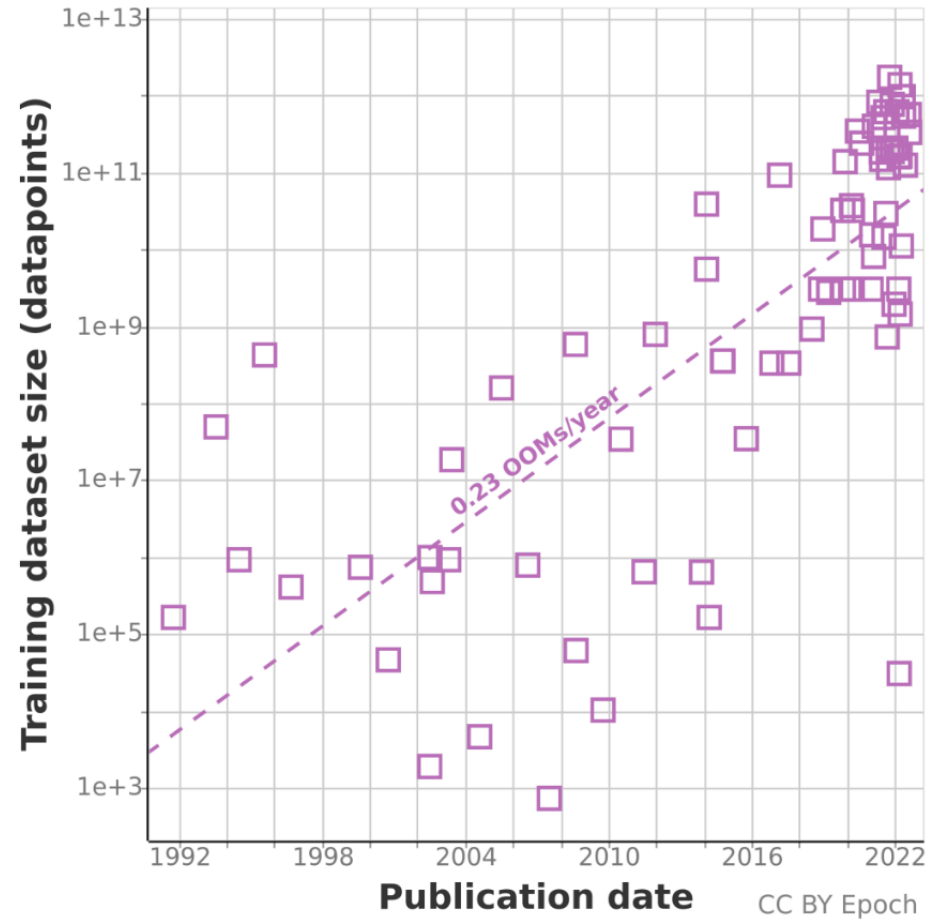# Supporting Secure Multi-GPU Computing with Dynamic and Batched Metadata Management

**Seonjin Na**[1], Jungwoo Kim[2], Sunho Lee[2], Jaehyuk Huh[2]

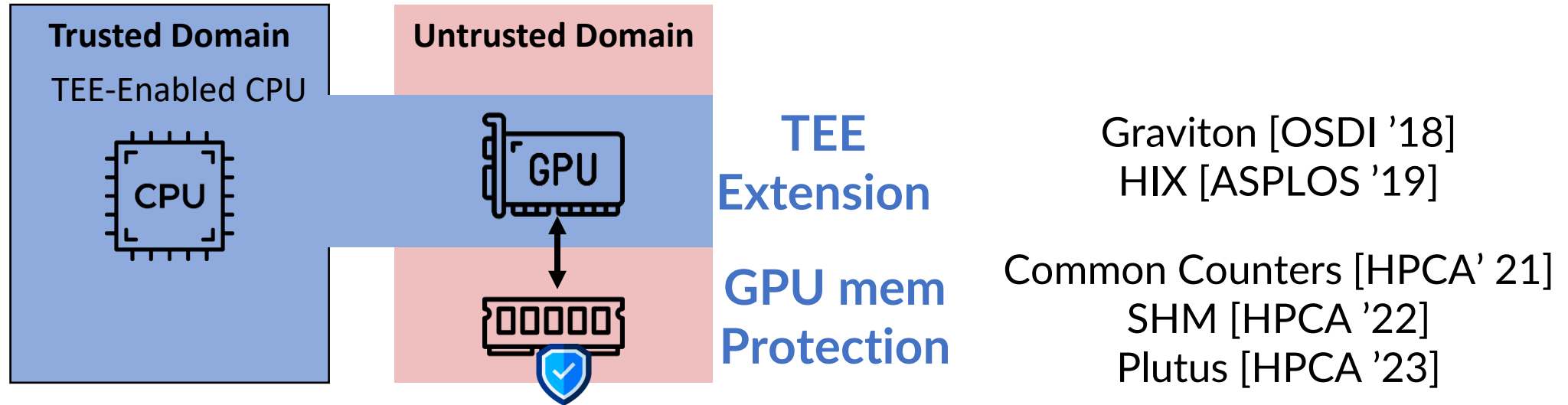[1]Georgia Institute of Technology,  [2]KAIST
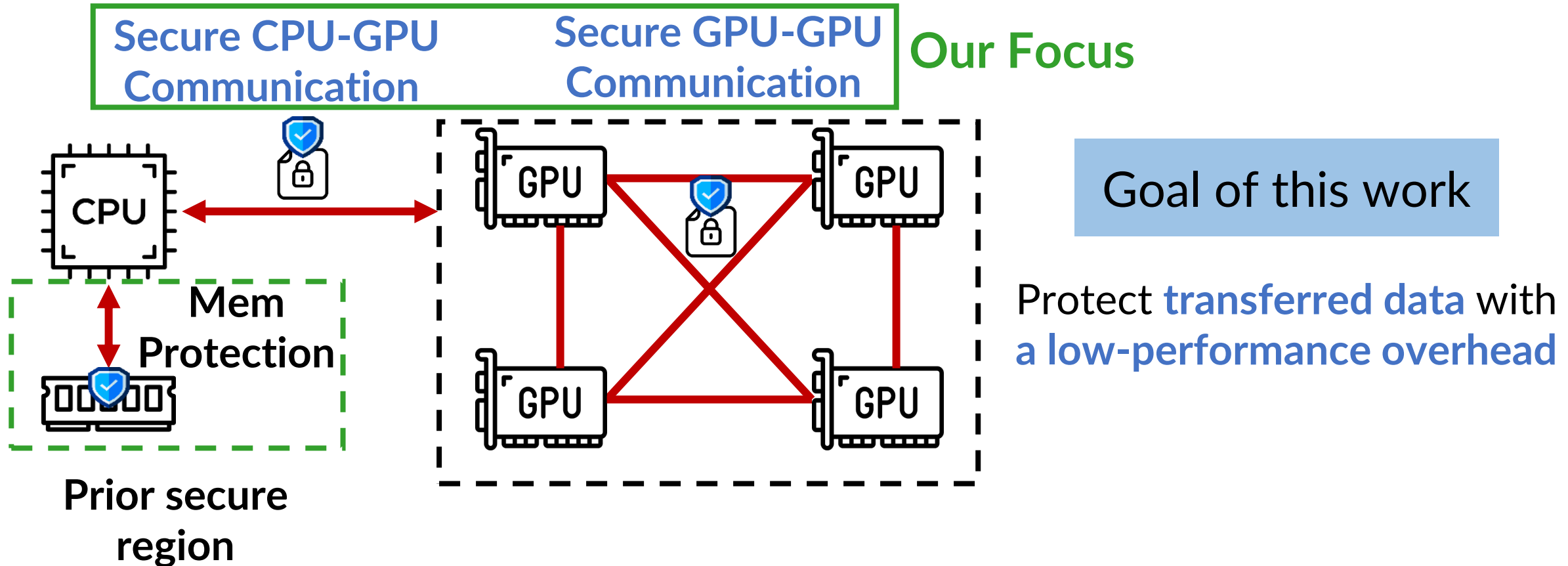
# Importance of Multi-GPU Computing

# Secure GPU Computing Efforts in Academia and Industry



**Trusted Domain**

TEE-Enabled CPU

**Untrusted Domain**

CPU

GPU

**TEE Extension**

**GPU mem Protection**

Graviton [OSDI '18]
HIX [ASPLOS '19]

Common Counters [HPCA' 21]
SHM [HPCA '22]
Plutus [HPCA '23]

## Lack of **data protection mechanism** optimized for multi-GPU systems

Secure CPU-GPU Communication

Secure GPU-GPU Communication

**Our Focus**

CPU

Mem Protection

**Prior secure region**

GPU   GPU

GPU   GPU

Goal of this work

Protect **transferred data** with **a low-performance overhead**

# Contents

- Introduction

- **Background and Motivation**

- Key insights and Main Idea

- Evaluation

**Confidentiality & Integrity**

Authenticated en/decryption

**Freshness**

Replay attack protection

Shared secret key

Sender

Receiver

CPU

Untrusted Interconnect

GPU

| Encrypted data | Additional metadata | MAC |

ID(Sender), Ctr

Re-computed MAC

=?

MAC

Sender

Receiver

CPU

Encrypted data with metadata

GPU

Intercept packet

Older Encrypted Data with metadata

# Background: Protecting Transferred Data through Interconnect

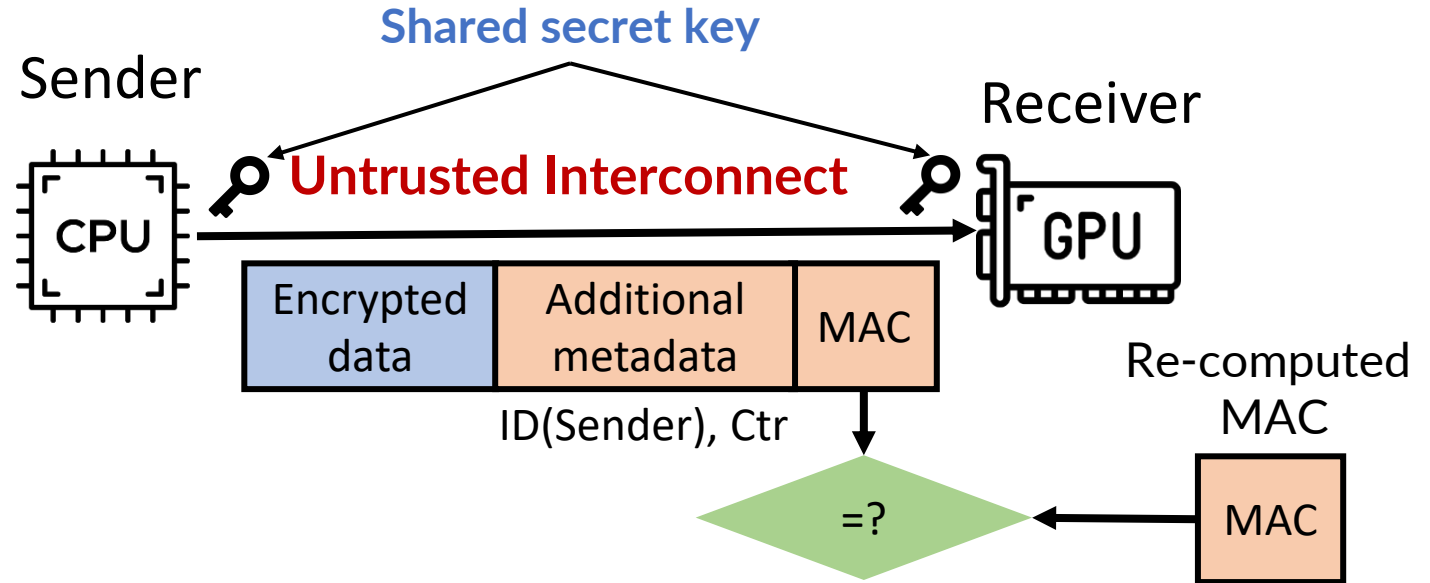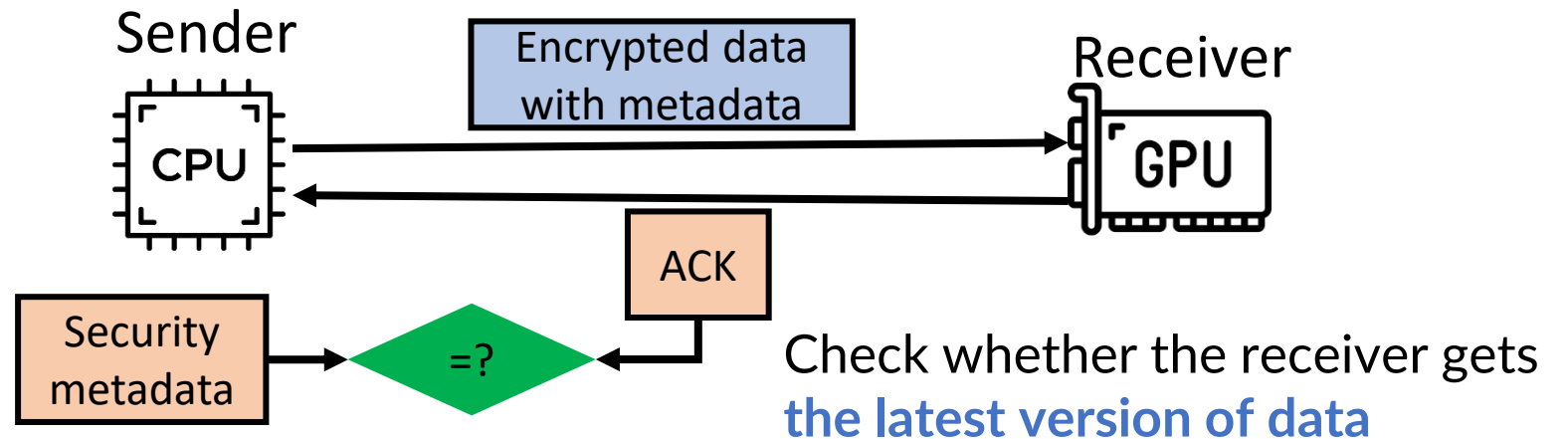**Confidentiality & Integrity**

Authenticated en/decryption

**Freshness**

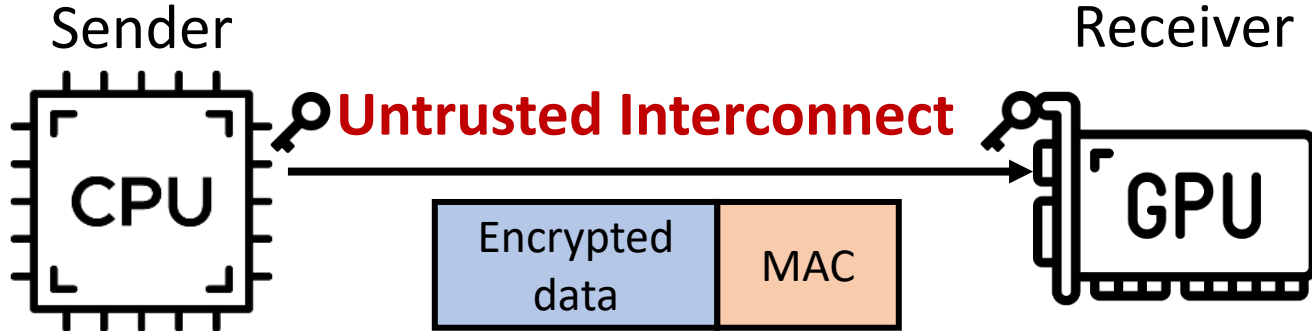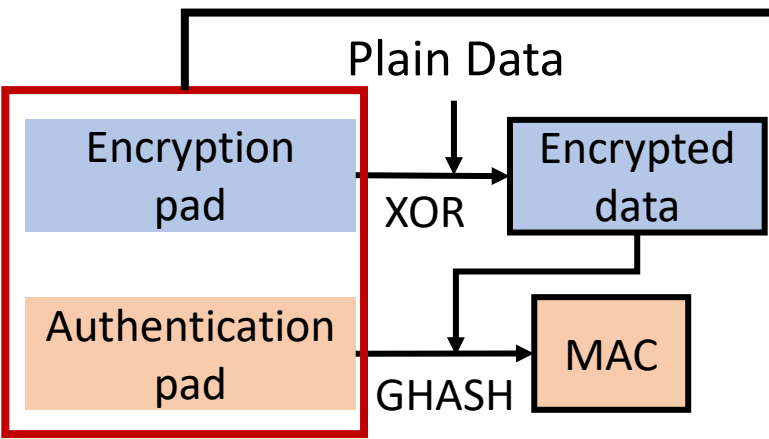Replay attack protection

Sender

Shared secret key

Untrusted Interconnect

CPU

Receiver

GPU

| Encrypted data | Additional metadata | MAC |

ID(Sender), Ctr

Re-computed MAC

=?

MAC

Sender

Encrypted data with metadata

CPU

Receiver

GPU

ACK

Security metadata

=?

Check whether the receiver gets **the latest version of data**

# Authenticated En/Decryption with Pre-Computation [1]

Sender

Receiver

**Untrusted Interconnect**

CPU → GPU

| Encrypted data | MAC |

**Store pre-generated Enc.pad/ Auth. pad**

Plain Data

Naïve approach

Encryption · Generate MAC

Time

Pre-computed pads

XOR · GHASH

Time

| Encryption pad | → XOR → | Encrypted data |
| Authentication pad | → GHASH → | MAC |

| Valid | Enc. pad | Auth. pad | Ctr |
|---|---|---|---|
| 1bit | 512bit | 64bit | 64bit |

**Pad table**
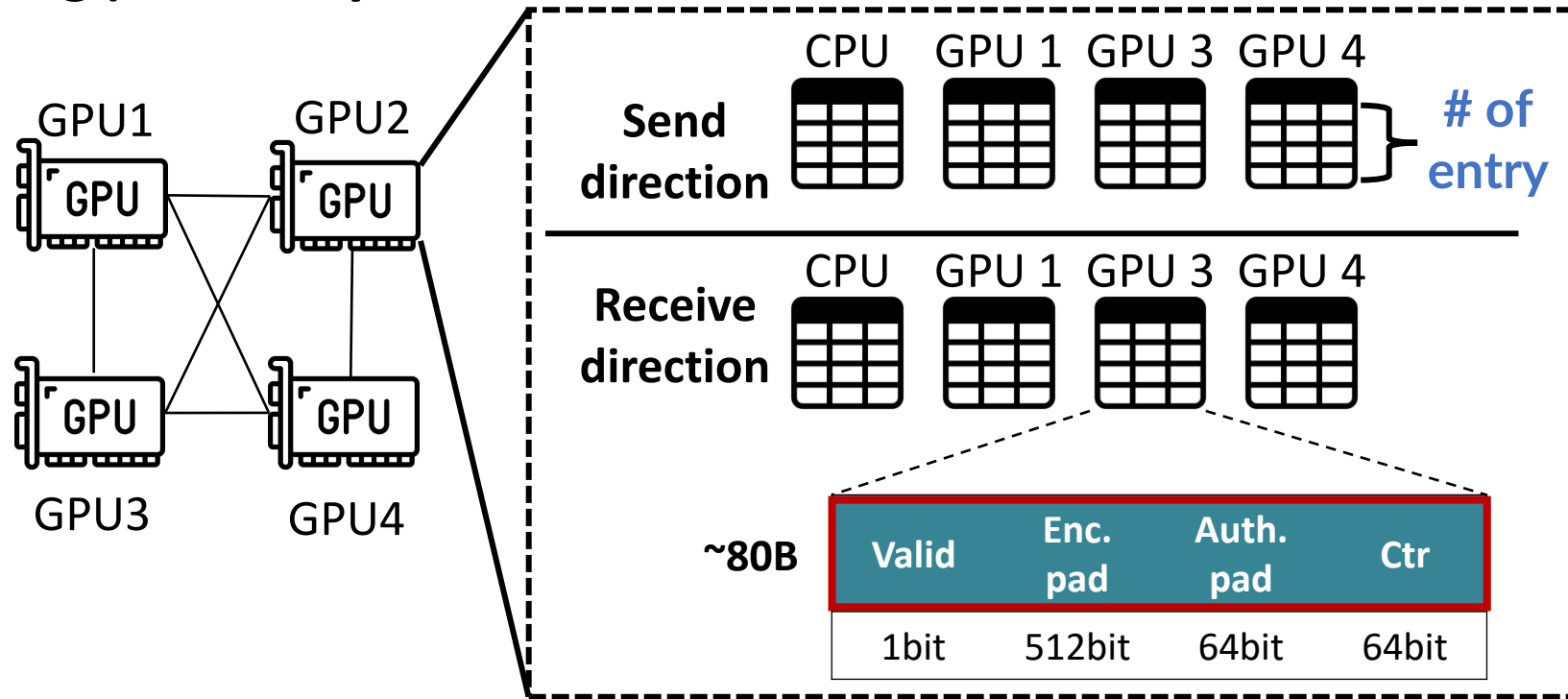
**Stored in on-chip area**

[1] Efficient data protection for distributed shared memory multiprocessors, PACT'06

# Prior Pad Table Management (Private) [1]

- Maintains **same # of pad entries** for all commu. pairs in a system

**E.g.) 4-GPU System**


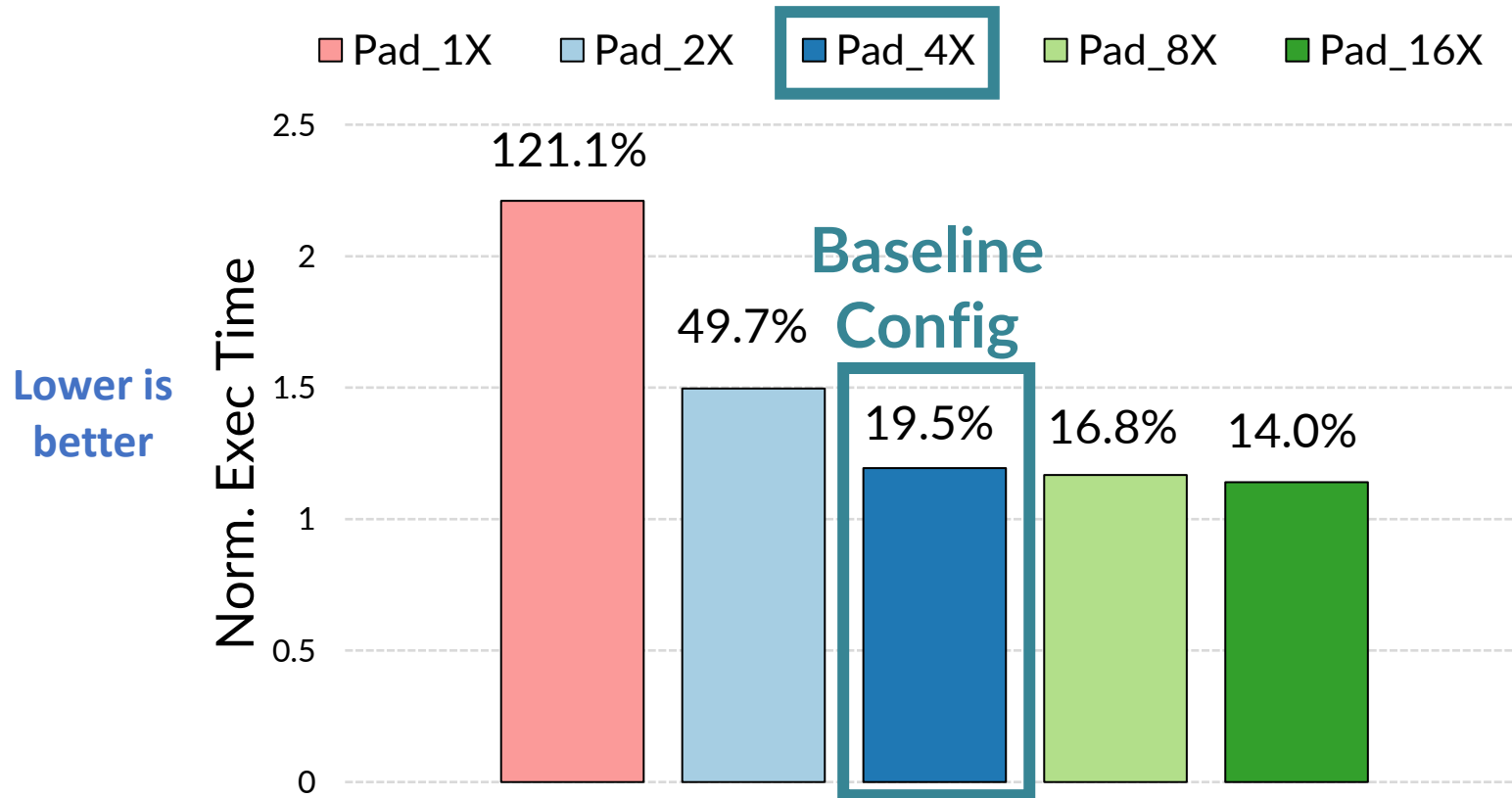
Increasing **# of pad table entries**

Performance / On-chip storage overhead

| Valid | Enc. pad | Auth. pad | Ctr |
|-------|----------|-----------|-----|
| 1bit | 512bit | 64bit | 64bit |

~80B

[1] Efficient data protection for distributed shared memory multiprocessors, PACT'06

Baseline:
**Unsecure 4 GPU**

**Use 4 pad table entries for all commu. pairs**

■ Pad_1X  ■ Pad_2X  ■ Pad_4X  ■ Pad_8X  ■ Pad_16X

**Lower is better**

121.1%

49.7%

**Baseline Config**

19.5%  16.8%  14.0%

Norm. Exec Time

2.5
2
1.5
1
0.5
0

[1] Efficient data protection for distributed shared memory multiprocessors, PACT'06

# Performance Breakdown Analysis

- Secure multi-GPU incurs average **19.5%** performance degradation
  - Auth. en/decryption: **8.2% slowdown**, Metadata traffic: **11.3% slowdown**

**+ 8.2%**      **+ 11.3%**
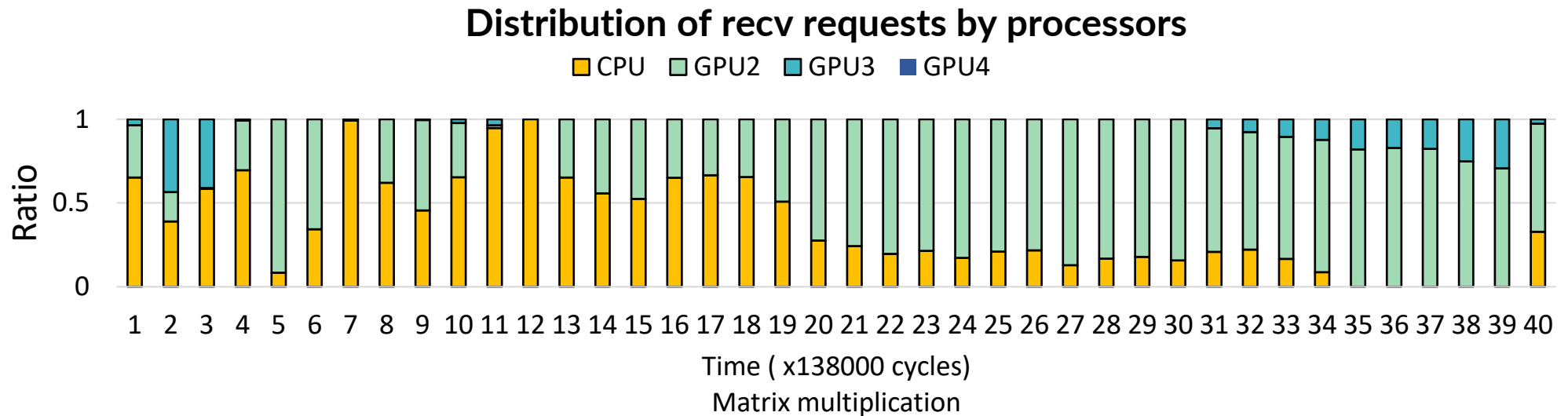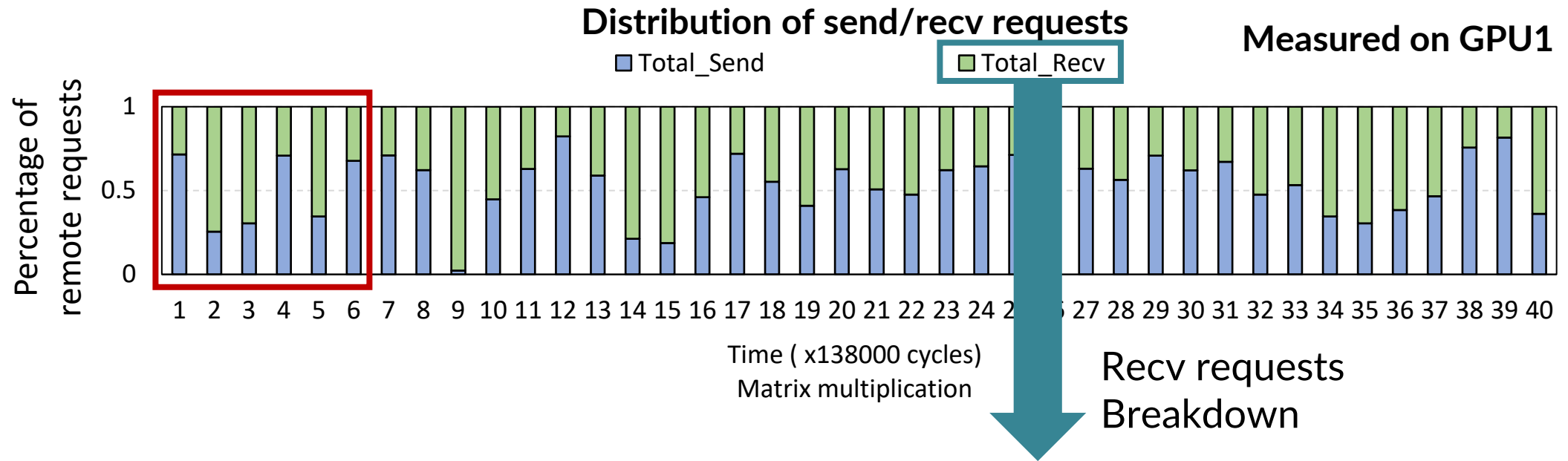
**Performance bottlenecks of secure communication**
**1. Authenticated en/decryption**
**2. Additional security metadata traffic**

better  Norm. Ex

0.6

0.4

0.2

0

# Contents

- Introduction

- Background and Motivation

- **Key insights and Main Idea**

- Evaluation

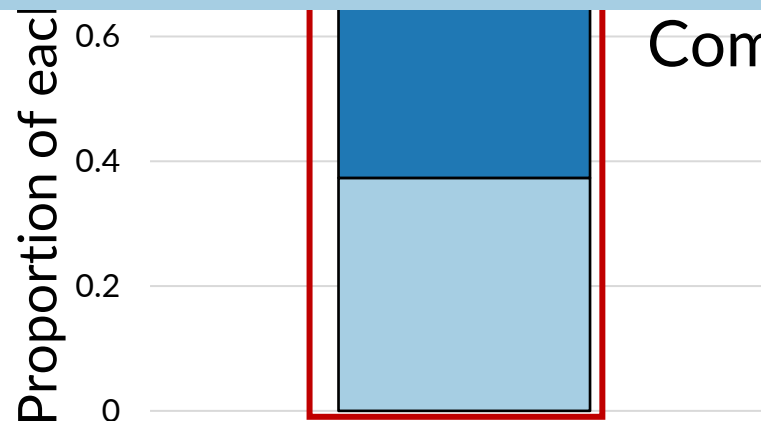# Key Insight 1: Dynamic Behavior of Communication Patterns



**Distribution of send/recv requests**

Measured on GPU1

Total_Send    Total_Recv

Percentage of remote requests

Time ( x138000 cycles)
Matrix multiplication

Recv requests
Breakdown

**Distribution of recv requests by processors**

CPU   GPU2   GPU3   GPU4

Ratio

Time ( x138000 cycles)
Matrix multiplication

- Analyze distribution of cycles for gathering 16 transmitted data blocks

Cycle distribution

**Our Key Observations**
1. Dynamic behavior of communication patterns
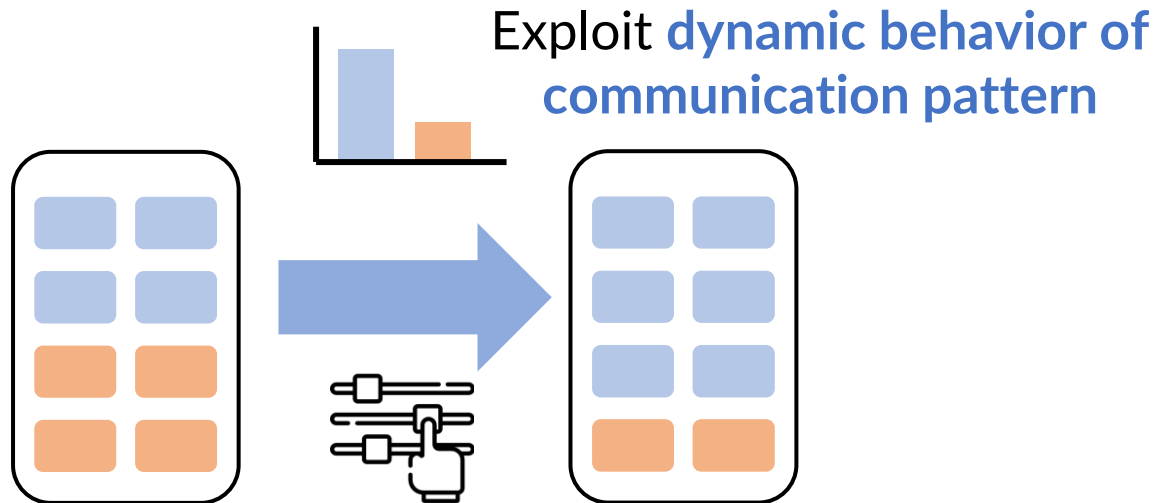2. Burstiness of communication

Communication between processors occurs **within a short period**

# Main Idea of This Work
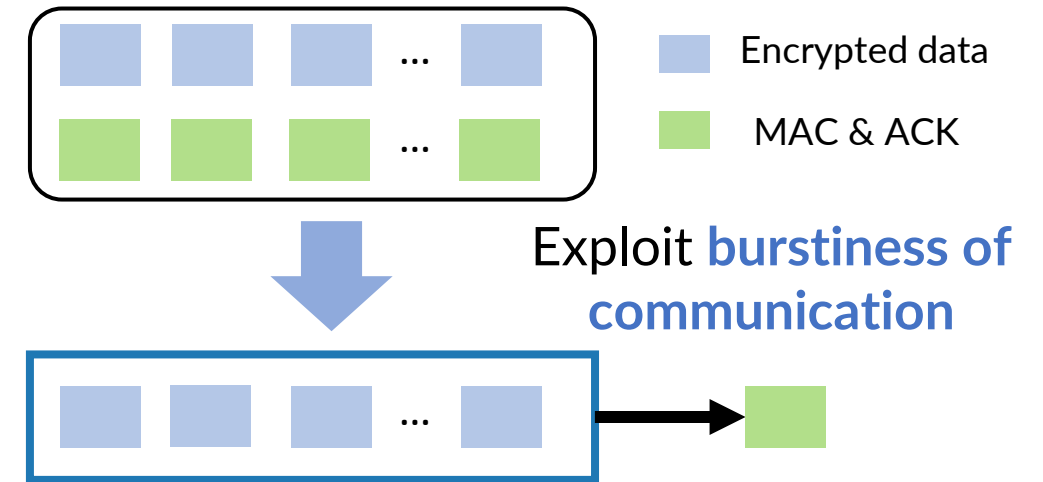
**Challenge 1: authenticated en/decryption overhead**

**Dynamic** pad table management

Exploit **dynamic behavior of communication pattern**

Increase opportunity to hide authenticated en/decryption latency
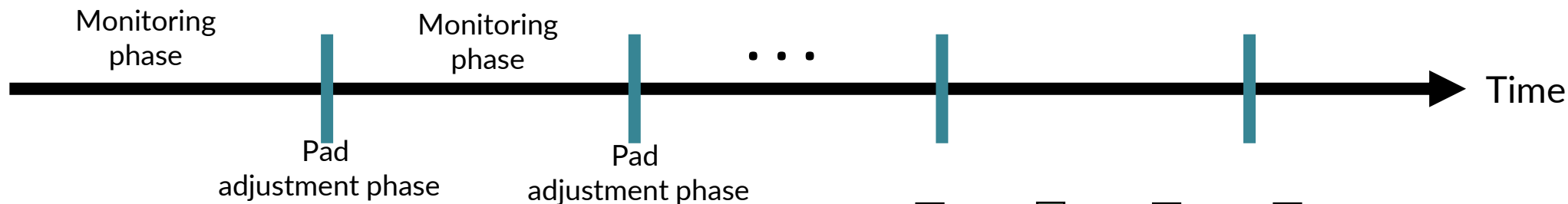
**Challenge 2: additional bandwidth by security metadata**

**Batched** MAC generation& verification

Encrypted data

MAC & ACK

Exploit **burstiness of communication**
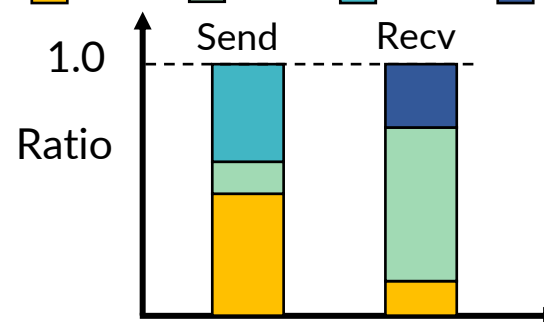
Reduce security metadata traffic

# Dynamic Pad Table Management

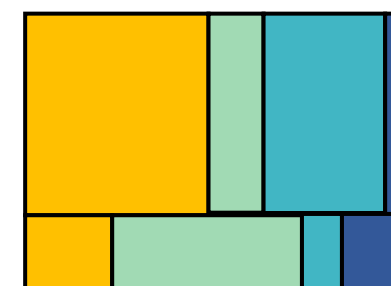- **Dynamically adjust pad table entries** based on communication pattern
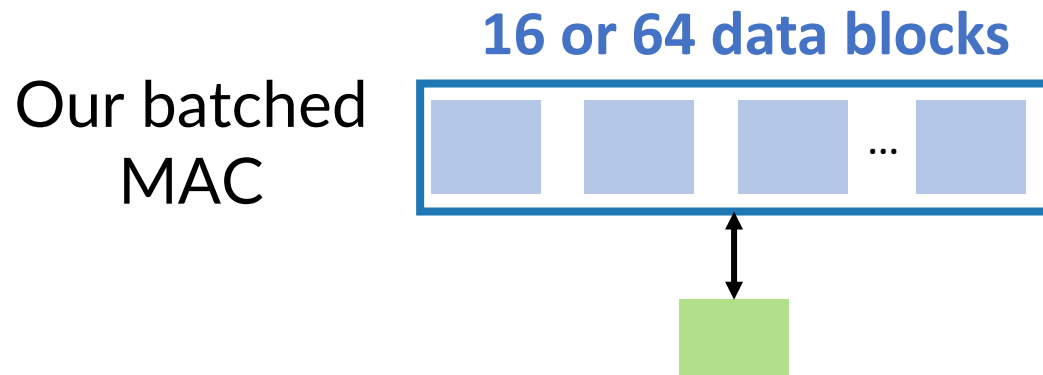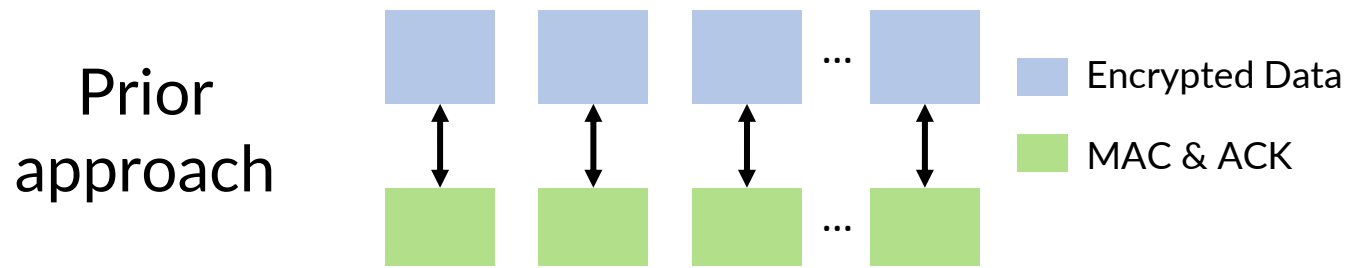
# Batched MAC Generation & Verification

- Generate **coarse-grained MAC** to reduce metadata bandwidth

**Prior approach**



Encrypted Data

MAC & ACK

Example: **4KB page migration**

**Traffic**   64 Data + 64 Metadata (for decryption)
**+ 64 MAC & ACK**

**Our batched MAC**

**16 or 64 data blocks**



**Traffic**   64 Data + 64 Metadata (for decryption)
**+ 1 Batch info + 1 MAC & ACK**

# Contents

- Introduction

- Background and Motivation

- Key insights and Main Idea

- **Evaluation**

# Evaluation Methodology

- Simulator: MGPUSim [ISCA '19]

- Workloads: 17 apps from various benchmark suites
  - AMD APP SDK, DNN Mark, Hetero Mark, Polybench, SHOC benchmark suites

- System configuration: Models 4 GPU system (AMD R9Nano GPU)

| GPU Configuration | |
|---|---|
| Compute Unit | 64 CUs per GPU, 1.0 GHz |
| L1 Inst / Vector / Scalar Caches<br>Shared L2 Cache | 16 KB / 32KB / 16KB<br>2MB |
| DRAM | 4GB HBM Memory, 512 GB/s |
| CPU-GPU, GPU-GPU Interconnect | 32 GB/s, 50 GB/s |
| Security Configuration | |
| Authenticated encryption/decryption | 40 cycles [1,2] |

[1] Adaptive Security Support for Heterogenous Memory on GPUs, HPCA '22
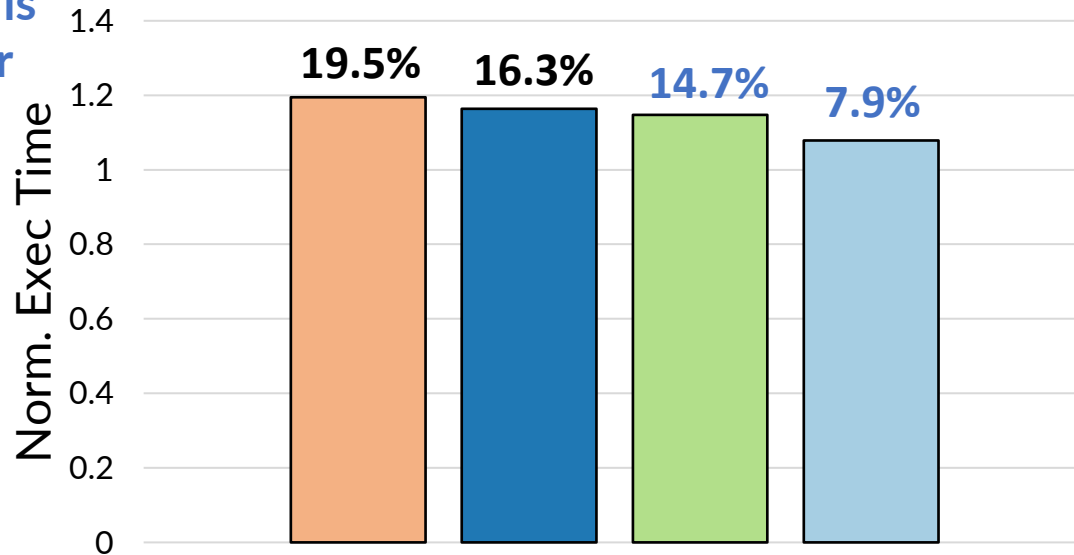[2] Plutus: Bandwidth-Efficient Memory Security for GPUs, HPCA'23

# Performance Comparison Result

- Compared with two different mechanisms
  - **Private**[1]: Uses same number of pad entries for all pairs
  - **Cached**[1]: Allocates pad table entries like LRU cache
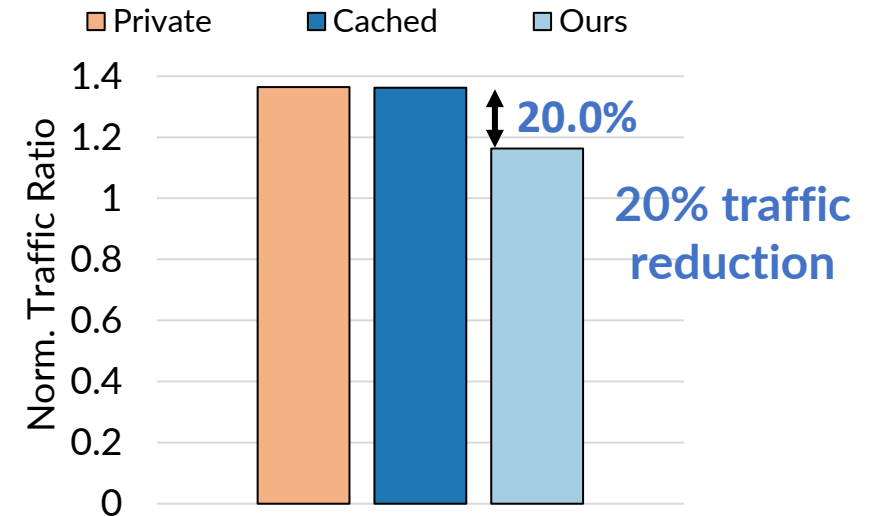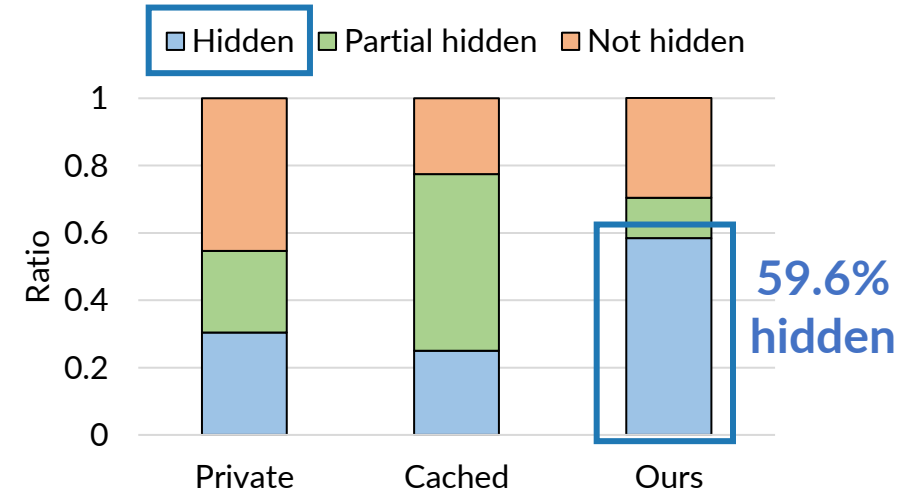
All configs have the **same total # of pad table entries**



**Higher is better**

59.6% hidden

20.0%

20% traffic reduction

Lower is better

19.5%  16.3%  14.7%  7.9%

[1] Efficient data protection for distributed shared memory multiprocessors, PACT'06

# More Results in the Paper

- Scalability study to the number of GPUs

- Sensitivity study to authenticated encryption/decryption latency

- Hardware overhead of our design

# Summary

- **Problem**
  - Secure communication degrades multi-GPU system performance

- **Key Idea**
  - **Dynamic pad table management** exploit dynamic communication patterns
  - **Batched MAC generation** leverage burstiness nature of communication

- **Evaluation results**
  - Reduces perf. overhead by **11.6%**, **8.4%** compared to Private, Cached

# Backup Slides

# Performance Comparison Result

- Compared with two different mechanisms
  - **Private**: Uses fixed number of pad entries
  - **Cached**: Manages pad table entries like LRU cache

**Lower is better**