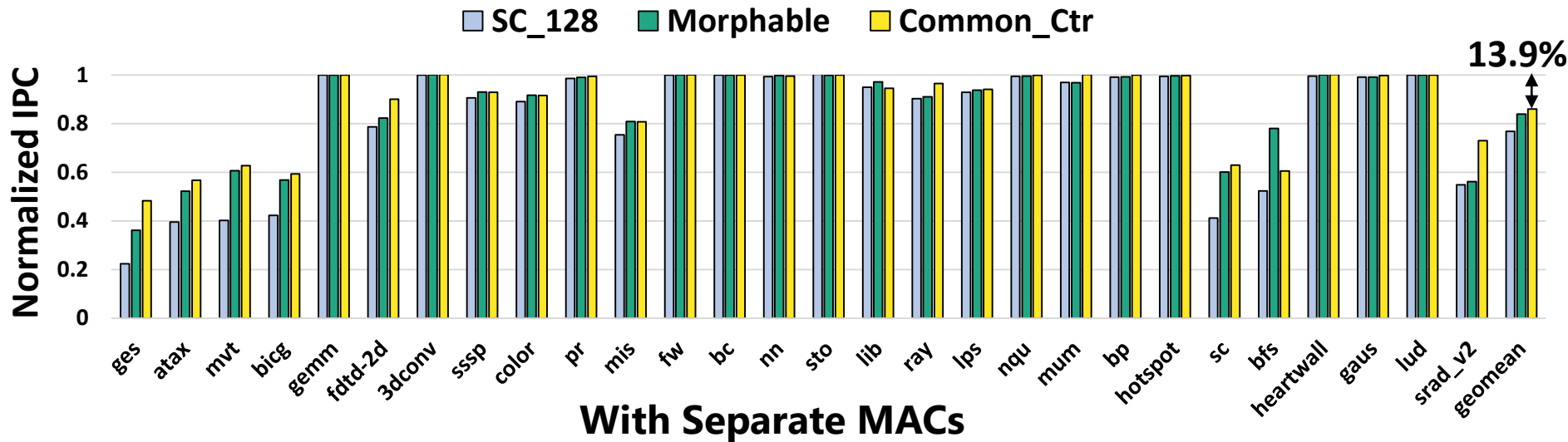
- **Evaluation**

# Methodology

- Simulator: GPGPU-Sim
- Workloads: ISPASS, Rodinia, Polybench, Pannotia
- System configuration: Models NVIDIA TITAN X Pascal GPU

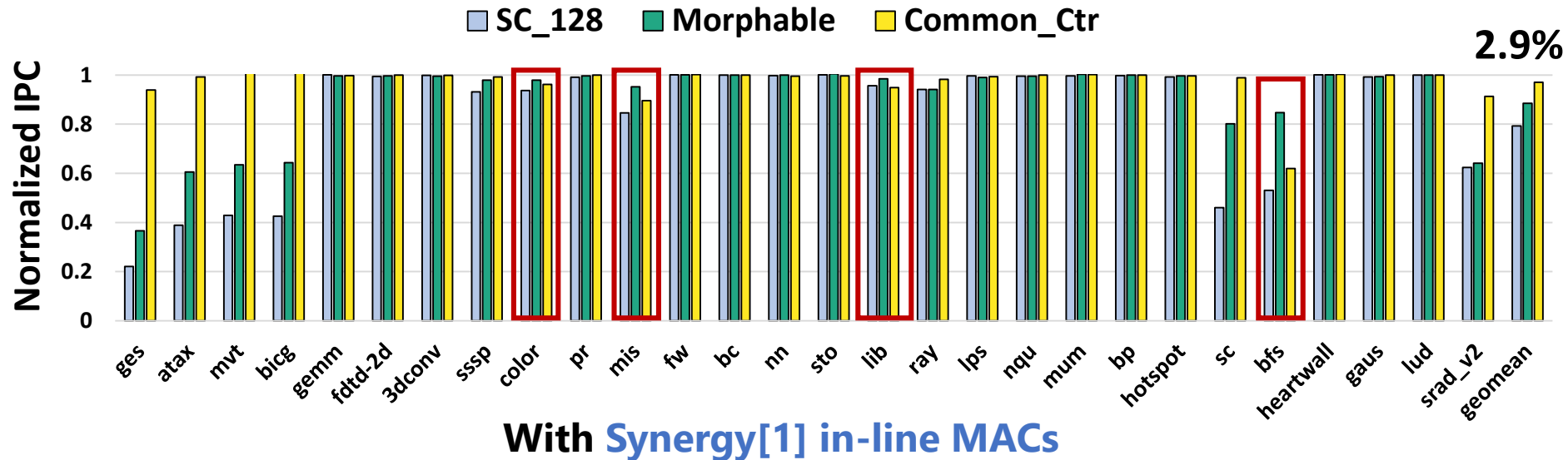| GPU Core Configuration | |
|---|---|
| System overview | 28 SMs, 64 warps per SM |
| Shader core | 1,417 MHz, GTO Scheduler |
| **Cache & Memory Configuration** | |
| L1 cache | 48 KB |
| L2 cache | 3 MB |
| DRAM | GDDR5X 1,251 MHz, 12 GB |
| Counter cache, Hash cache | **16 KB** |
| CCSM cache | **1 KB** |
| Segment size | **128 KB** |
| Number of common ctrs | **15** |

# Performance: Separate MACs

- Performance overhead analysis (Baseline: Non-secure GPU)
  - SC-128: **128-arity** split counter
  - Morphable Counter [MICRO '18]: **256-arity** split counter
  - Common_Ctr: Implemented on top of SC-128 (**128-arity**)



**With Separate MACs**

# Performance: Synergy In-line MACs

- **Common counter** reduces the performance degradation to **2.9%**

**color, mis,lib,bfs** : As kernel runs, # of requests served by **common counters** decreases.



**With Synergy[1] in-line MACs**

[1] :SYNERGY: Rethinking Secure-Memory Design for Error-Correcting Memories, HPCA'18

# More Results in the Paper

- Uniformly updated ratios of real-world GPU Applications

- Hardware area/energy cost for common counter mechanism

- Ratios of LLC misses served by common counters

- Scanning Overheads

- Counter cache sensitivity experiments

**Please Refer to our paper for more details!**

# Conclusion

- **Result**
  - **Common Counter** reduces the performance degradation to **2.9%**

- **Problem**
  - **Memory encryption** is one of the critical bottlenecks for secure GPU memory

- **Key Observation**
  - GPU programs tend to uniformly update memory
  - The number of distinct common counters is small

- **Our Approach**
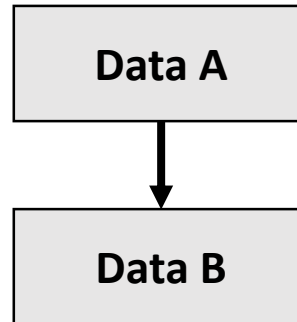  - **Common Counter** provides compressed representation of per-block counters

# Backup Slides

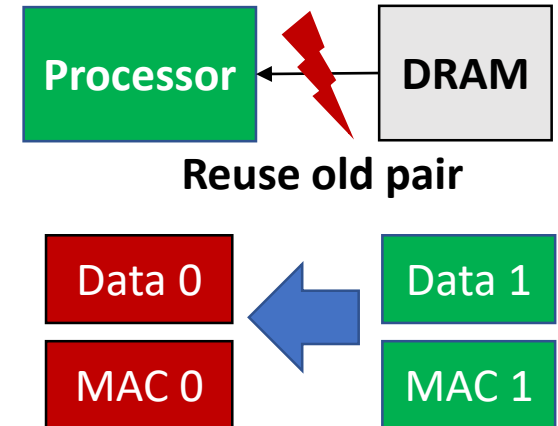# Type of Physical Attacks

| Unauthorized Data Leak | Unauthorized Data Modification | Replay Attack |
|---|---|---|



**Cold Boot Attack**

Data A → Data B

**DMA Attack**

Processor ⚡ DRAM
**Reuse old pair**

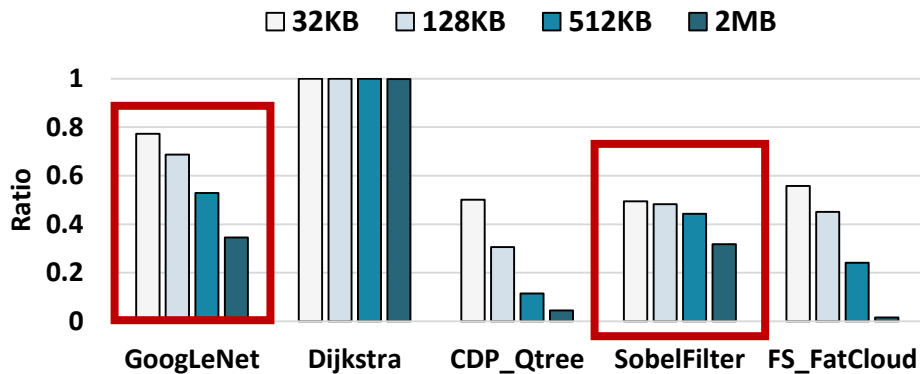Data 0 ← Data 1
MAC 0     MAC 1

**Man in the middle attack**

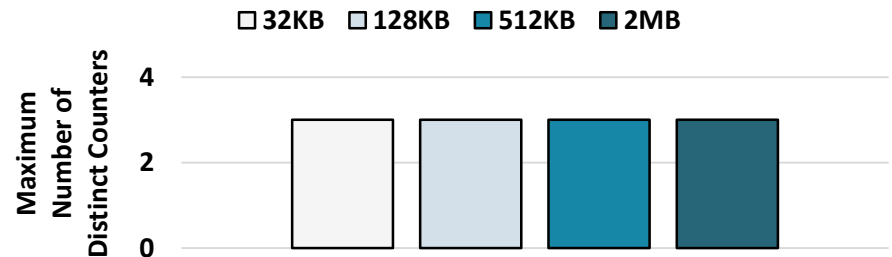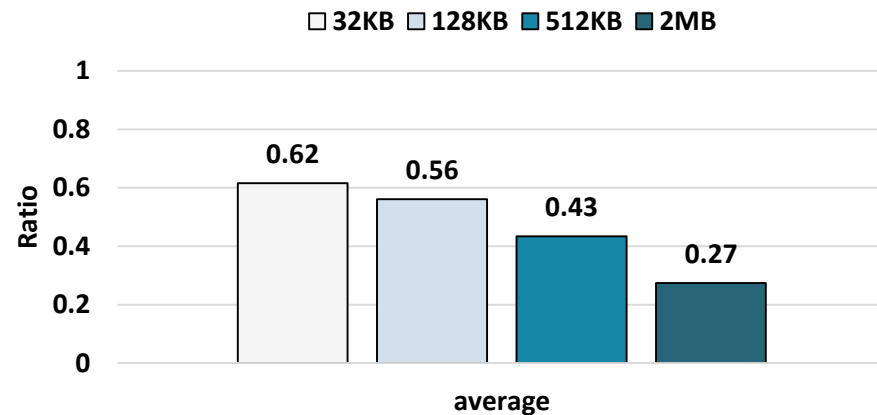**Secure memory needs to protect against these attacks !**

# GPU SW Memory Write Patterns

- Analyze memory read/write behavior by using **NVBit [1]**
  - Collect traces for load/store instructions

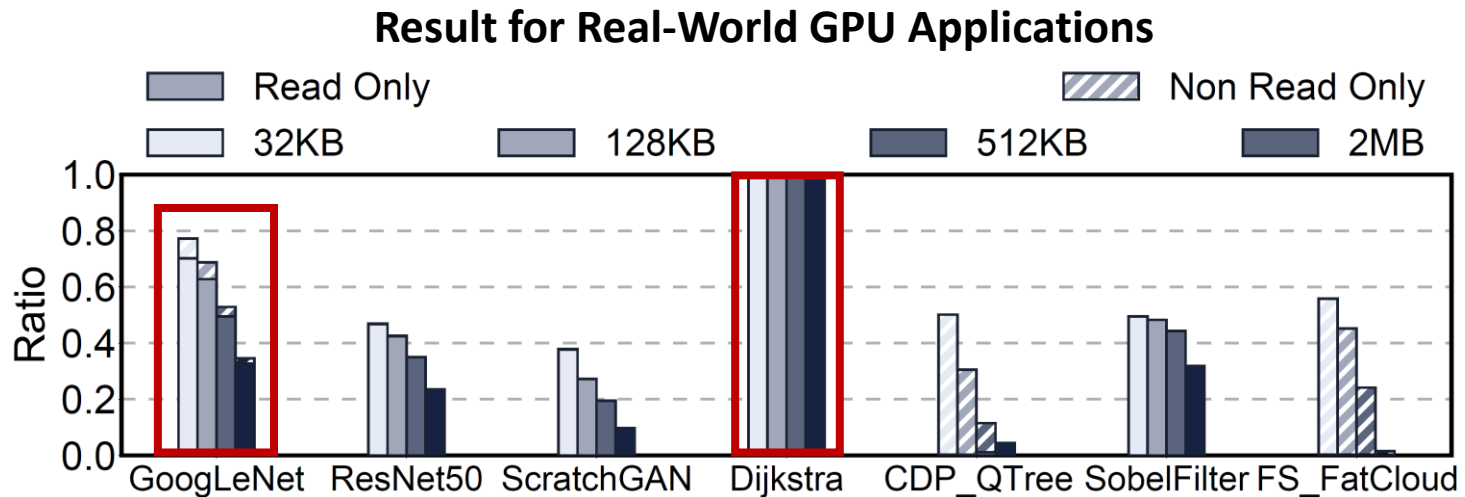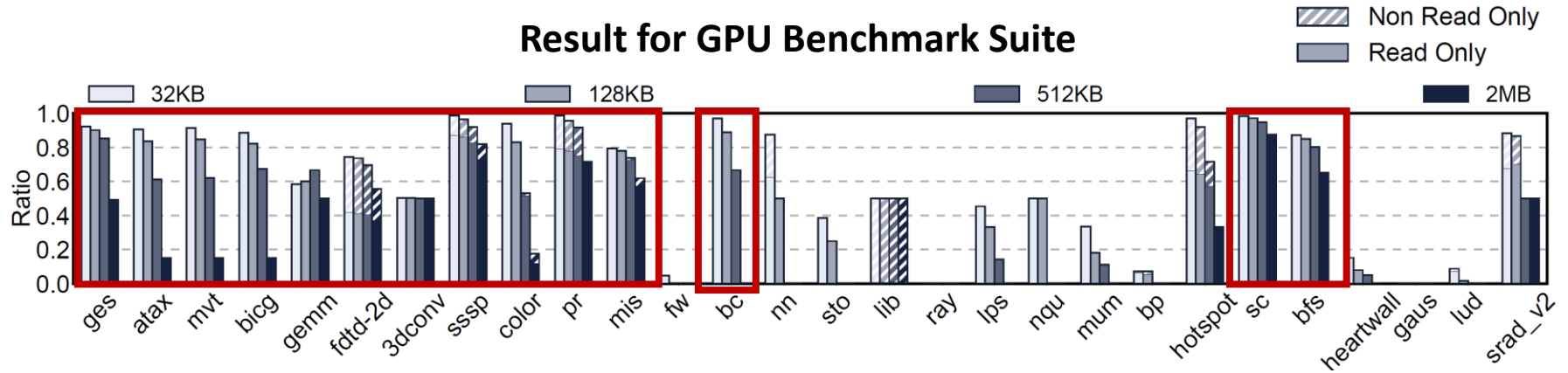**Result of Real-World GPU Applications**
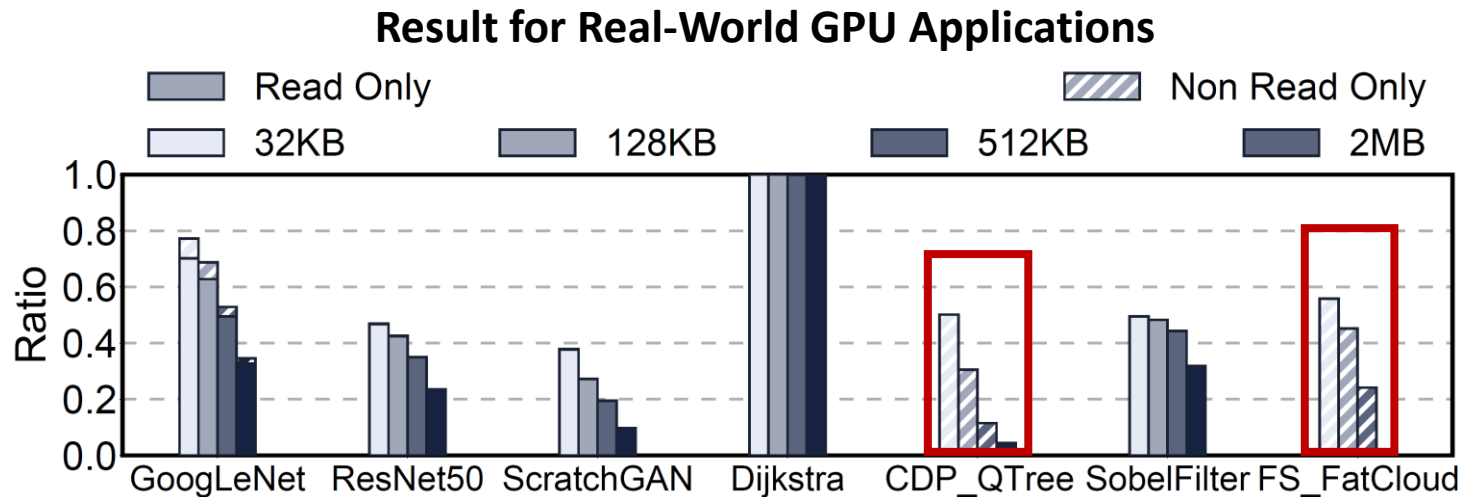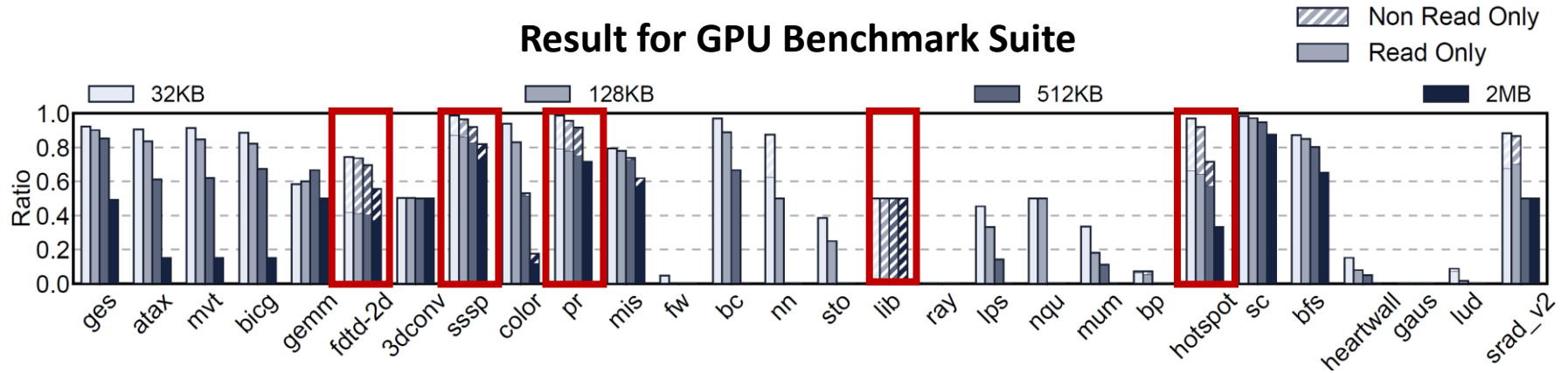


**Result of GPU Benchmark Suite**



[1] : NVBit: A Dynamic Binary Instrumentation Framework for NVIDIA GPUs, MICRO'19

# GPU SW Memory Write Patterns

- Analyze



Result for GPU Benchmark Suite

Result for Real-World GPU Applications

# GPU SW Memory Write Analysis



**Result for GPU Benchmark Suite**

Non Read Only
Read Only

32KB  128KB  512KB  2MB

ges, atax, mvt, bicg, gemm, fdtd-2d, 3dconv, sssp, color, pr, mis, fw, bc, nn, sto, lib, ray, lps, nqu, mum, bp, hotspot, sc, bfs, heartwall, gaus, lud, srad_v2

**Result for Real-World GPU Applications**

Read Only       Non Read Only
32KB      128KB      512KB      2MB

GoogLeNet  ResNet50  ScratchGAN  Dijkstra  CDP_QTree  SobelFilter  FS_FatCloud

# Coverage

- Coverage comparison

| Scheme | Granularity | Per-block Coverage |
|---|---|---|
| split counter | 128B data block | 128 * 16KB = **16KB data** |
| Common counter | 128 KB data block | 256 * 128KB = **32MB data** |

**2048x** efficient coverage

# Performance Result & Counter Coverage



**With Synergy in-line MACs**

Legend (top chart): SC_128, Morphable, Common_Ctr

Legend (bottom chart): Read Only, Non Read Only

# Counter Cache Miss ratio

# Scanning Overhead

- Evaluate scanning procedure

Neglig

| Workload | # of Executed Kernels | Total Scan Size | Ratio* |
|----------|----------------------|-----------------|--------|
| 3dconv | 254 | 32,256 MB | 0.372 % |
| gemm | 1 | 32 MB | 0.090 % |
| bfs | 24 | 4,108 MB | 0.004 % |
| bp | 2 | 390 MB | 0.372 % |
| color | 28 | 5,650 MB | 0.081 % |
| fw | 255 | 2,040 MB | 0.114 % |

*Scanning overhead over total kernel execution time