

# LÆGIS: Pinpointing and Addressing Performance Overheads of GPU-based Confidential Computing

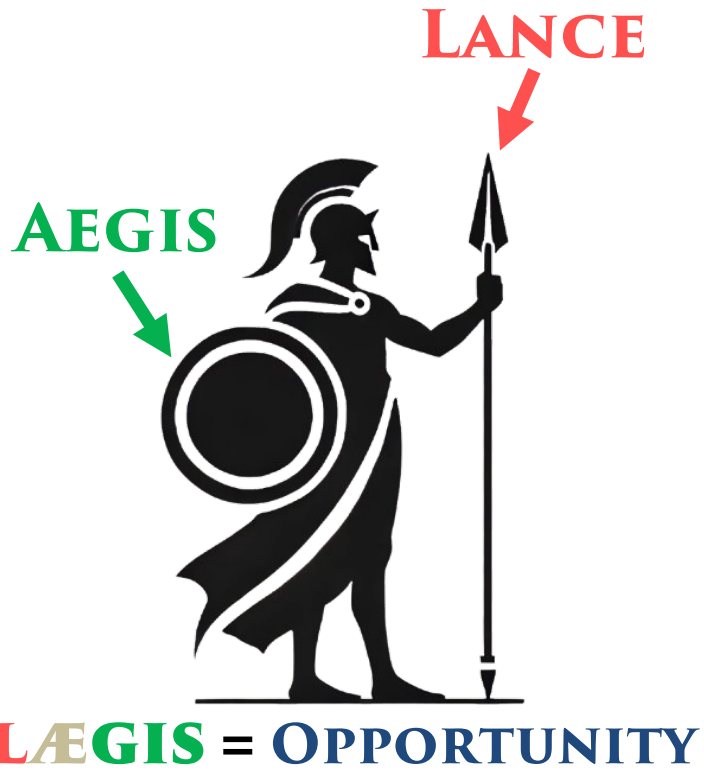
Yang Yang, Adwait Jog  
{yangyang, ajog} @virginia.edu

*Insight Computer Architecture Lab*



ISCA 2026

# Summary



- ❖ We identify three inefficiencies in GPU-based CC when it comes to UVM.
- ❖ We propose **LÆGIS**, a design that opportunistically pre-encrypts pages in CC and uses secure **HBM** for IV management.
- ❖ **LÆGIS** introduces an in-memory IV Bank in 3D-stacked HBM, decoupling encryption from CPU–GPU synchronization and removing the need for integrity trees.

**Performance**



**Security**

Average speedup of **2.22×** (up to **3.13×**)  
and **2.74×** (up to **5.05×**) under  
default and aggressive prefetching.

# Outline

- **INTRODUCTION**
- **BACKGROUND**
- **UVM PERFORMANCE DISSECTION UNDER CC**
- **LÆGIS: OVERVIEW & IMPLEMENTATION**
- **EVALUATION**
- **TAKEAWAY**

# Outline

- **INTRODUCTION**
- **BACKGROUND**
- **UVM PERFORMANCE DISSECTION UNDER CC**
- **LÆGIS: OVERVIEW & IMPLEMENTATION**
- **EVALUATION**
- **TAKEAWAY**

# Needs for Security and Privacy

## □ Data is Private

- ▣ Client privacy (LLMs, ...)
- ▣ Regulated (GDPR, PII, ...)

## □ Data is Massive

- ▣ Outsource compute & data to public cloud Infra.

**What techniques can be deployed to strengthen trust in cloud services?**

# Protection Principles

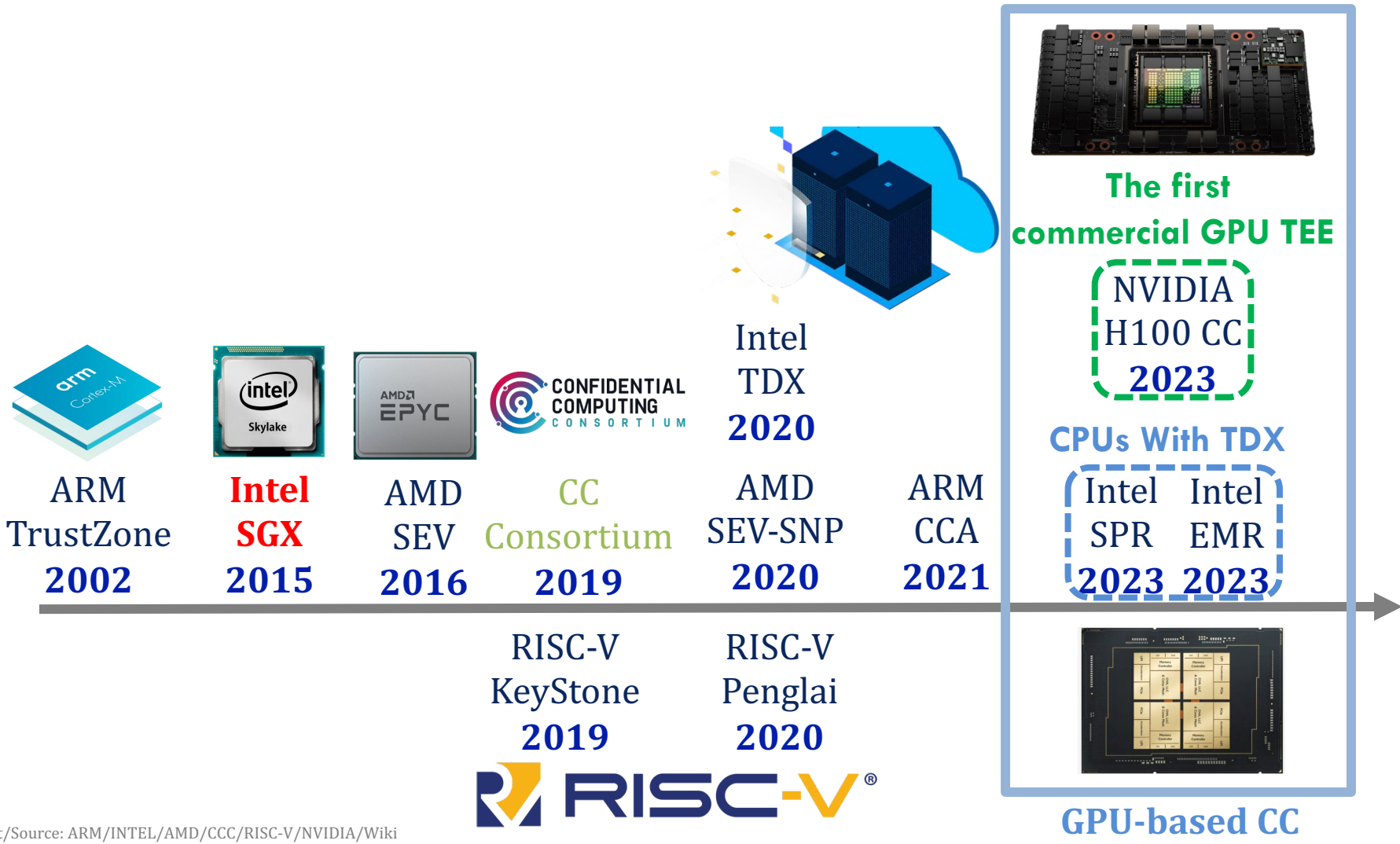
## Hardware-based TEE



**Isolation + Encryption**

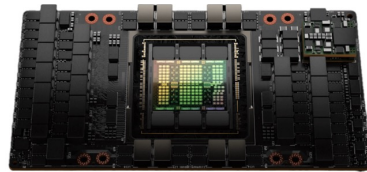
**Trusted/Confidential Computing**

# Trusted Computing



Credit/Source: ARM/INTEL/AMD/CCC/RISC-V/NVIDIA/Wiki

# Trusted Computing

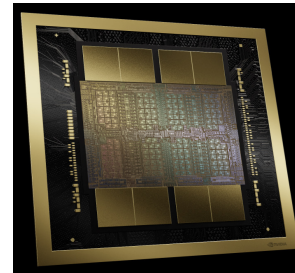


The first commercial GPU TEE

NVIDIA  
H100 CC  
2023

CPU's With TDX

Intel SPR  
Intel EMR  
2023 2023



NVIDIA  
Blackwell  
2024-2025



NVIDIA  
Vera-Rubin  
2026

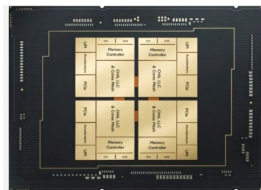
Rack-scale CC

Intel  
TDX  
2020

AMD  
SEV-SNP  
2020

ARM  
CCA  
2021

RISC-V  
Penglai  
2020



GPU-based CC



Credit/Source: ARM/INTEL/AMD/CCC/RISC-V/NVIDIA/Wiki

# CC in Real World

Google Cloud

Overview Solutions Products Pricing Resources

Security Threat Intel SecOps Cloud security Consulting Security resources

## Confidential Computing

Protect data in-use with Confidential VMs, Confidential GKE, Confidential Dataflow, Confidential Dataproc, and Confidential Space.

Go to console

- ✓ Secure your data by keeping
- ✓ Simple easy-to-use deployment
- ✓ Confidential collaboration wh

What is Confidential Computing?

Learn / Azure / Confidential Computing /

Ask Learn

## Azure Confidential GPU options

In this article

- Sizes
- Azure CLI

on EPYC processors with SEV-SNP  
this VM SKU Trusted Execution  
GPU and attached GPU, enabling  
the GPU.

# Azure Confidential Computing

Increase data privacy by protecting data in use

Get started with Azure

Apple Security Research

Overview Blog Bounty Research Device Submit a Report

## Blog

June 10, 2024

# Private Cloud Compute: A new frontier for AI privacy in the cloud

Written by Apple Security Engineering and Architecture (SEAR), User Privacy, Core Operating Systems (Core OS), Services Engineering (ASE), and Machine Learning and AI (AIML)

Get started with a Nitro-based instance today

to console Create account

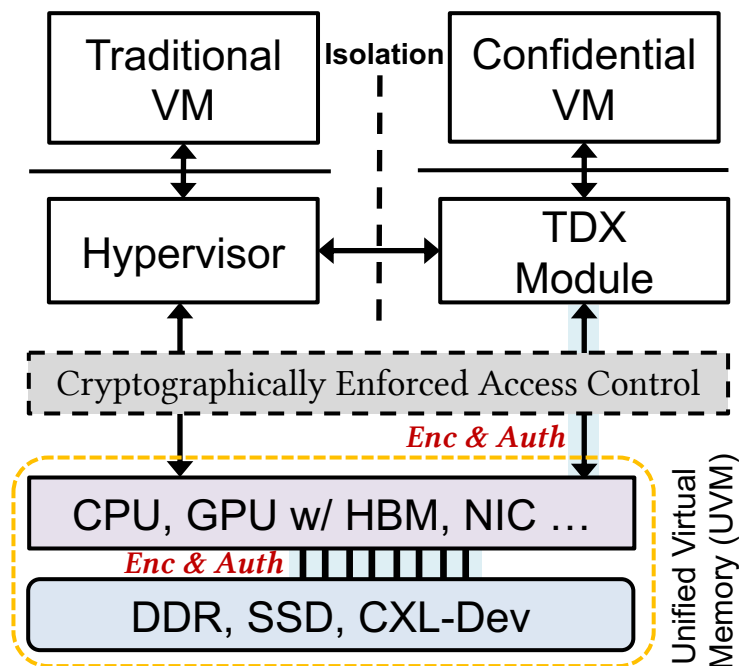
## Protecting data in use

AWS confidential computing is always on. There is no mechanism for any AWS operator to access customers' Amazon Elastic Compute Cloud (Amazon EC2) instances within the AWS Nitro System.

# CC w/ GPU is practical

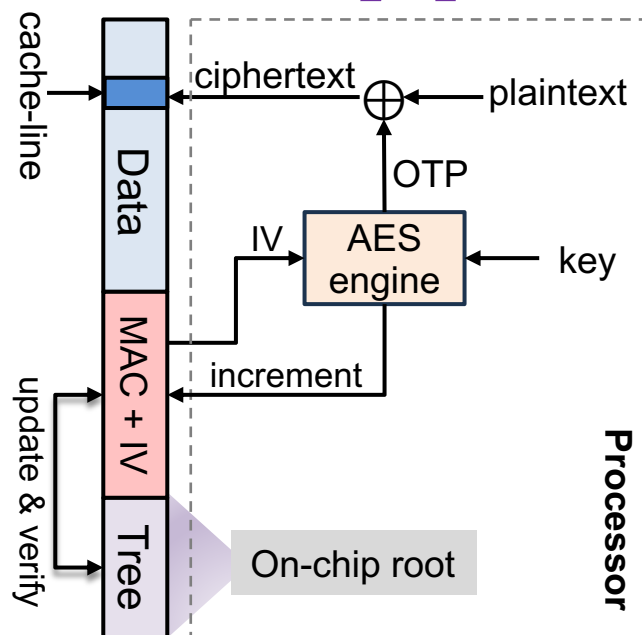
# Closer Look

## Isolation



(a) CC with UVM

## Encryption



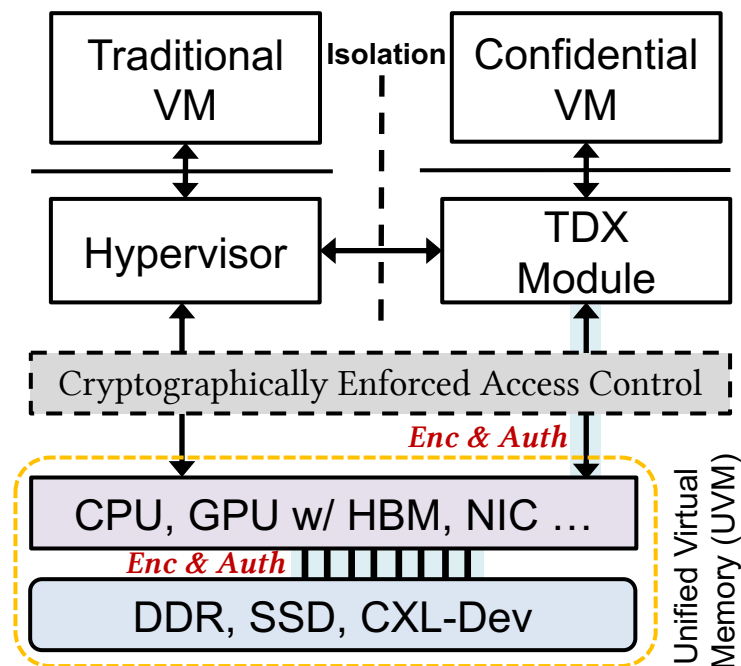
(b) Traditional memory encryption

IV: Initialization Vector

$$\text{IV}++; \text{OTP} = \text{AES}_K(\text{IV} || \text{CTR}); \text{Ct} = \text{Pt} \oplus \text{OTP}$$

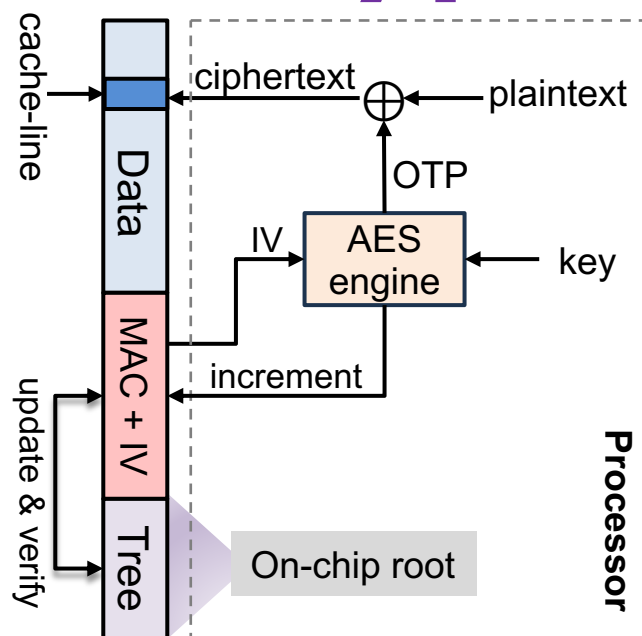
# Closer Look

## Isolation



(a) CC with UVM

## Encryption

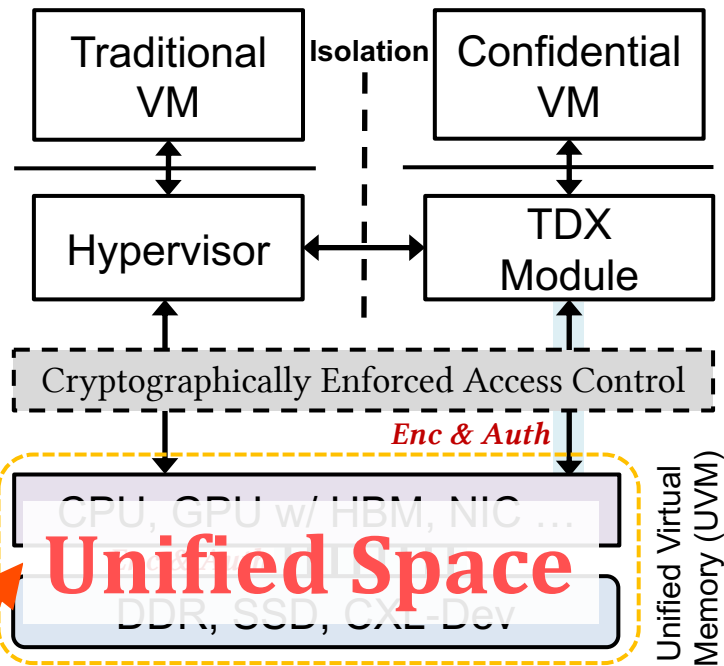


(b) Traditional memory encryption

**The memory problem!**

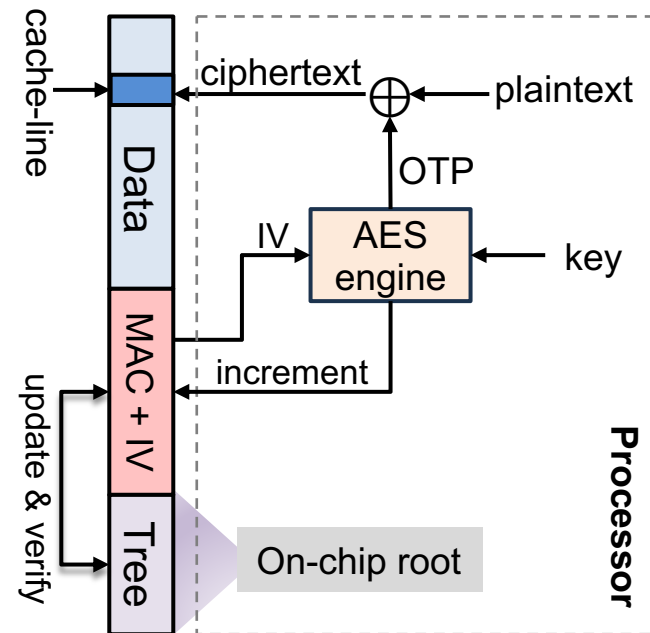
# Closer Look

## Isolation



(a) CC with UVM

## Encryption

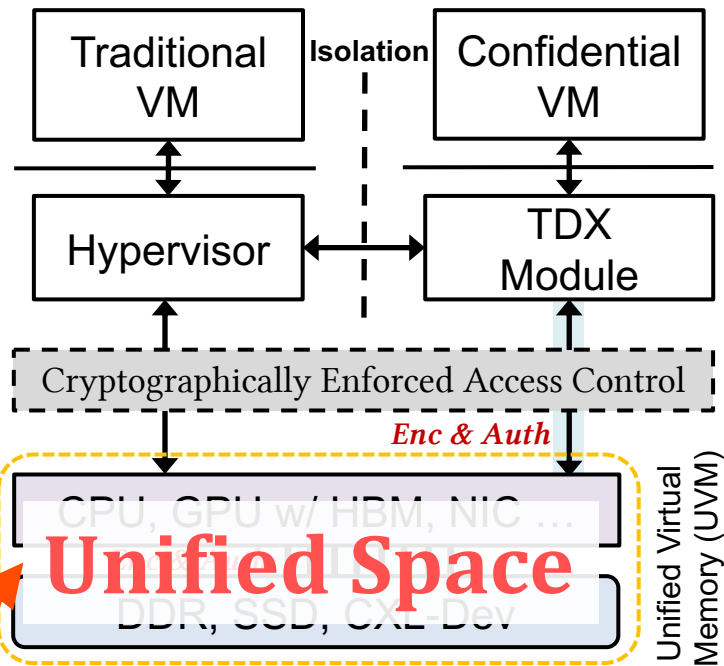


(b) Traditional memory encryption

**The memory problem!**

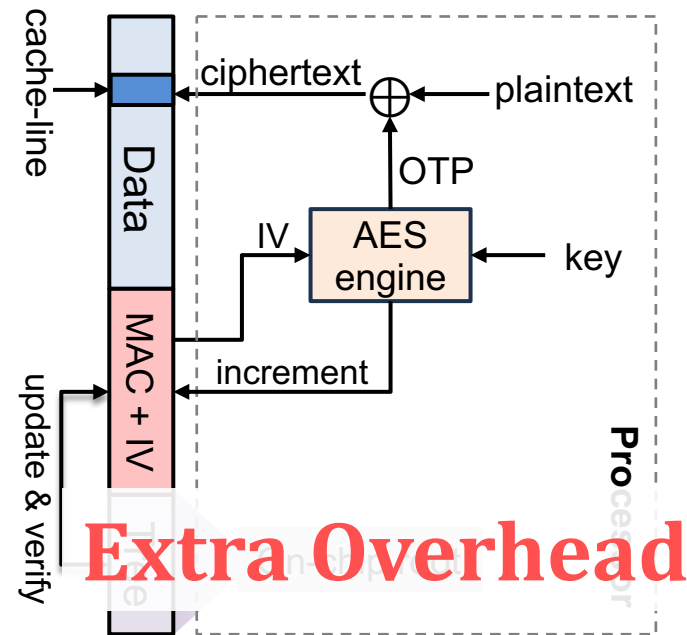
# Closer Look

## Isolation



(a) CC with UVM

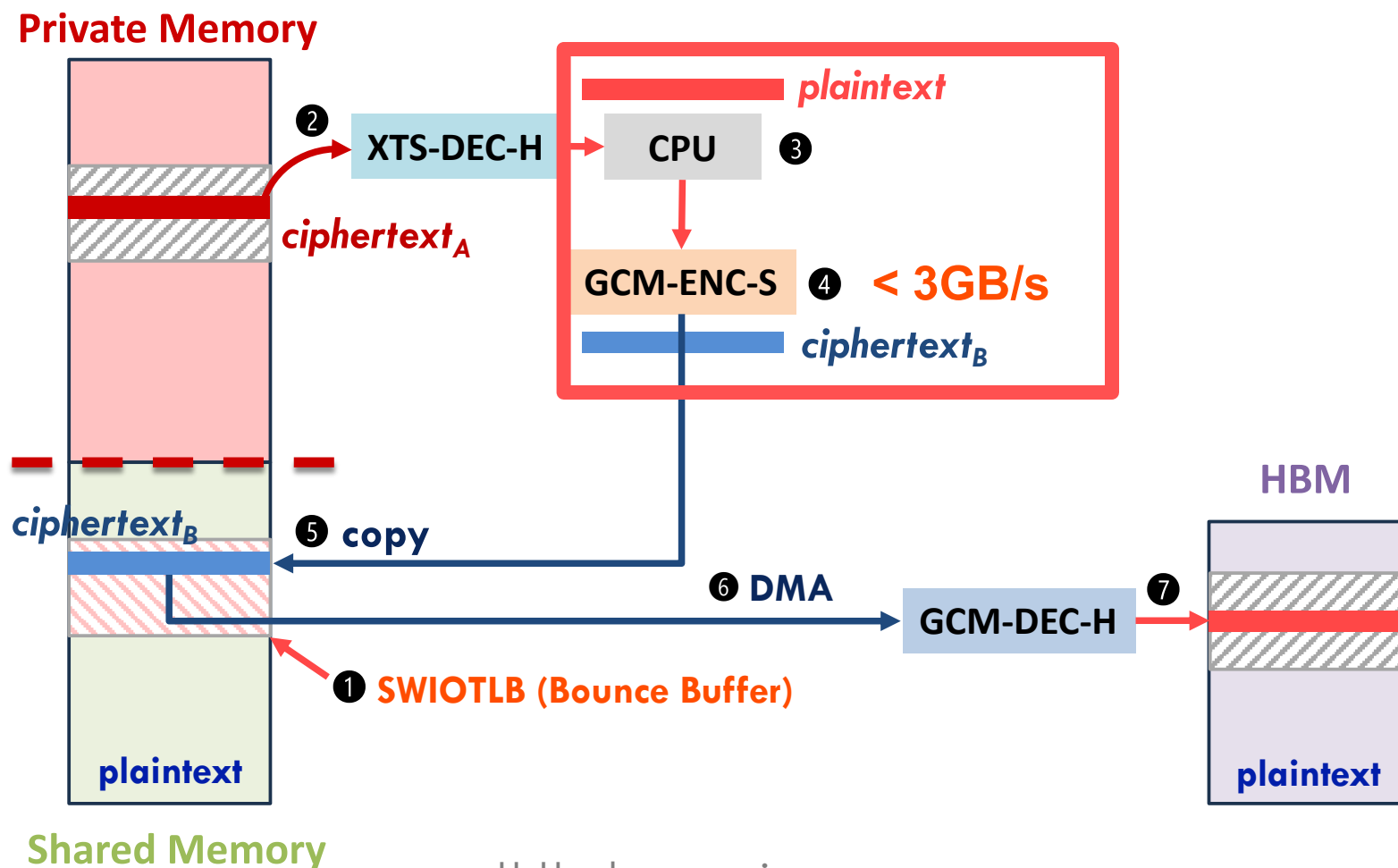
## Encryption



(b) Traditional memory encryption

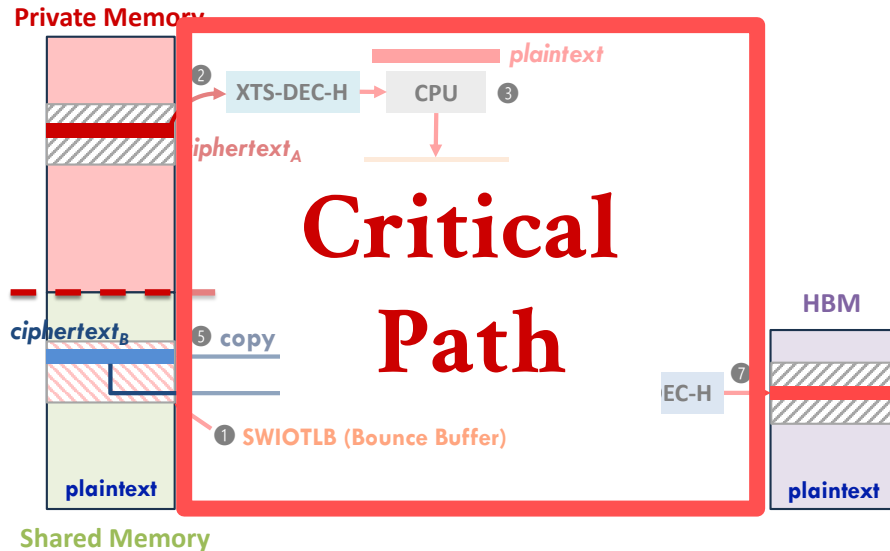
**The memory problem!**

# Problem: Bounce Buffer CC



-H: Hardware engine  
-S: Software implementation

# Problem



UVM brings frequent CPU-GPU interaction



[ISPASS'25] Low performance under CC



Lack of explicit IV  
(Encryption on the critical path)

The missing IV management makes encryption synchronized with runtime access.

# This Work: LÆGIS

## Bring flexible IV management back to GPU-based CC

Consideration	Prior Counter Mode	GPU-based CC (UVM)
IV / MAC	✓ Explicit per-CL	✗ Implicit, access-based
Integrity Tree	✗ Required	✓ Deprecated
OTP Timing	✓ Decoupled	✗ On-access only
Granularity	✗ Cache-line (64B/128B)	✓ Page-level (64KB, 2MB)
Memory	✗ Encrypted DRAM	✓ Plaintext HBM

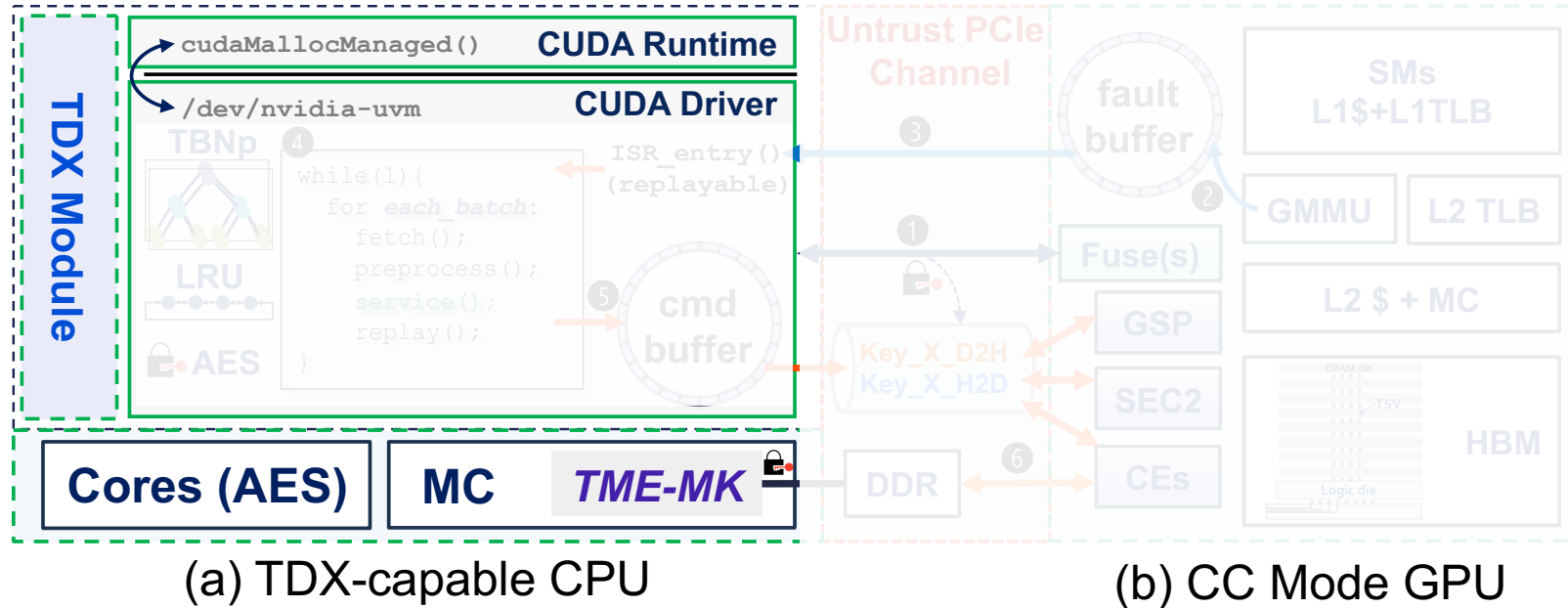
*Encryption context management is more critical.*

- ❖ Performance dissection of UVM under CC
- ❖ Observations + Opportunities
- ❖ LÆGIS: opportunistically pre-encryption + HBM for IV management

# Outline

- INTRODUCTION
- **BACKGROUND**
- UVM PERFORMANCE DISSECTION UNDER CC
- LÆGIS: OVERVIEW & IMPLEMENTATION
- EVALUATION
- TAKEAWAY

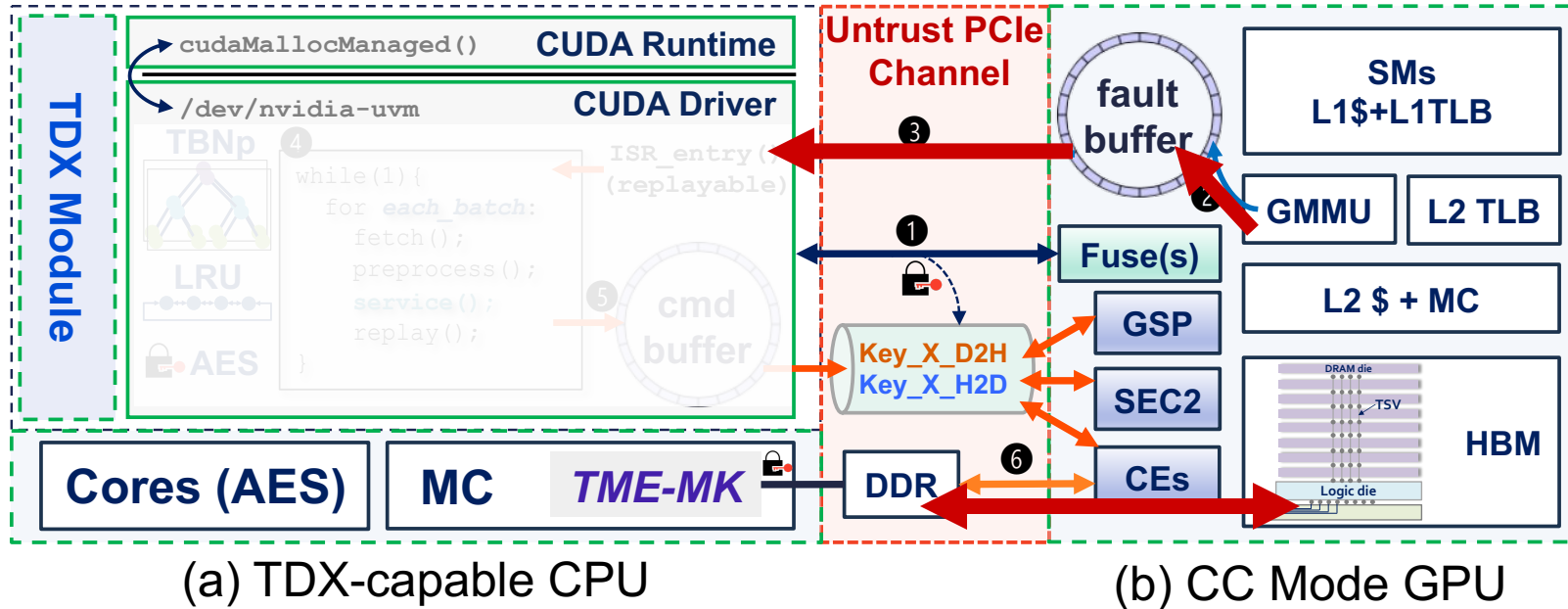
# UVM in TDX-based System



(a) TDX-capable CPU

(b) CC Mode GPU

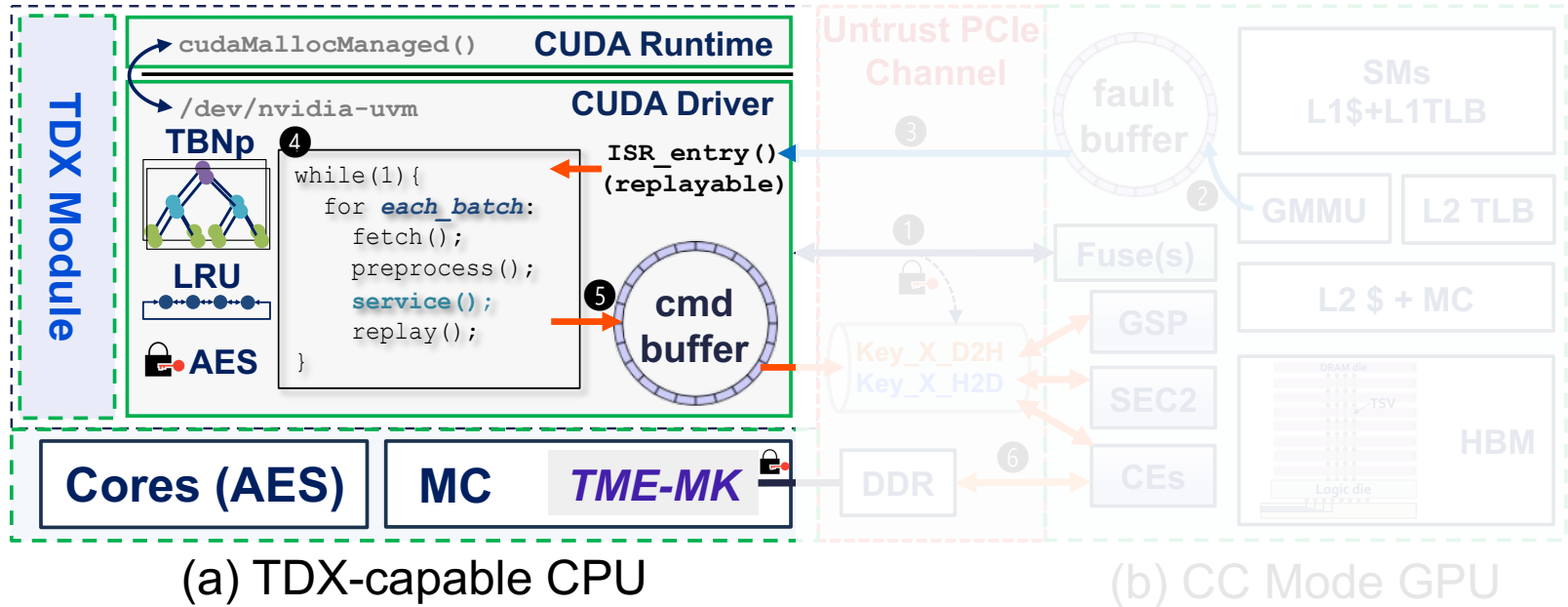
# UVM in TDX-based System



(a) TDX-capable CPU

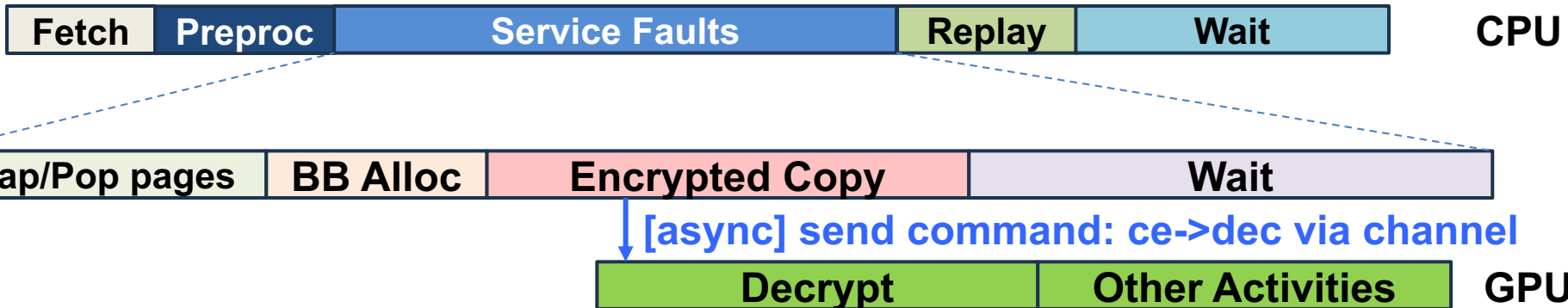
(b) CC Mode GPU

# UVM in TDX-based System

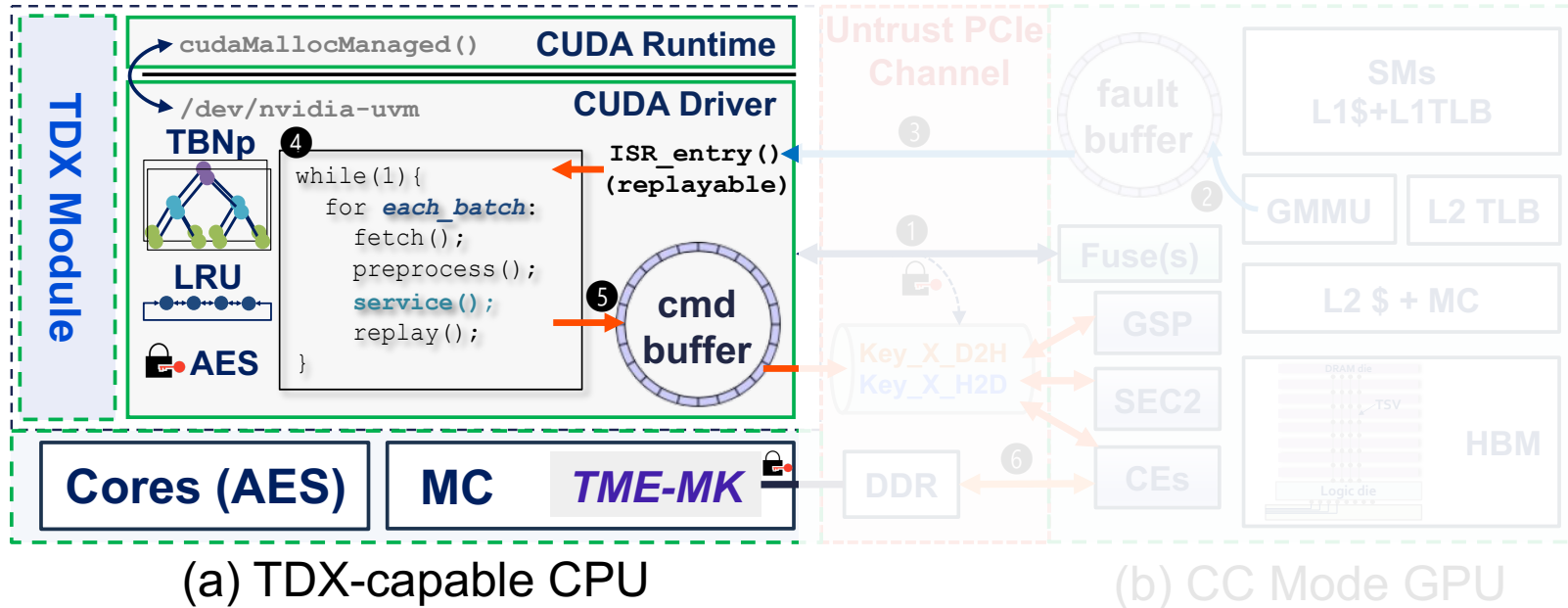


(a) TDX-capable CPU

(b) CC Mode GPU

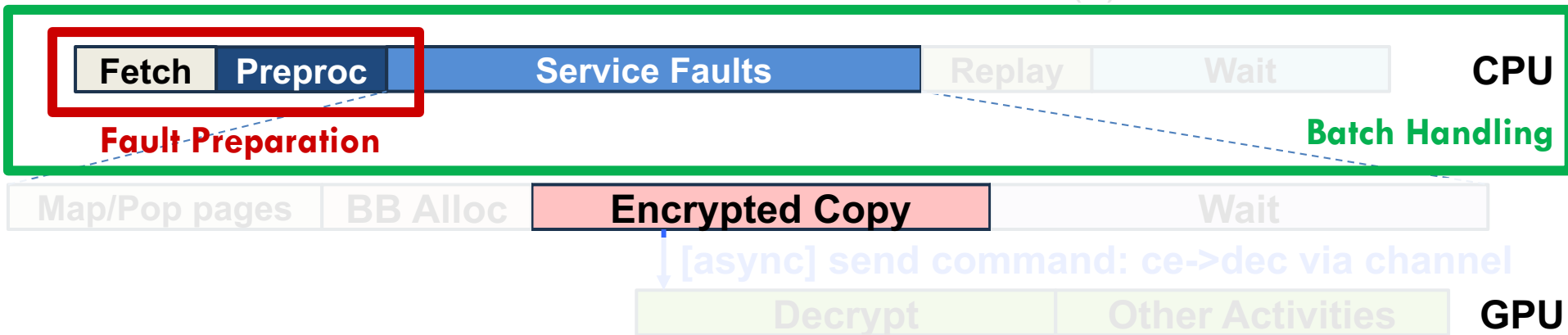


# UVM in TDX-based System

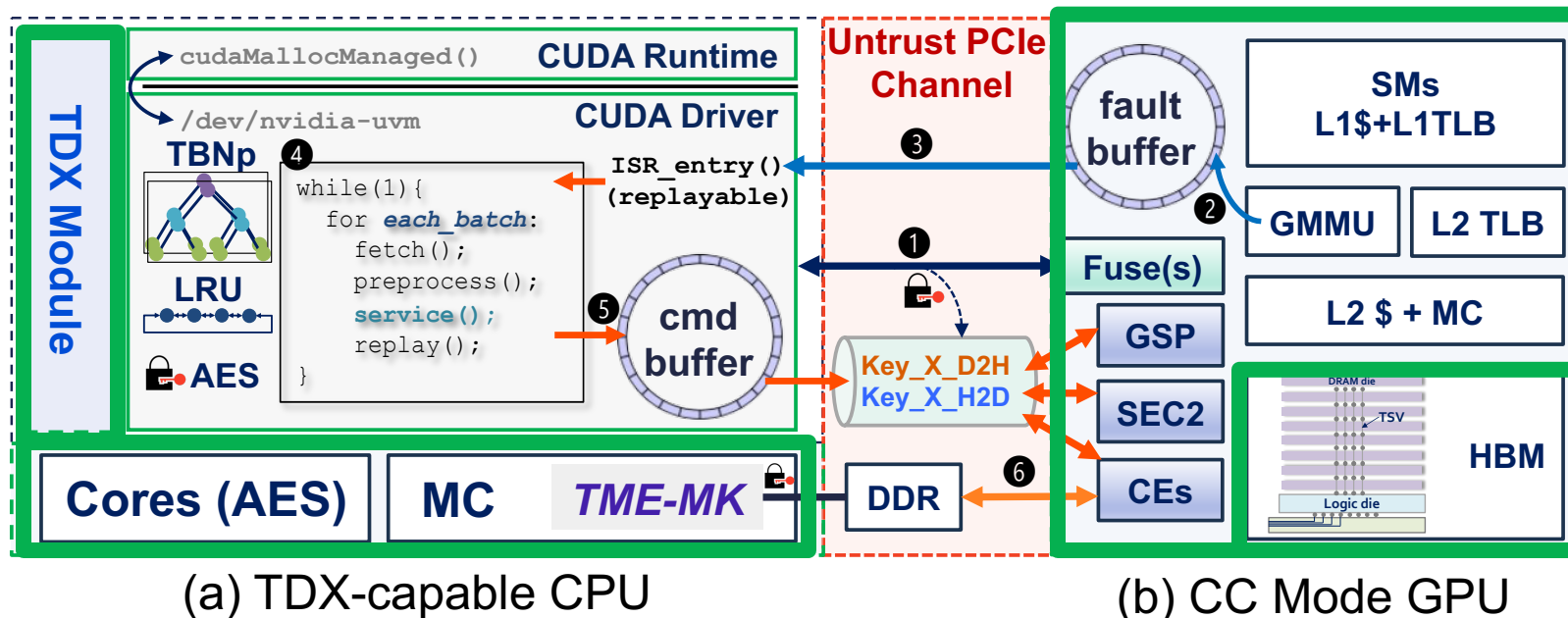


(a) TDX-capable CPU

(b) CC Mode GPU



# Threat Model



**Signed software**  
**Chip packages**  
**Internal structures (HBM)**

**Check our paper: Section IX for more discussions on threat model**

# Outline

- INTRODUCTION
- BACKGROUND
- **UVM PERFORMANCE DISSECTION UNDER CC**
- LÆGIS: OVERVIEW & IMPLEMENTATION
- EVALUATION
- TAKEAWAY

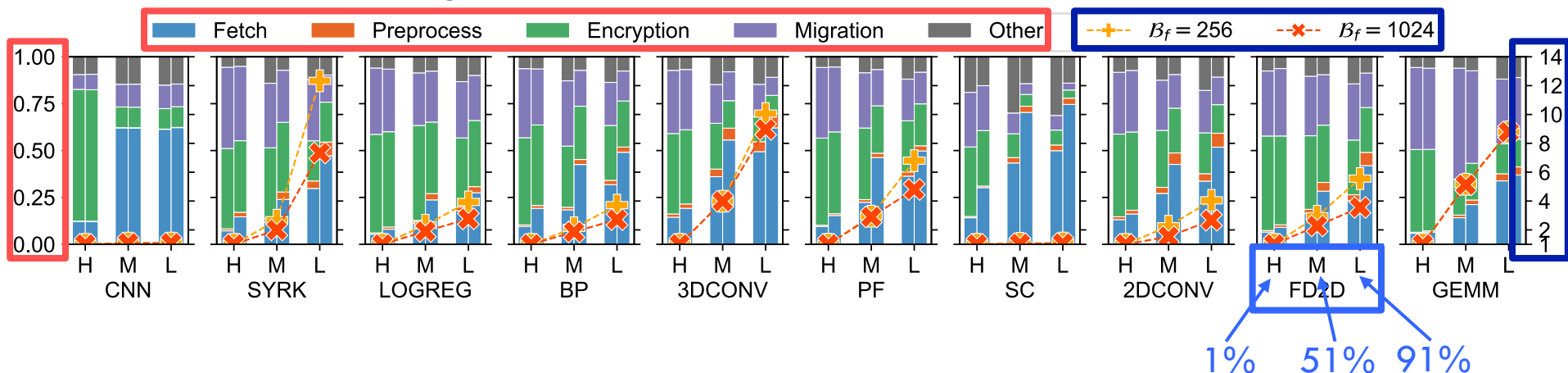
# Setup

- ❖ **5th Gen Intel Xeon 6530 Gold**
- ❖ **NVIDIA H100 NVL**
- ❖ **CUDA 12.4**
- ❖ **OpenSSL v3.0.2**
- ❖ **Linux 6.2.0-mvp10v1+8-generic**

**UVM performance evaluated using an instrumented NVIDIA driver (version 550.163.01)**

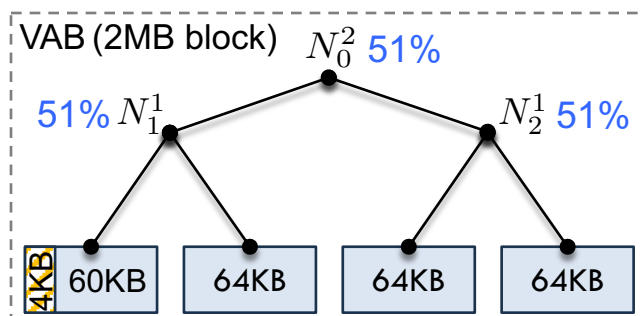
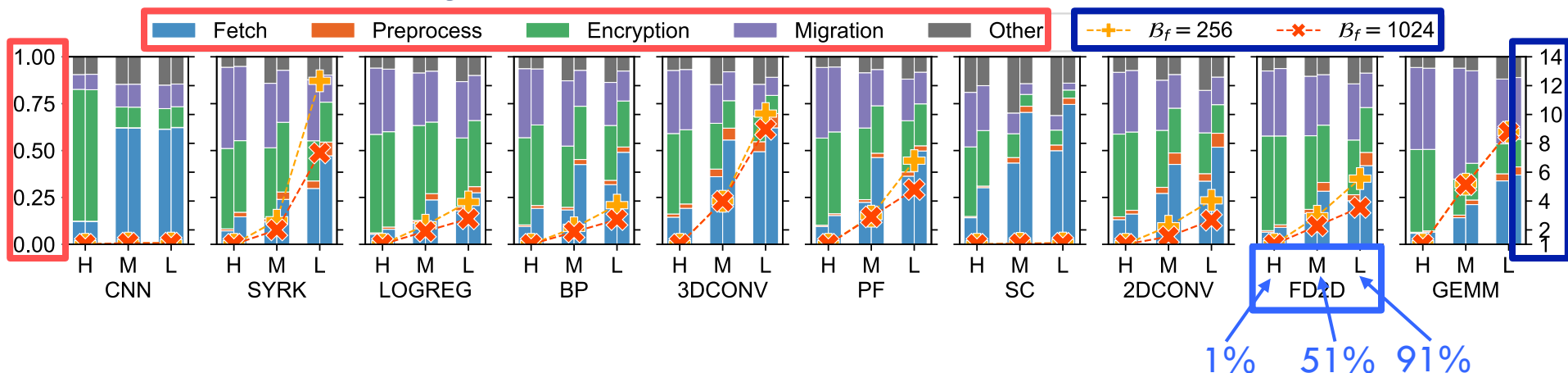
# Batch Handling Performance

❖ Evaluated using an instrumented NVIDIA driver (version 550.163.01)



# Batch Handling Performance

❖ Evaluated using an instrumented NVIDIA driver (version 550.163.01)

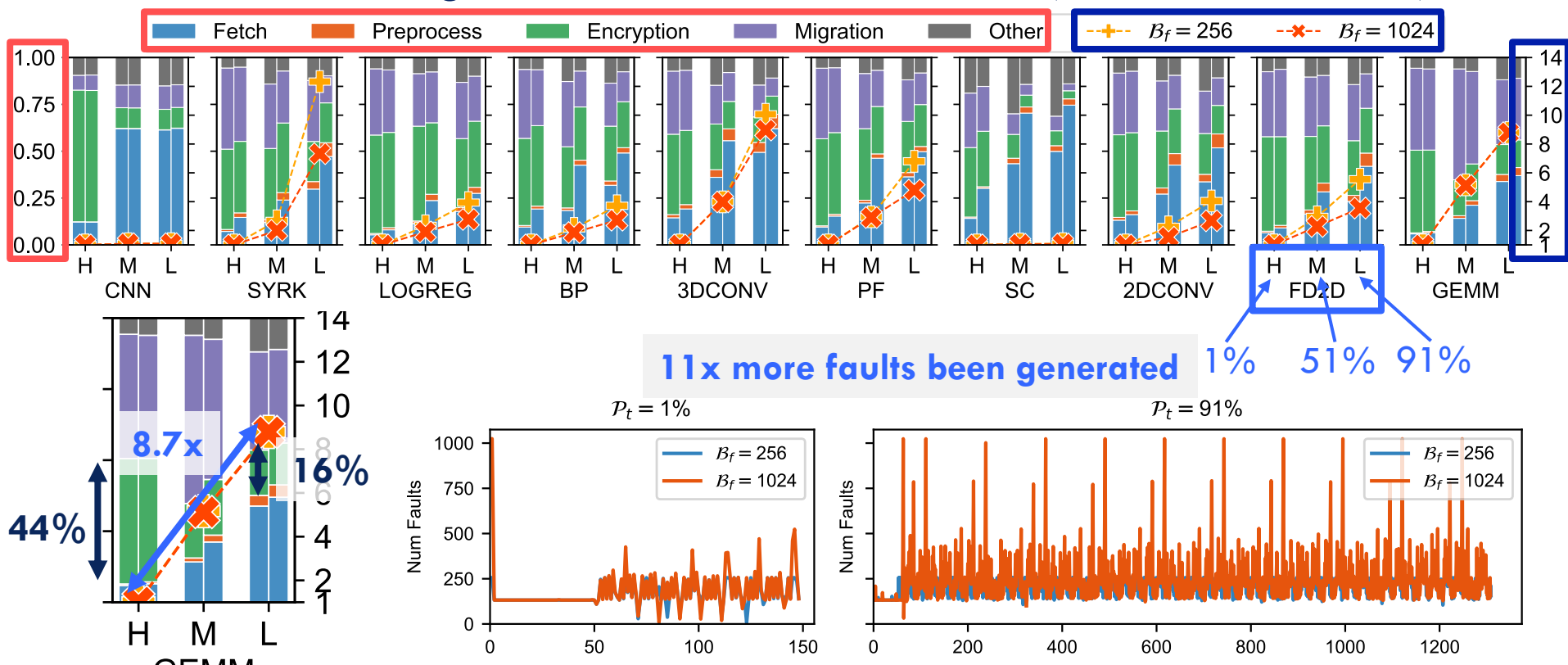


Tree-based Neighborhood Prefetcher (TBNp)

When more than 51% pages been faulted, the sub-tree will be prefetched.

# Batch Handling Performance

❖ Evaluated using an instrumented NVIDIA driver (version 550.163.01)



**Observation 1:** Aggressive prefetching ( $P_t = 1\%$ ) reduces costly GPU-CPU interactions. However, it shifts the burden to encryption, which dominates the overall handling time and can exceed 70%. Thus, encryption is the critical bottleneck.

# Encryption Mechanism

```
1 va_bb = dma_alloc_coherent();
2 for_each_va_block_page_in_region(page_index, region) {
3     void *va_src = kmap(src_page);
4     uvm_cc_cpu_encrypt(push->channel, va_bb, va_src, iv,
5     ↪ PAGE_SIZE, va_tag);
6     kunmap(src_page);
7     gpu->parent->ce->decrypt(push, gpu_dst, va_bb, PAGE_SIZE,
8     ↪ va_tag);
9 }
```

# Encryption Mechanism

@[aead\_enc,

crypto\_aead\_encrypt+1

crypto\_aead\_encrypt+57

lkca\_aead\_internal+243

libspdm\_aead\_preallocated+961

libspdm\_aead\_aes\_gcm\_encrypt\_prealloc+157

ccslEncrypt\_KERNEL+201

nvGpuOpsCcsIEncrypt+32

nvUvmInterfaceCsIEncrypt+64

uvm\_conf\_computing\_cpu\_encrypt+112

uvm\_channel\_end\_push+3012

uvm\_push\_end+14

.....

service\_fault\_batch+182

uvm\_gpu\_service\_replayable\_faults+356

.....

replayable\_faults\_isr\_bottom\_half\_entry+39

\_main\_loop+137

kthread+235

ret\_from\_fork+41

Linux Kernel Crypto APIs

libspdm/lkca

NVIDIA  
Cryptography Services Library

UVM Fault Handling  
/dev/nvidia-uvm

Kernel Threading

]: 1

# Encryption Mechanism

integrity/freshness/...? **iv\_space** [page\_src]  
 (decoupled encryption) `uvm_cc_cpu_encrypt` (page\_src, **iv++**)

```

1 va_bb = dma_alloc_coherent();
2 for_each_va_block_page_in_region(page_index, region) {
3     void *va_src = kmap(src_page);
4     uvm_cc_cpu_encrypt push->channel, va_bb, va_src, iv,
      ↪ PAGE_SIZE, va_tag);
5     kunmap(src_page);
6     gpu->parent->ce->decrypt (push, gpu_dst, va_bb, PAGE_SIZE,
      ↪ va_tag);
7 }

```

**Observation 2:** Under GPU-based CC, UVM paging encryption is a kernel-space, Linux Kernel Crypto API (LCA)-backed path that operates **synchronously** at **4KB page granularity** with IV synchronization. This differs from user-space OpenSSL-based bulk encryption.

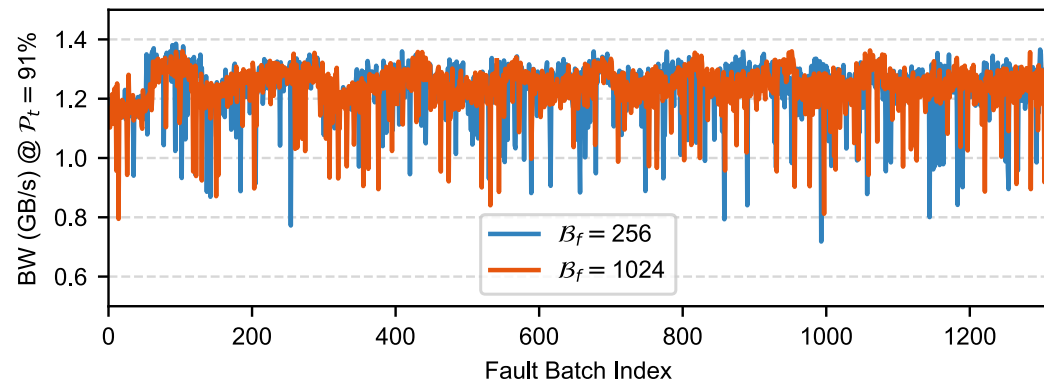
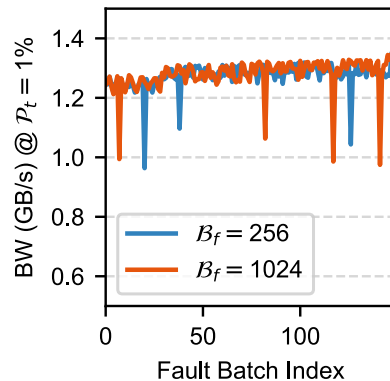
# Encryption Performance

```

1 va_bb = dma_alloc_coherent();
2 for_each_va_block_page_in_region(page_index, region) {
3     void *va_src = kmap(src_page);
4     uvm_cc_cpu_encrypt(push->channel, va_bb, va_src, iv,
5     ↪ PAGE_SIZE, va_tag);
6     kunmap(src_page);
7     gpu->parent->ce->decrypt(push, gpu_dst, va_bb, PAGE_SIZE,
8     ↪ va_tag);
9 }

```

## Encrypting a 4KB page takes 3 $\mu$ s

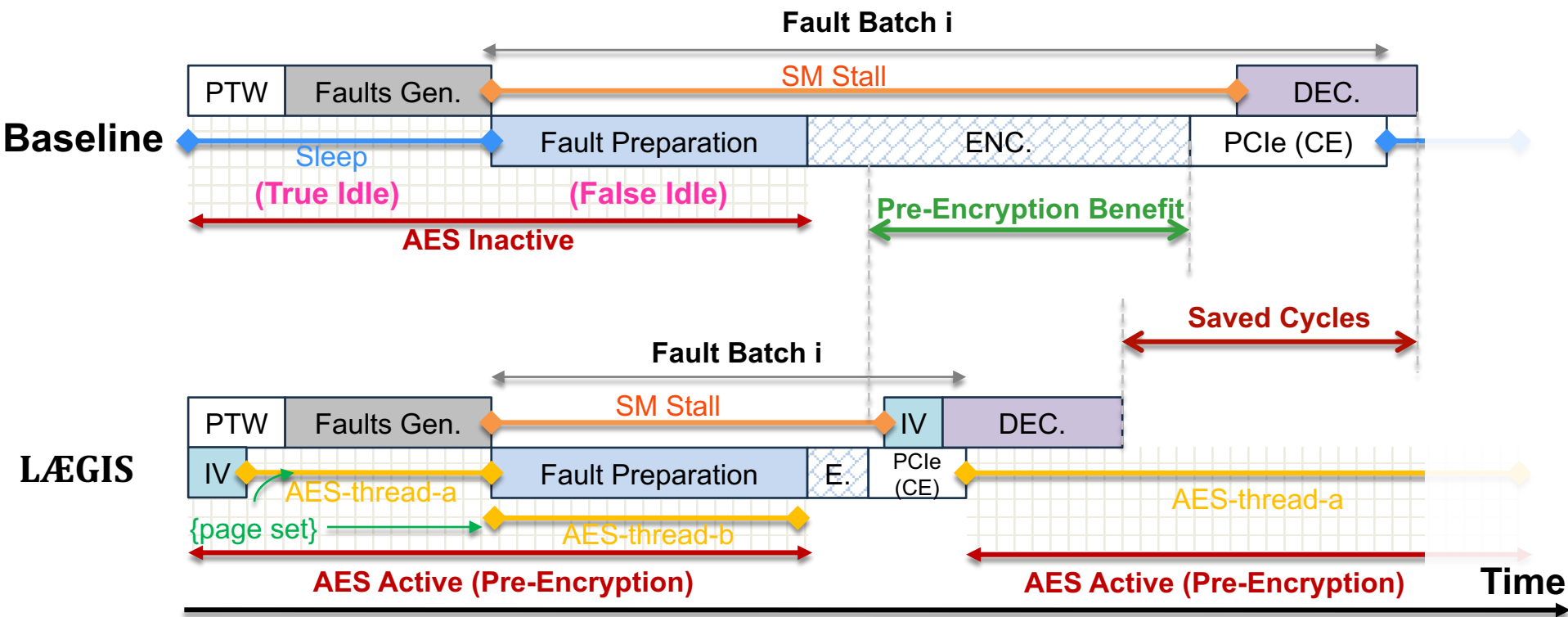


**Observation 3:** GPU CC presents slow (1.3 GB/s) and synchronous software encryption in UVM, placing it on the critical path. Such low throughput presents significant performance overhead.

# Outline

- INTRODUCTION
- BACKGROUND
- UVM PERFORMANCE DISSECTION UNDER CC
- **LÆGIS: OVERVIEW & IMPLEMENTATION**
- EVALUATION
- TAKEAWAY

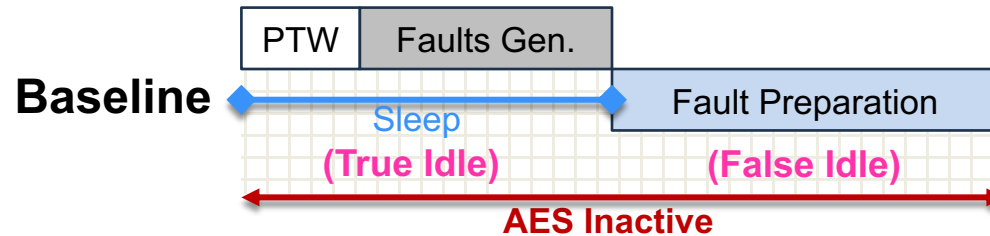
# Overview



\* Assume IV can be explicitly managed by LÆGIS

**PRE-ENCRYPTION**  
**When, How and What**

# Pre-Encryption Definition

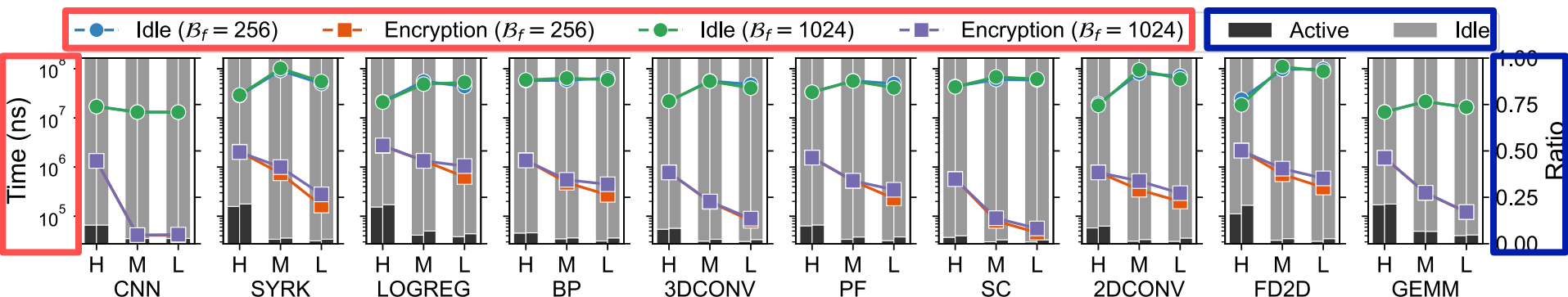


## INVARIANT

In LÆGIS, all pre-encrypted pages will be *committed* in the end.

**LÆGIS tolerates wrong prediction**

# Idle Service Routine: When to Encrypt



On average, the driver remains idle for **87%** of its execution time, during which encryption threads are underutilized.

**Opportunity 1:** Between fault batches, the service thread enters true idle, and during fault preparation the thread experiences false idle. Both leave CPU encryption resources underutilized. These intervals can be used to hide (some) encryption latency if pre-encryption is applied.

# LÆGIS Goal

**RQ1**: How to design an encryption scheme tailored for **GPU-based CC** that can support out-of-order encryption;

**RQ2**: With this scheme, how to achieve efficient batch handling.

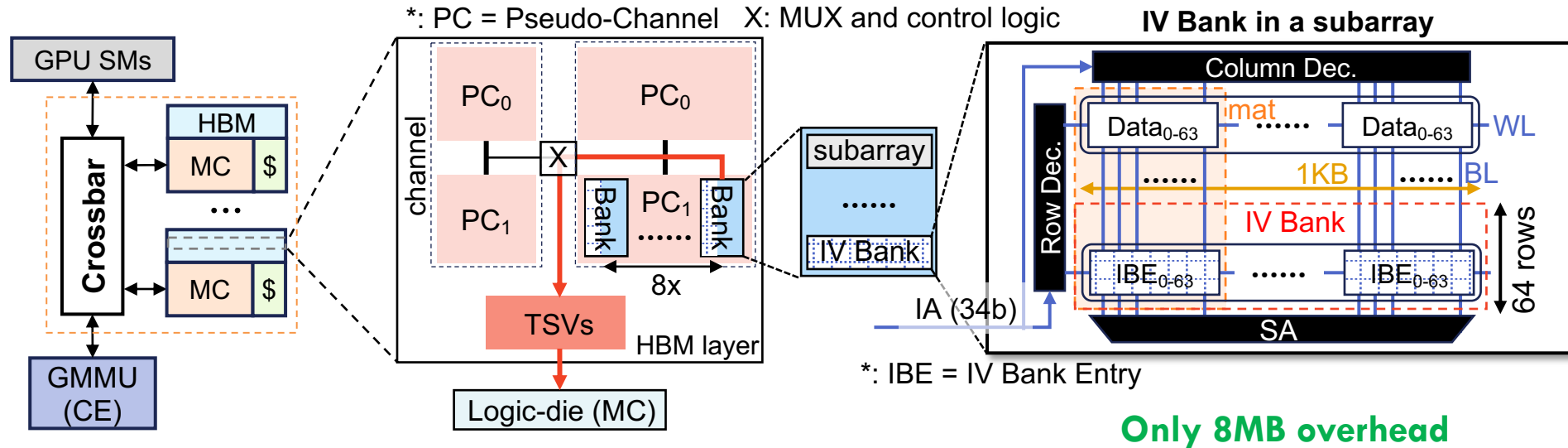
- ✓ Leverages HBM for metadata (e.g., IV) storage.
- ✓ Decouples encryption from access ordering by introducing explicit metadata management.
- ✓ Proposes an out-of-order pre-encryption scheme exploiting idle intervals.

# Outline

- INTRODUCTION
- BACKGROUND
- UVM PERFORMANCE DISSECTION UNDER CC
- **LÆGIS: OVERVIEW & IMPLEMENTATION**
- EVALUATION
- TAKEAWAY

# IV Bank Concept

Reserved logical region, termed the IV Bank, in each HBM stack



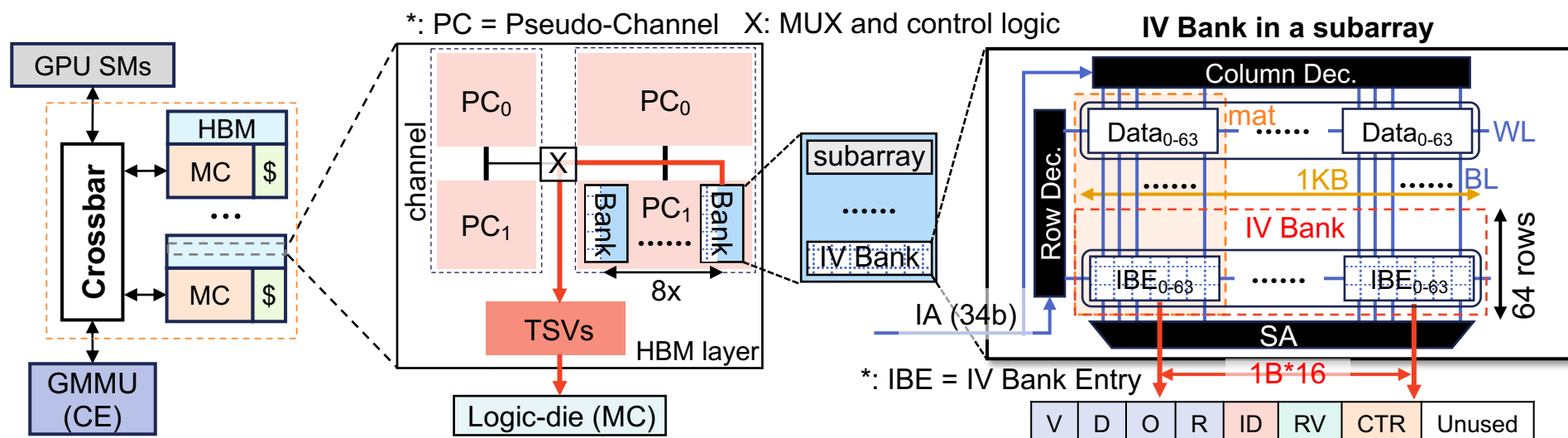
**More details please check the paper.**

The IV Bank stores per-VABlock  
IV Bank Entries (IBEs)

**IBEs are interleaved to all physical banks**

# IV Bank Concept

Reserved logical region, termed the IV Bank, in each HBM stack

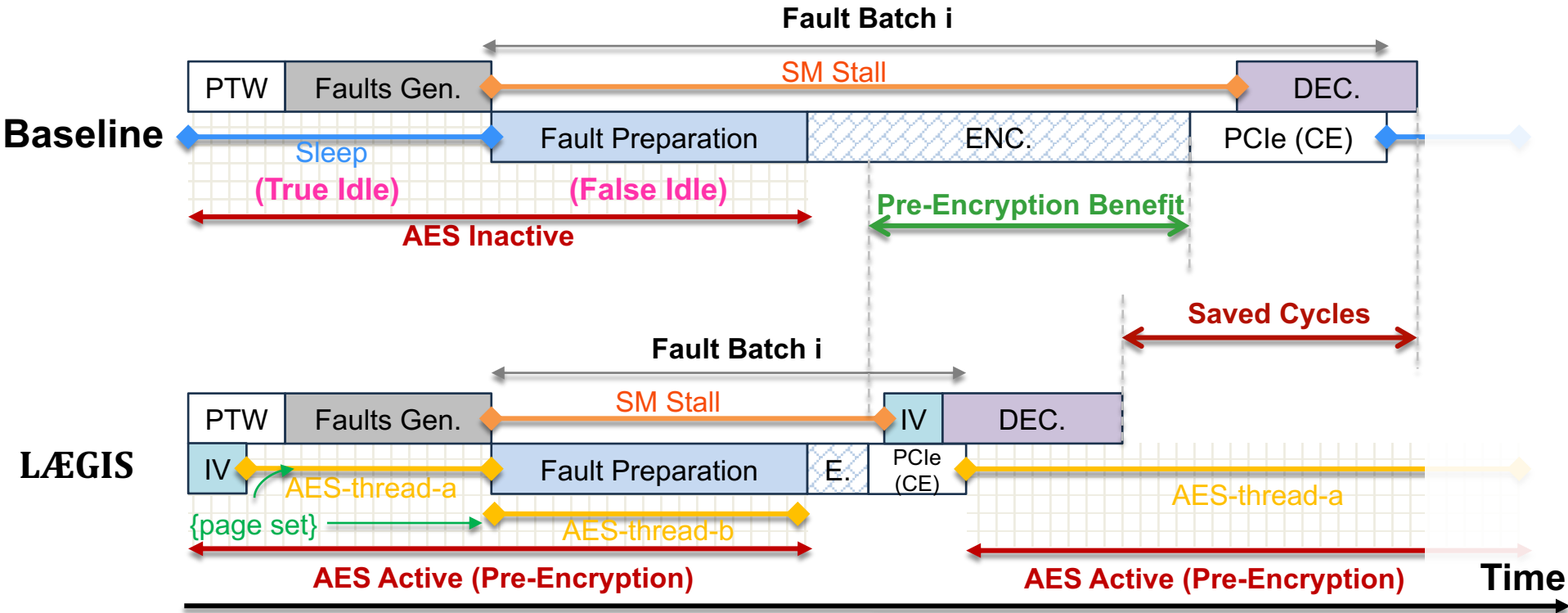


More details please  
read the paper.

To avoid IV collisions, the 96-bit IV into two parts:

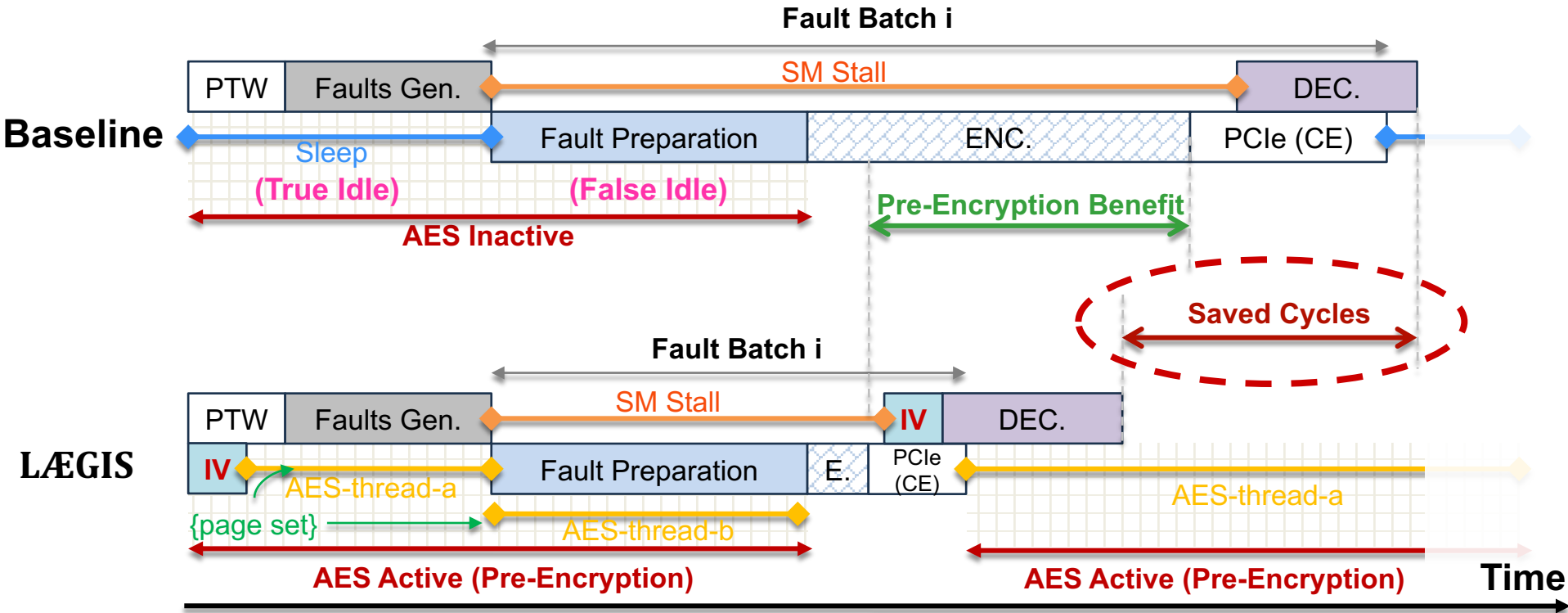
- ✓ 19-bit identifier (ID)
- ✓ 77-bit random value (RV)

# HBM Assisted Encryption



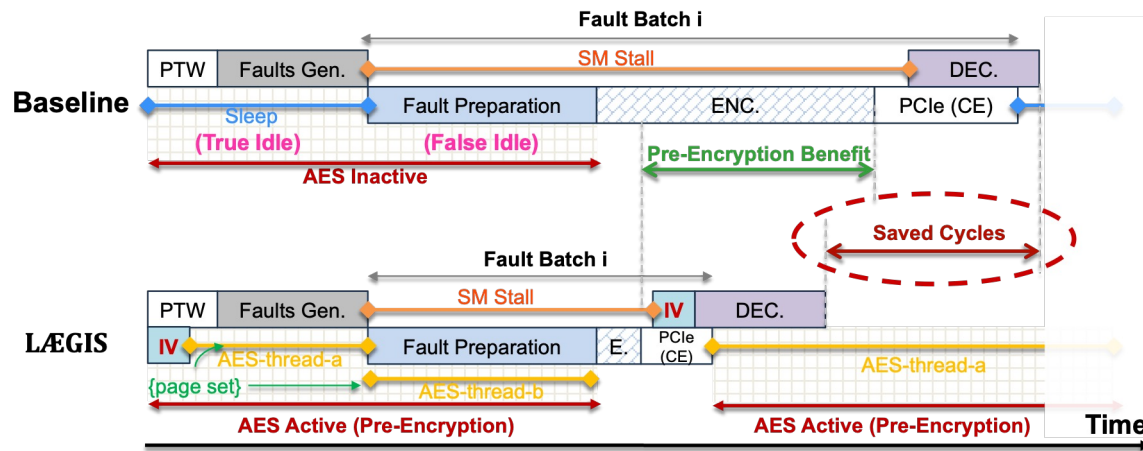
**\* Assume IV can be explicitly managed by LÆGIS**

# HBM Assisted Encryption



**\* Now IV can be explicitly managed by LÆGIS !**

# HBM Assisted Encryption



**\* Now IV can be explicitly managed by LÆGIS !**

**Full walkthrough example please read the paper.**

Opportunistically pre-encrypt pages with the help of IV Bank during batch handling.

Even random order encryption will work:  
IV Bank decouples encryption from access.

# Outline

- INTRODUCTION
- BACKGROUND
- UVM PERFORMANCE DISSECTION UNDER CC
- LÆGIS: OVERVIEW & IMPLEMENTATION
- **EVALUATION**
- TAKEAWAY

# Simulation Setup

## GPGPU-Sim v4.2 + UVMSmart



80 SMs, 64 warps/SM; shared GMMU; private TLB

8GB HBM2 (1 stack);

Page size: 4 KB base, 2 MB big;

PCIe: 64 GB/s per direction;

Pt: 51% (default), 1% (aggressive);

Thread switch: 2  $\mu$ s

Fault preparation and idle time: profiled from hardware setup

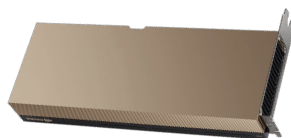


8GB 4-Hi stack, 8 channels, 16 pseudo channels

FR-FCFS; 16 banks, 4 bank groups/channel; 128-bit interface

256 GB/s via 8 channels @1 GHz;

Static mapping



5th Gen Intel Xeon 6530 Gold

NVIDIA H100 NVL

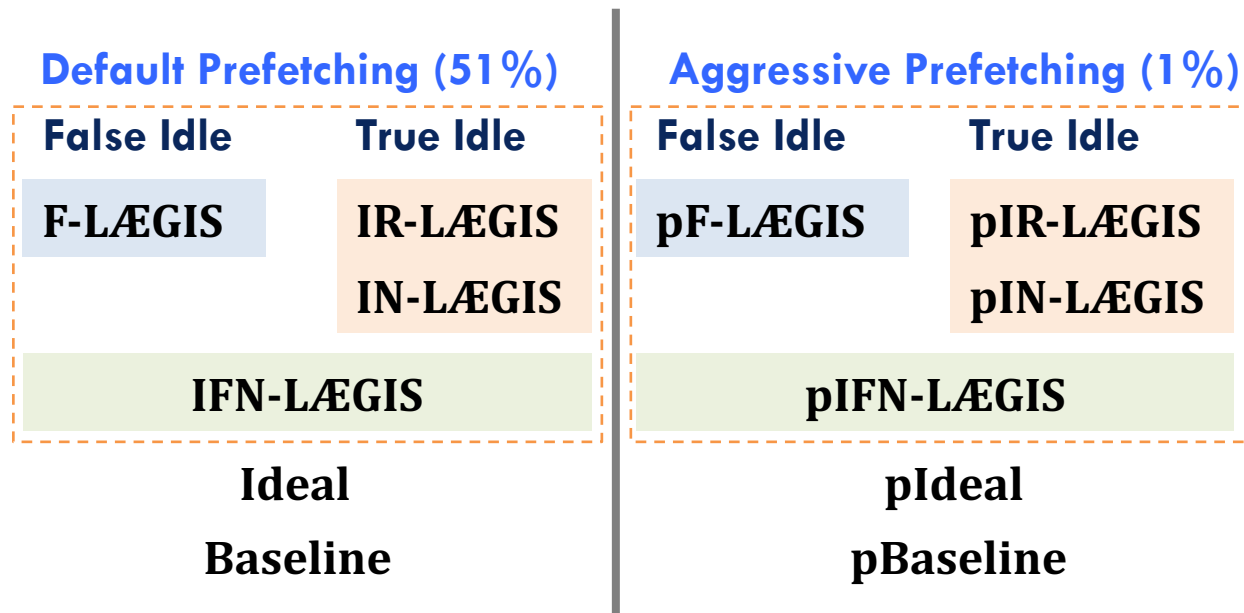
CUDA 12.4

OpenSSL v3.0.2

Linux 6.2.0-mvp10v1+8-generic

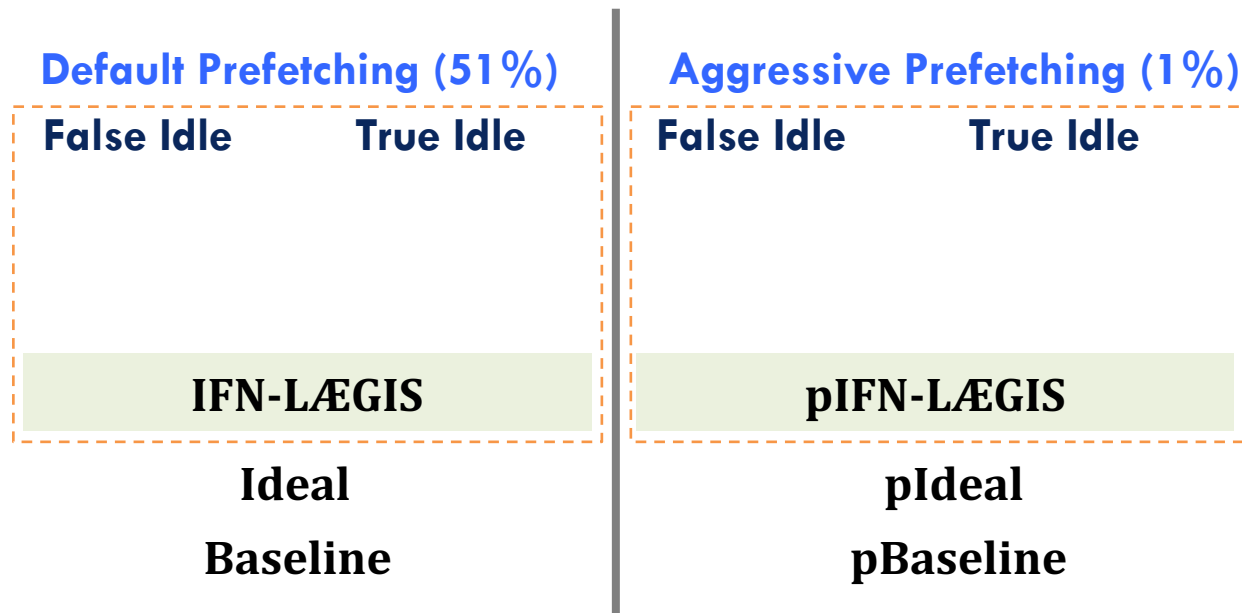
# Method

16 applications from widely adopted benchmark suites



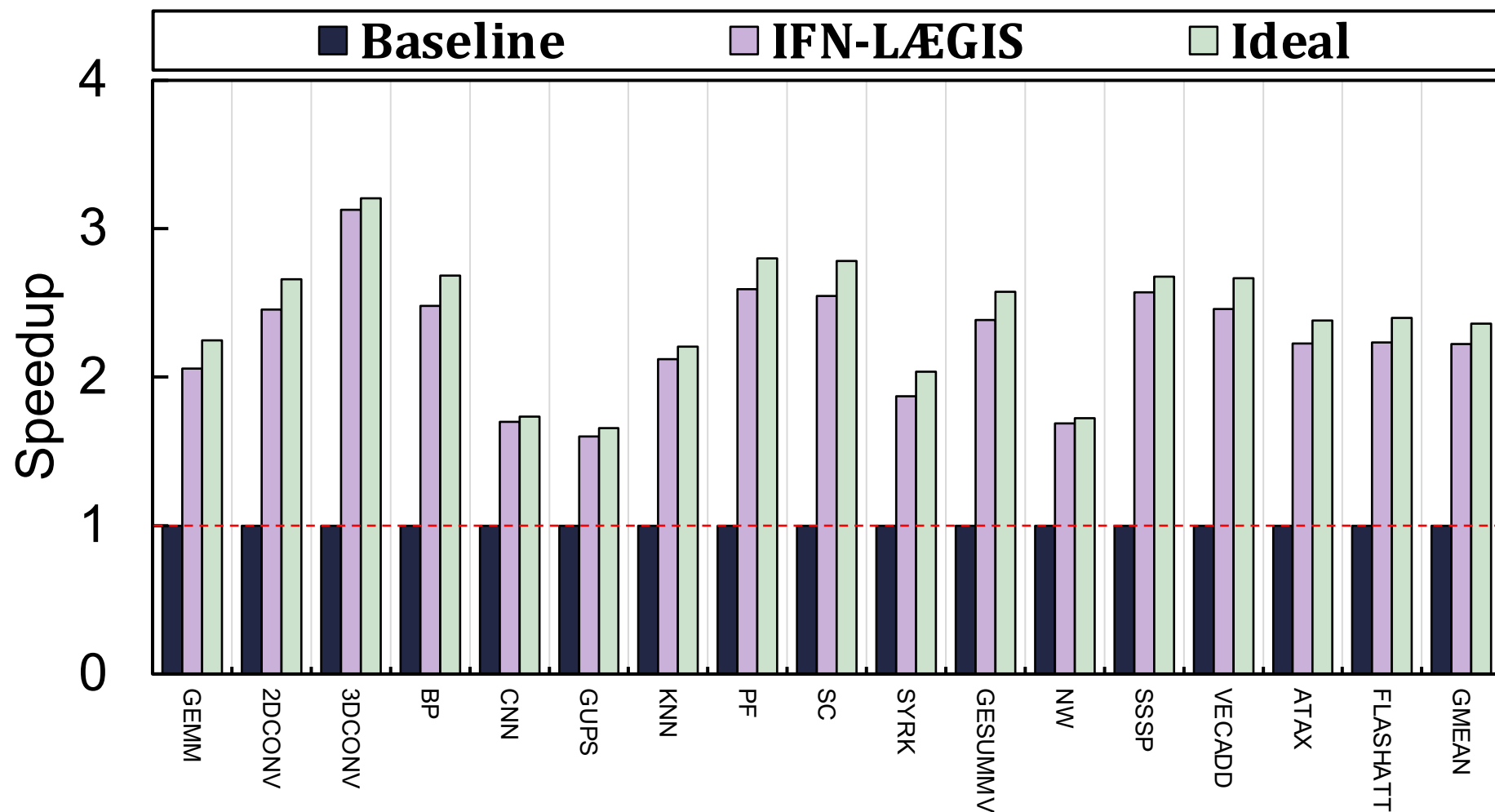
# Method

16 applications from widely adopted benchmark suites

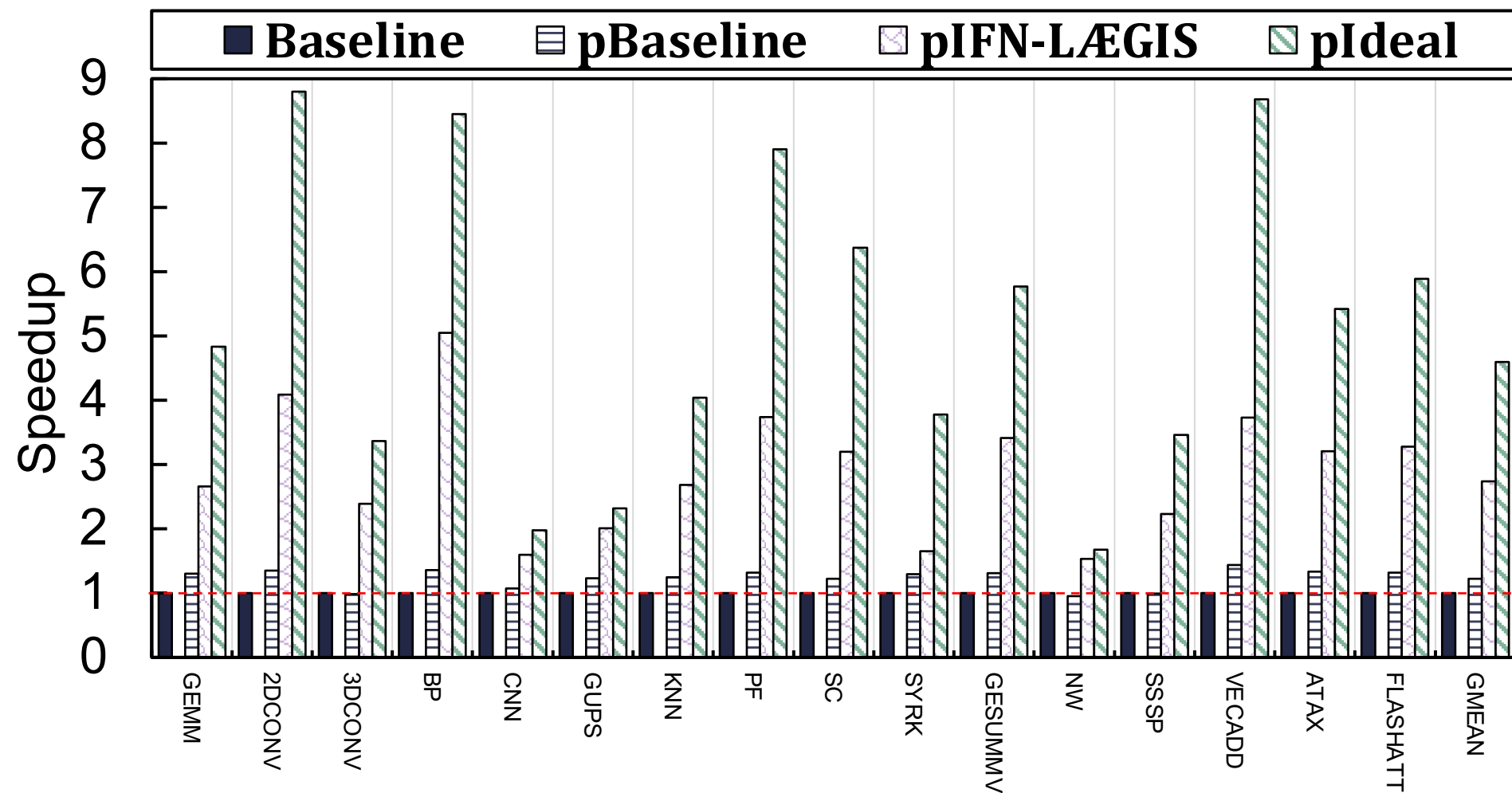


**More results please check the paper.**

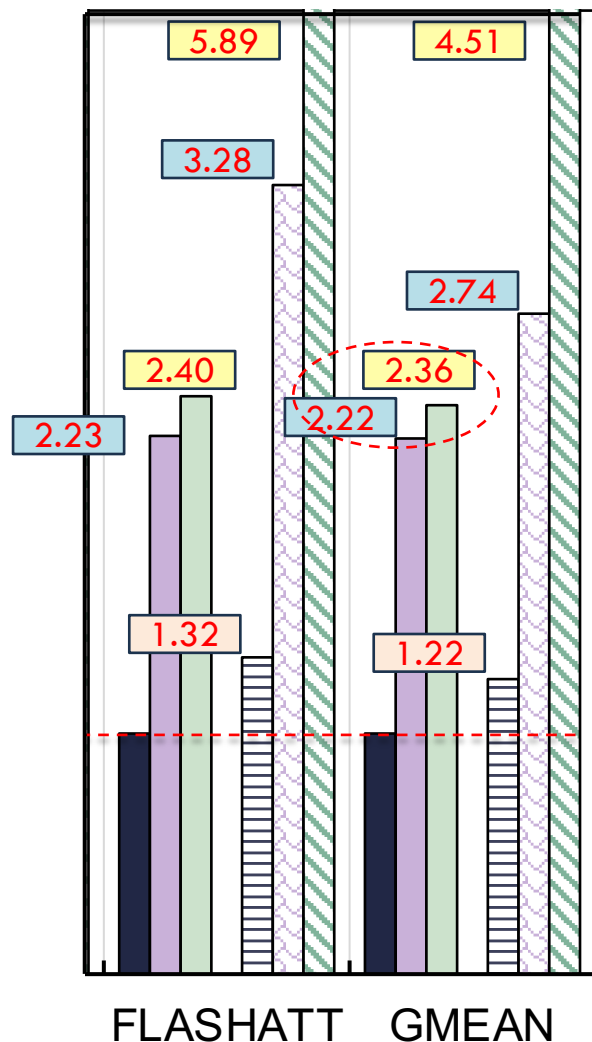
# Main Results



# Main Results



# Main Results



On average, **IFN-LÆGIS** achieves a **2.22×** speedup, with a maximum of **3.13×**, reducing the gap to Ideal to only **5.8%**.

**pIFN-LÆGIS** further improves performance under aggressive prefetching. It achieves a maximum speedup of **5.05×** and an average of **2.74×**, narrowing the gap with **pIdeal** to 40.34%.

# More Details

- Threat model and HBM security discussion
- Security analysis
- Pre-Encryption: how and what
- IV Bank details
- Pre-Encryption walkthrough example
- More LÆGIS variants results
- Future directions
- Scalability
- Related works
- More baselines for CPU encryption:
  - Multi-threading
  - Advanced hardware
  - OpenSSL versions

.....

# Outline

- INTRODUCTION
- BACKGROUND
- UVM PERFORMANCE DISSECTION UNDER CC
- LÆGIS: OVERVIEW & IMPLEMENTATION
- EVALUATION
- TAKEAWAY

# Take Away Messages

- ❖ We identify three inefficiencies in GPU-based CC when it comes to UVM.
  - CPU–GPU IV synchronization
  - Driver idleness
  - Slow software encryption
- ❖ We propose **LÆGIS**, a design that opportunistically pre-encrypts pages in CC and uses secure HBM for IV management.
- ❖ **LÆGIS** introduces an in-memory IV Bank in 3D-stacked HBM, decoupling encryption from CPU–GPU synchronization and removing the need for integrity trees.

**Average speedup of 2.22× (up to 3.13×) and 2.74× (up to 5.05×) under default and aggressive prefetching.**

[PDF] [https://adwaitjog.github.io/docs/pdf/yang\\_laegis\\_isca26.pdf](https://adwaitjog.github.io/docs/pdf/yang_laegis_isca26.pdf)  
[GitHub] <https://github.com/insight-cal-uva/laegis-isca26-artifact>

# Thank You!

# Questions?

**LÆGIS: Pinpointing and Addressing  
Performance Overheads of  
GPU-based Confidential Computing**

**Yang Yang, Adwait Jog**  
{*yangyang, ajog*} @ *virginia.edu*



ISCA 2026