### Dissecting Performance Overheads of Confidential Computing in GPU-based Systems

### Yang Yang, Mohammad Sonji, Adwait Jog

{yangyang, npv2tk, ajog} @virginia.edu

**Insight Computer Architecture Lab** 



# **Executive Summary**

### **Problem:**

- Confidential Computing with GPU is emerging
- Existing work lacks detailed performance dissection
- Optimization insights remain unclear due to limited characterization

**Goal:** Dissect the **performance overheads** of GPU-based confidential computing that can motivate various optimizations to address them.

### Our results on real hardware show that GPU-based CC comes with:

- Data Movement Overhead (up to 19.7×)
- Low Encryption Throughput (3.36 GB/s)
- Kernel Execution Time (up to 164,030.65× with UVM)
- Kernel Launch Overhead (up to 5.31×)
- **Queueing** (1.43× and 2.23× on average)
- Fusion, overlapping and quantization



#### https://github.com/insight-cal-uva/hcc-ispass25-artifact

### **Graphics Processing Units (GPUs)**





Credit/Source: AMD, NVIDIA, Wiki

### **High Throughput + Energy Efficiency**

Introduction

Background

Performance Model

**Evaluation** 

# **Needs for Security and Privacy**

### Data is Private

Client privacy (chat bot, ...)
Regulated (GDPR, PII, ...)





### Data is Massive

Outsource compute & data to public cloud Infra.

### **Can cloud service providers be trusted?**

Introduction

### **Trusted Computing**



# **Confidential Computing**



6



MEE: Memory Encryption Engine

	Introduction	Background	Performance Model	Evaluation	Conclusion
--	--------------	------------	-------------------	------------	------------

# **Confidential Computing**



7



MEE: Memory Encryption Engine

Introduction	Background	Performance Model	Evaluation	Conclusion
muouuenom	Duchstound	I CHOI Munee Mouel		Gonerasion

8

# **Confidential Computing**



MEE: Memory Encryption Engine

Introduction	Background	Performance Model	Evaluation	Conclusion
	<u>C</u>			







Communication (MMIO/DMA) will trigger #VE, TD will context-switch



# **Real World Setup**



https://github.com/insight-cal-uva/hcc-ispass25-artifact

System Configuration

Configuration	Details
CPU	2× 5th Gen Intel Xeon 6530 Gold @2.1GHz, 32 cores
TME-MK	Auto bypass enabled
OS	Ubuntu 22.04.5 LTS (Linux 6.2.0, tdx patched)
Hypervisor	QEMU 7.2.0 (tdx patched)
TDX Tools	TDX 1.5 (tag 2023ww15)
GPU	NVIDIA H100 NVL, 94GB HBM3, PCIe 5.0 ×16 CUDA 12.4, Driver 550.127.05

### **Problem:**

- Confidential Computing with GPU is emerging
- Existing work lacks detailed performance dissection
- **Optimization insights** remain unclear due to limited characterization

Performance Mode

### **Performance Model**

ALLOC

### **Performance Model**



### **Performance Model**









D2H: device (GPU) to host (CPU) copy H2D: host (CPU) to device (GPU) copy

Introduction



D2H: device (GPU) to host (CPU) copy H2D: host (CPU) to device (GPU) copy

Introduction





$$\boldsymbol{P} = \mathbf{T}_{\text{mem}} + \Sigma(\text{KLO} + \text{LQT}) + \Sigma[(1 - \beta_i)(\text{KET} + \text{KQT})] + T_{\text{other}}$$

Time on memory copy

D2H: device (GPU) to host (CPU) copy H2D: host (CPU) to device (GPU) copy

🔀 Kernel Queueing Time (KQT)	💹 Launch Queueing Time (LQT)
Kernel Launch Overhead (KLO)	Launch Execution Time (KET)



$$P = T_{\text{mem}} + \sum(\text{KLO} + LQT) + \sum[(1 - \beta_i)(\text{KET} + KQT)] + T_{\text{other}}$$

D2H: device (GPU) to host (CPU) copy H2D: host (CPU) to device (GPU) copy

Introduction



$$P = T_{\text{mem}} + \Sigma(\text{KLO} + LQT) + \Sigma[(1 - \beta_i)(\text{KET} + KQT)] + T_{\text{other}}$$

D2H: device (GPU) to host (CPU) copy H2D: host (CPU) to device (GPU) copy

Introduction



$$P = T_{\text{mem}} + \Sigma(\text{KLO}+\text{LQT}) + \Sigma[(1-\beta_i)(\text{KET}+\text{KQT})] + T_{\text{other}}$$

D2H: device (GPU) to host (CPU) copy H2D: host (CPU) to device (GPU) copy

Kernel Queueing Time (KQT)
Kernel Launch Overhead (KLO)
Launch Execution Time (KET)

$$P = T_{\text{mem}} + \Sigma(\text{KLO}+\text{LQT}) + \Sigma[(1-\beta_i)(\text{KET}+\text{KQT})] + T_{\text{other}}$$

D2H: device (GPU) to host (CPU) copy H2D: host (CPU) to device (GPU) copy

Introduction

Kernel Queueing Time (KQT)
Kernel Launch Overhead (KLO)
Launch Execution Time (KET)

$$P = T_{\text{mem}} + \Sigma(\text{KLO} + LQT) + \Sigma[(1 - \beta_i)(\text{KET} + KQT)] + T_{\text{other}}$$

### **Focus of this presentation**

- Data movement (H2D, D2H, page migration)
- Encryption
- Kernel Execution
- Kernel Launch
- Queuing

### **Transfer Bandwidth**



#### **Observation 1.**

PCIe bandwidth utilization in CC mode drops significantly compared to non-CC.

### **Transfer Bandwidth**



#### **Observation 1.**

- ✓ PCIe bandwidth utilization in CC mode **drops significantly** compared to non-CC.
- ✓ Bandwidth gap between pageable and pinned memory observed in non-CC mode disappears in CC mode, suggesting that pinned memory relies on pageable mechanisms in CC mode.

Introduction

Performance Model

**Evaluation** 

29

# **Crypto Throughput**



#### **Observation 2.**

- ✓ The absence of **dedicated hardware AES engines** results in low encryption throughput, even when using **AES-NI** acceleration.
- ✓ While *alternative cryptographic algorithms* may offer higher throughput, they often come at the cost of weaker security guarantees.

AES is on the critical path.

	Introduction	Background	Performance Model	Evaluation	Conclusion
--	--------------	------------	-------------------	------------	------------



Fig.3 Normalized kernel execution time.

**Observation 3.** 

CC has minimal impact on non-UVM kernels (0.48% increase).

**?** Execution is locked inside GPU, no interaction with CPU

Introduction

**Performance Model** 

**Evaluation** 



Fig.3 Normalized kernel execution time.





#### Source: NVIDIA

**Frequent CPU-GPU interaction!** 

Introduction

32

### **Kernel Execution Time**



#### Fig.3 Normalized kernel execution time.



#### **Frequent CPU-GPU interaction!**





Fig.3 Normalized kernel execution time.



Fig.3 Normalized kernel execution time.

Introduction



Fig.3 Normalized kernel execution time.

**Observation 3.** 

CC has minimal impact on non-UVM kernels (0.48% increase).

**UVM** in CC mode incurs an average slowdown of **188.87**×.

Introduction

Performance Model

### **Memory Transfer Time**



Fig.4 Time ( $\mu$ s) spent on memory copy.

#### **Observation 4.**

- ✓ On average, copy operations in CC mode take 5.80× longer compared to non-CC mode, with a maximum slowdown of 19.69×.
- ✓ Pinned memory is converted to UVM encrypted paging in CC mode which incurs high overhead.

Performance Mode

# Launch and Queuing

Results are normalized to non-CC time.



Introduction

37

### **Closer Look**



Fig.6 Simplified call stack of a cudaLaunchKernel call, obtained via perf.

	Observation 5.
✓	On average, CC increases KLO by $1.42 \times$ mostly due to TDX <i>hypercalls</i> . It increases LQT by $1.43 \times$ and KQT by $2.32 \times$ .
✓	For applications with a low number of kernel launches, KQT can be significantly amplified.

Performance Model

### **Case Study**



Fig.7 Distribution of events and their durations for representative applications, in  $\mu$ s.

### **Case Study**



Fig.7 Distribution of events and their durations for representative applications, in  $\mu$ s.

**40** 

### **Case Study**



Fig.7 Distribution of events and their durations for representative applications, in  $\mu$ s.

How to optimize? Insights: Fusion and overlapping!

41

### **Insight: Fusion**



Aggressive fusion is not working, a right fusion parameter helps.

### **More Details**

- Memory Management
- Kernel-to-Launch Ratio
- Overlapping
- CNNs
- LLMs
- Quantization
- ..



### More details in the paper!

# Conclusion

We present a **comprehensive performance evaluation** of GPU-based CC guided by a simple performance model, and considered several **optimizations** towards addressing the overheads of CC.

#### Our results on real hardware show that GPU-based CC comes with:



# Thank You! Questions?

Dissecting Performance Overheads of Confidential Computing in GPU-based Systems

Yang Yang, Mohammad Sonji, Adwait Jog

{yangyang, npv2tk, ajog} @virginia.edu



