# Graviton
## Trusted Execution Environments on GPUs

Stavros Volos,[†] Kapil Vaswani,[†] and Rodrigo Bruno[‡]

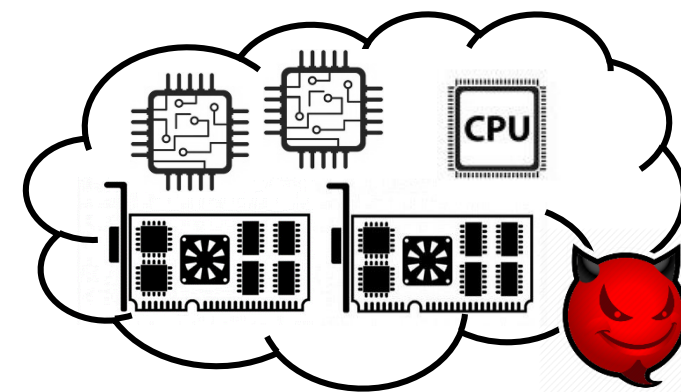[†]Microsoft Research          [‡]University of Lisbon

# Trends in Cloud Computing

Accelerators play pivotal role in cloud
- CPUs running out of steam due to *End of Moore's Law*
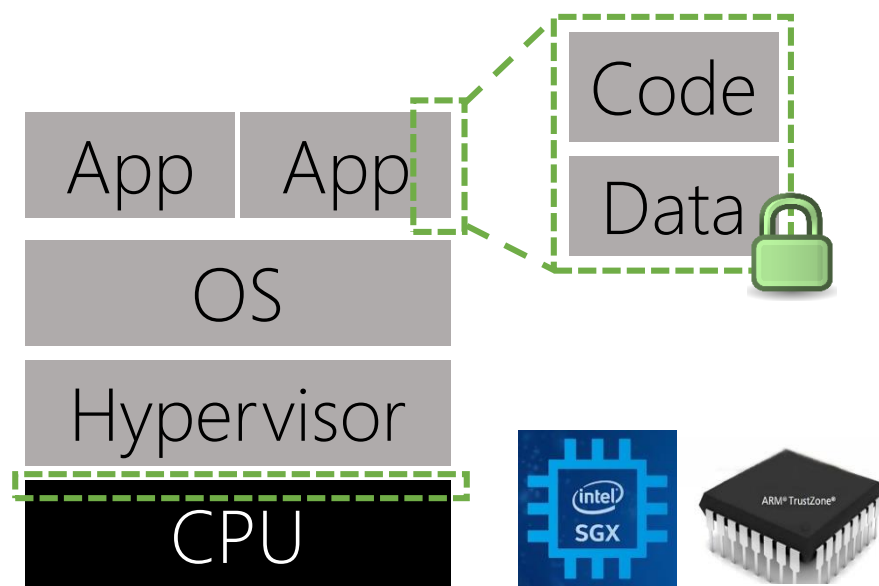- GPUs, FPGAs, custom silicon deliver 10-100x higher performance

Cloud privacy important but challenging
- Customers operate on sensitive data (e.g., patients, transactions)
- Increasing frequency and sophistication of data breaches

*Need strong security mechanisms for preserving data privacy in cloud*

# Confidential Cloud Computing



Trusted Execution Environments (TEE)

- Execution isolated from privileged attackers
- Remote attestation for establishing trust
- Examples: Intel SGX, ARM TrustZone
- Supported by major cloud providers (e.g. Azure Confidential Computing)

But, CPU TEEs cannot be used in apps that utilize accelerators

*Undesirable trade-off between performance and security*

# Our Proposal: Graviton

Graviton: Trusted Execution Environments on GPUs

- Execution isolated from system software and other co-tenants
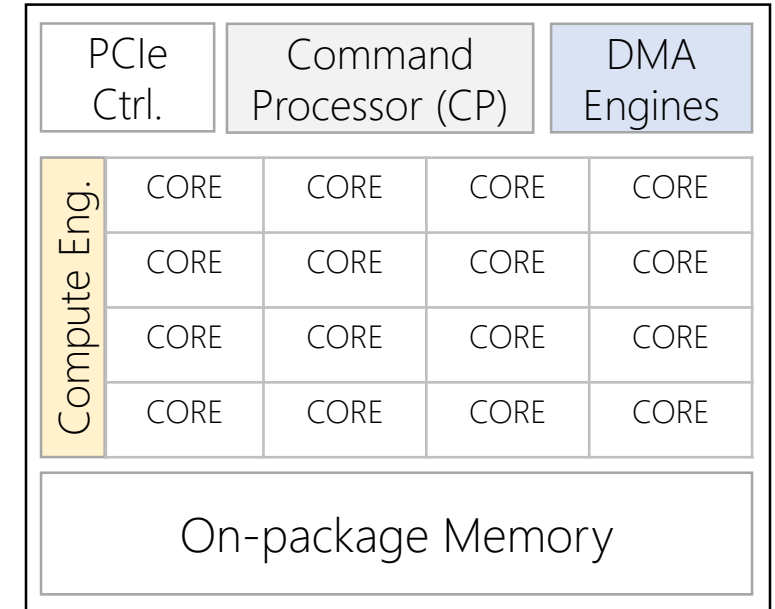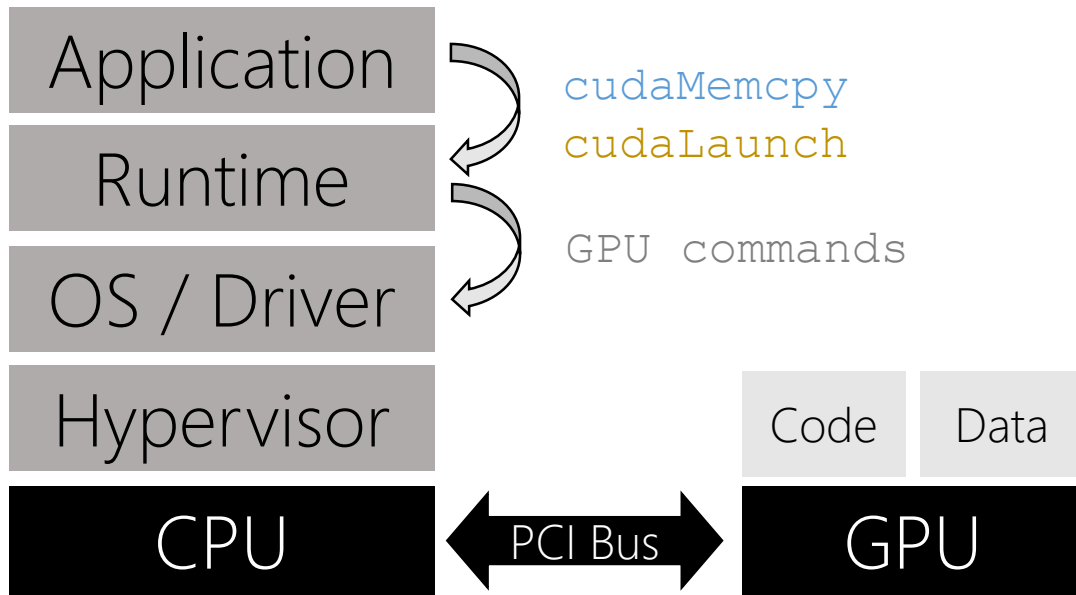- Remote attestation for establishing trust

Contributions

- Graviton architecture with minimal hardware extensions
- Extensions to CUDA runtime for end-to-end security
- Graviton implementation for demonstrating low performance overheads

# Outline

- Introduction
- GPUs & Threat Model
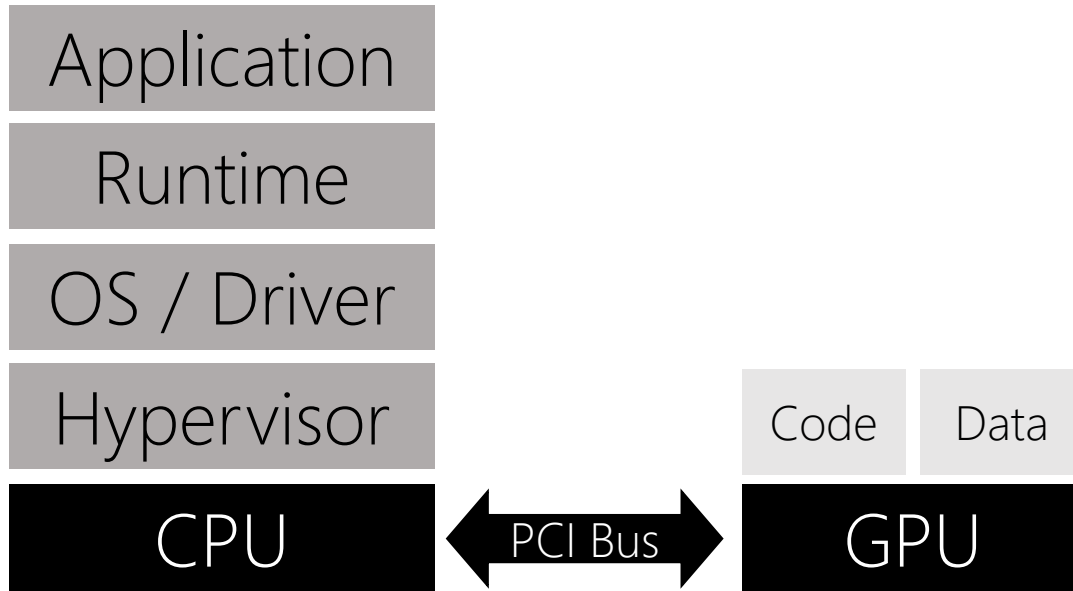- Graviton
- Evaluation
- Conclusion

# GPU 101: System Stack

Application

Runtime

`cudaMemcpy`
`cudaLaunch`

`GPU commands`

OS / Driver

Hypervisor

Code   Data

CPU   ← PCI Bus →   GPU

| PCIe Ctrl. | Command Processor (CP) | DMA Engines |
|---|---|---|
| Compute Eng. | CORE | CORE | CORE | CORE |
| | CORE | CORE | CORE | CORE |
| | CORE | CORE | CORE | CORE |
| | CORE | CORE | CORE | CORE |
| On-package Memory | | | |

GPU engines controlled via group of commands
- Generated by runtime and fetched by command processor

# GPU 101: Execution Model

Application

Runtime

OS / Driver

Hypervisor
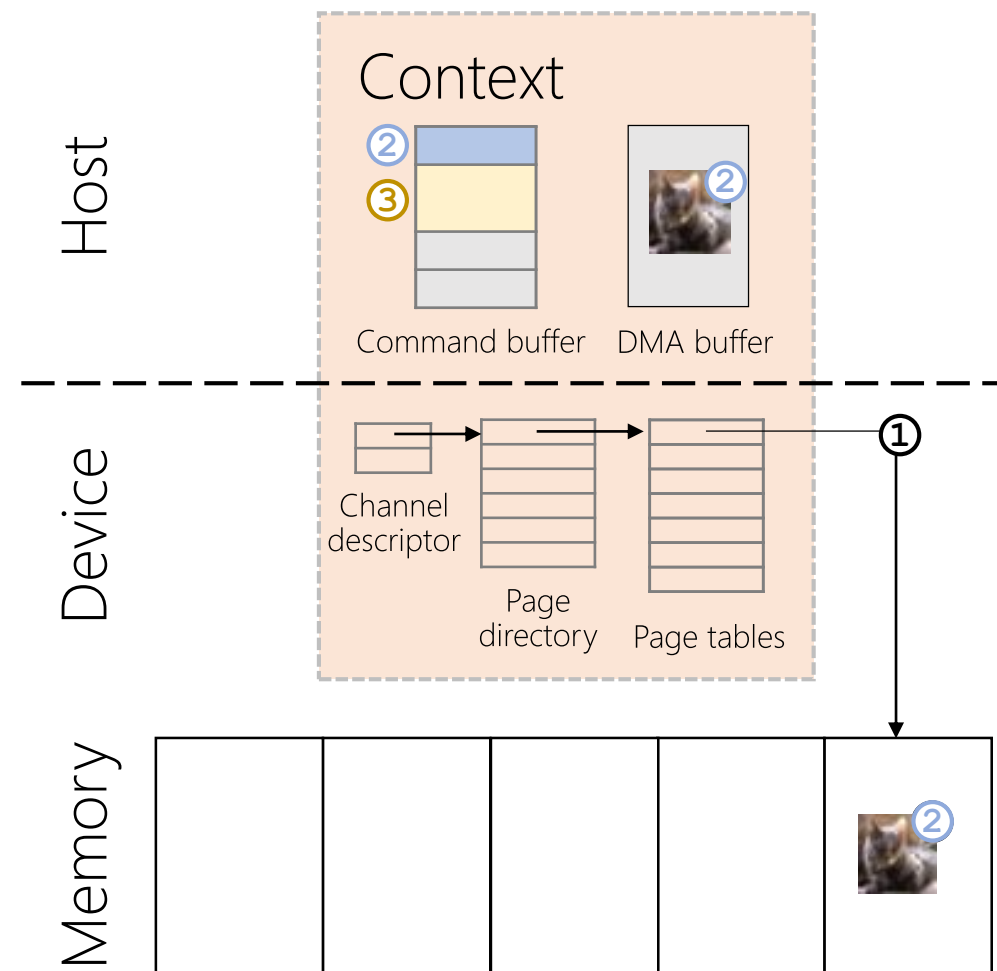
Code  Data

CPU  ◀— PCI Bus —▶  GPU

## Contexts supported by channels

- Implement *virtual memory* abstraction
- Expose command queues to runtime

Host

Context

Command buffer

Device

Channel descriptor

Page directory    Page tables

Memory

# GPU 101: Cat Classifier Example



① `cudaMalloc`
② `cudaMemcpy`
③ `cudaLaunch`

# GPU 101: Tampering with Commands & Data



Context 1

Malicious OS

Command buffer   DMA buffer

Channel descriptor

Page directory

Page tables

Host

Device

Memory

# GPU 101: Violating Context Isolation



Context 1

Command buffer

Channel descriptor

Page directory

Page tables

Malicious OS

Context 2

Command buffer

Channel descriptor

Page directory

Page tables

Host

Device

Memory

# Threat Model

Trusted computing base

- GPU package including on-package memory
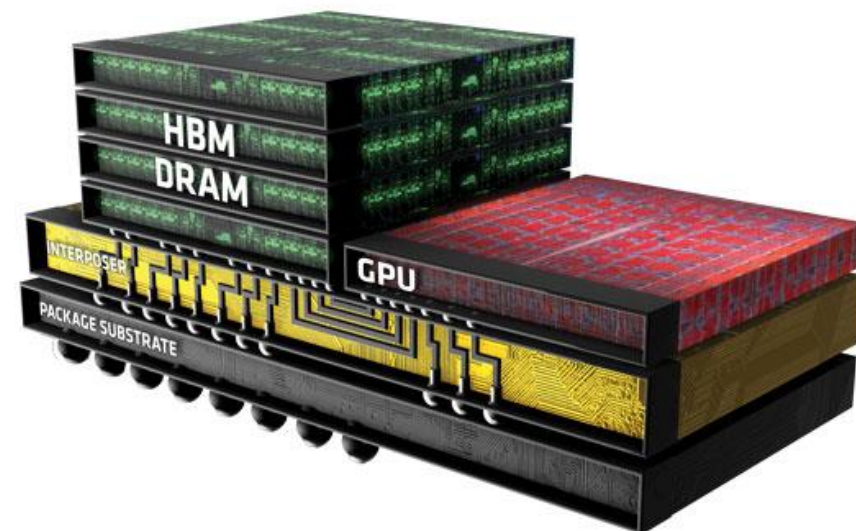- CPU package including TEE implementation
- GPU runtime hosted in CPU TEE

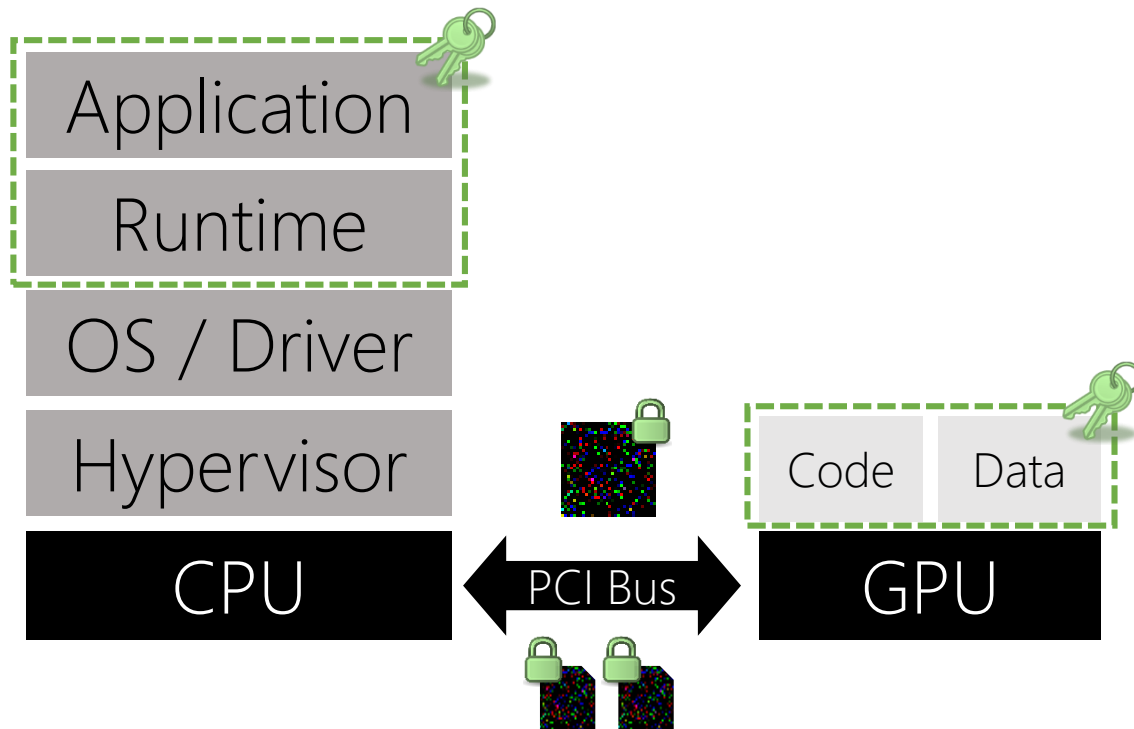Goal: Confidentiality and integrity of computation and data

Out of scope: side channels and package assembly attacks

# Outline

- Introduction
- GPUs & Threat Model
- Graviton
- Evaluation
- Conclusion

# Graviton: Overview

## Key concept: Redefined interface between hardware and software

Application

Runtime

OS / Driver

Hypervisor

CPU

PCI Bus

GPU

Code

Data

Hardware primitives in GPU

- Remote attestation for establishing trust
- Context isolation
- Secure command submission

Runtime abstractions

- Secure memory management
- Secure memory copy and task launch

# Graviton: Context Isolation
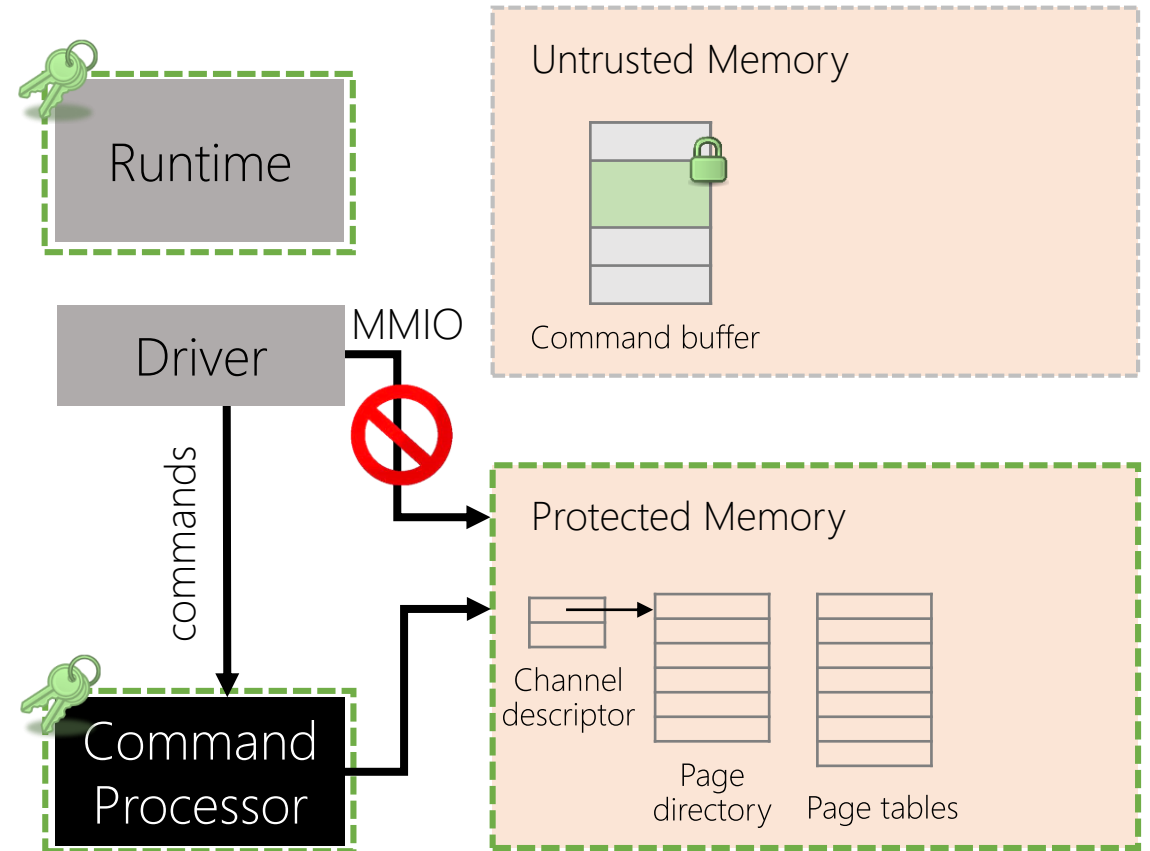
## Protected memory

- Hosts VM structures, code, and data
- CPU's MMIO accesses are blocked

## Virtual memory management via CP

- Ensures use of protected memory
- Exclusive use of context's memory resources

## Secure command submission

- Session key during context creation
- Only owner runtime can execute tasks

Runtime

Driver

MMIO

commands

Command Processor

Untrusted Memory

Command buffer

Protected Memory

Channel descriptor

Page directory

Page tables

# Graviton: Secure Memory Copy

## Key concept

- Data/code plaintext only inside TEEs
- Data/code ciphertext outside TEE (DMA buffer)

## Protocol



Runtime

Driver

commands

Command Processor

Untrusted Memory

Command buffer          DMA buffer

Protected Memory

Channel descriptor
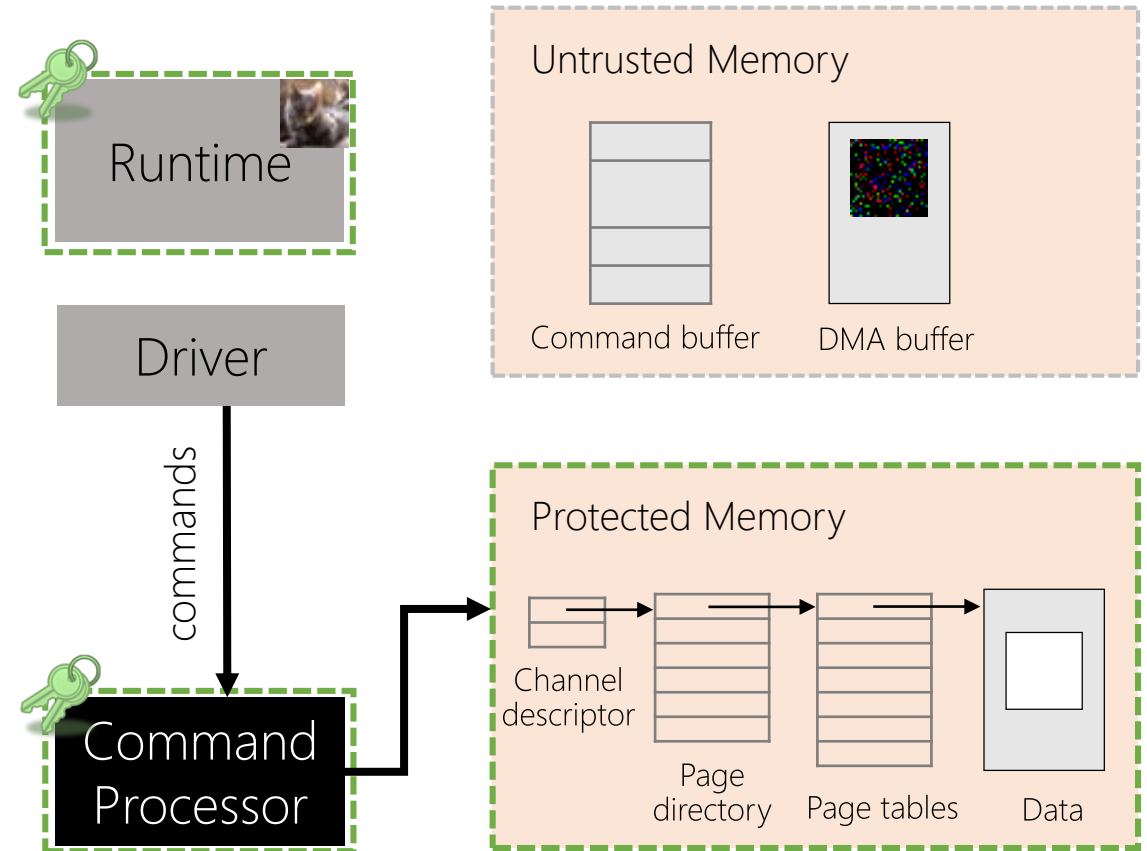
Page directory          Page tables          Data

# Graviton: Secure Memory Copy

## Key concept

- Data/code plaintext only inside TEEs
- Data/code ciphertext outside TEE (DMA buffer)

## Protocol

- Secure submission of *copy* task
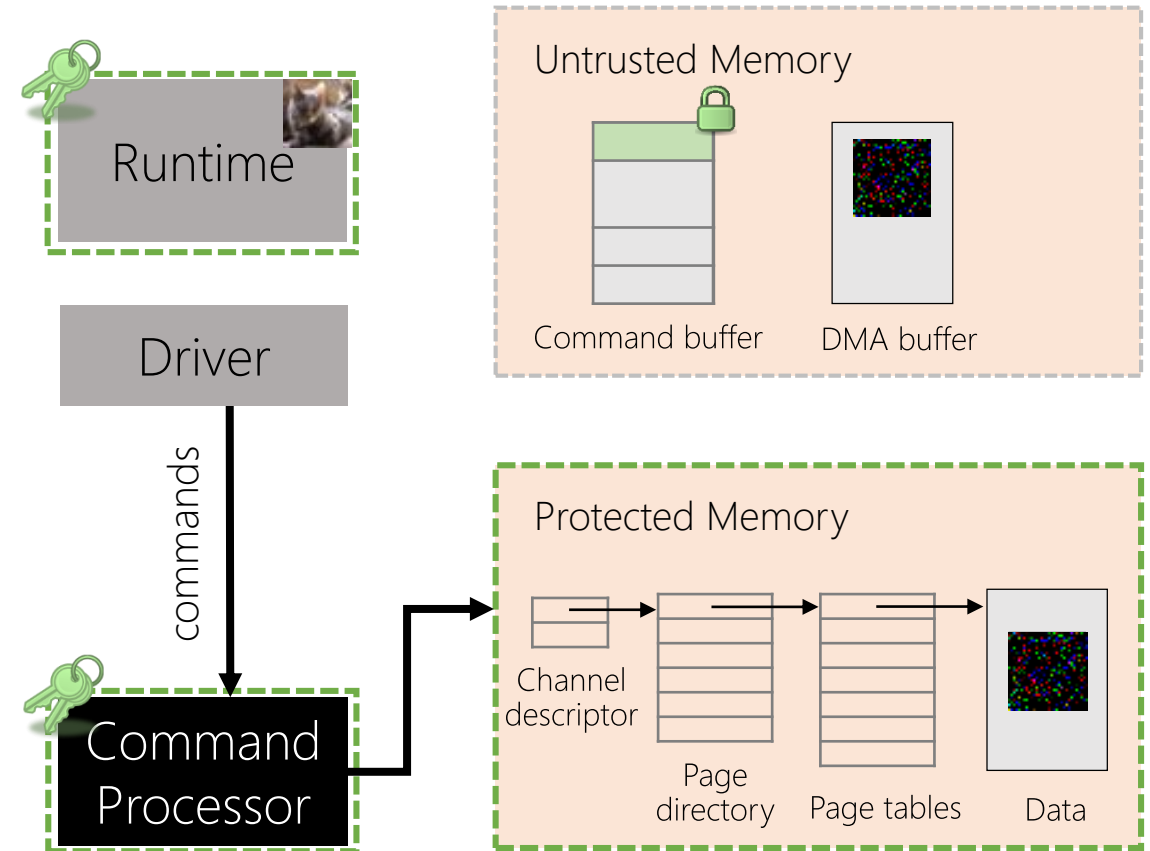
# Graviton: Secure Memory Copy

## Key concept

- Data/code plaintext only inside TEEs
- Data/code ciphertext outside TEE (DMA buffer)

## Protocol

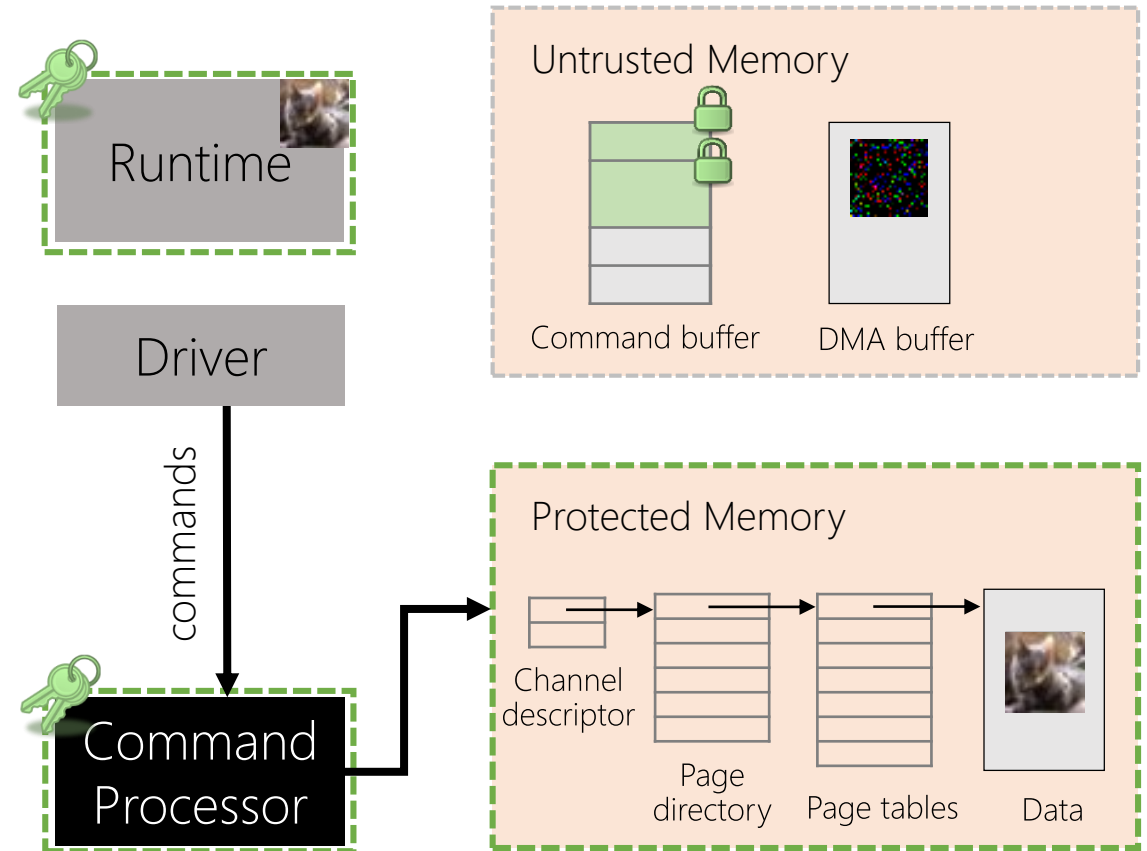- Secure submission of *copy* task
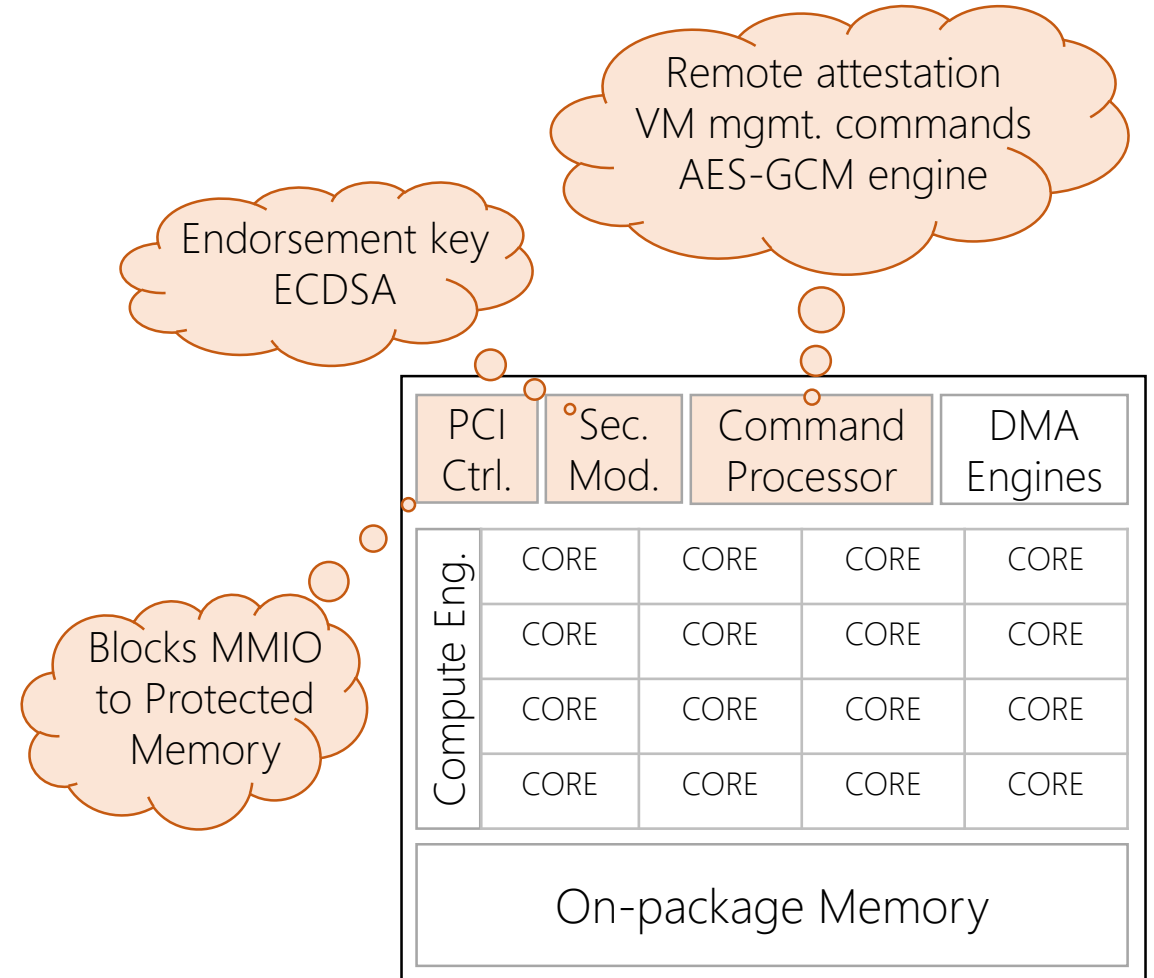- Secure submission for *authenticated decryption*

Runtime

Driver

commands

Command
Processor

Untrusted Memory

Command buffer      DMA buffer

Protected Memory

Channel
descriptor

Page
directory      Page tables      Data

# Graviton in a Nutshell

## Low hardware complexity

- Changes limited to peripheral components
- No changes to CPU, GPU cores and memory

## Transparent to developers

- GPU runtime abstractions
- Hidden behind GPU programming model

Remote attestation
VM mgmt. commands
AES-GCM engine

Endorsement key
ECDSA

Blocks MMIO
to Protected
Memory

| PCI Ctrl. | Sec. Mod. | Command Processor | DMA Engines |
|---|---|---|---|
| Compute Eng. | CORE | CORE | CORE | CORE |
| | CORE | CORE | CORE | CORE |
| | CORE | CORE | CORE | CORE |
| | CORE | CORE | CORE | CORE |

On-package Memory
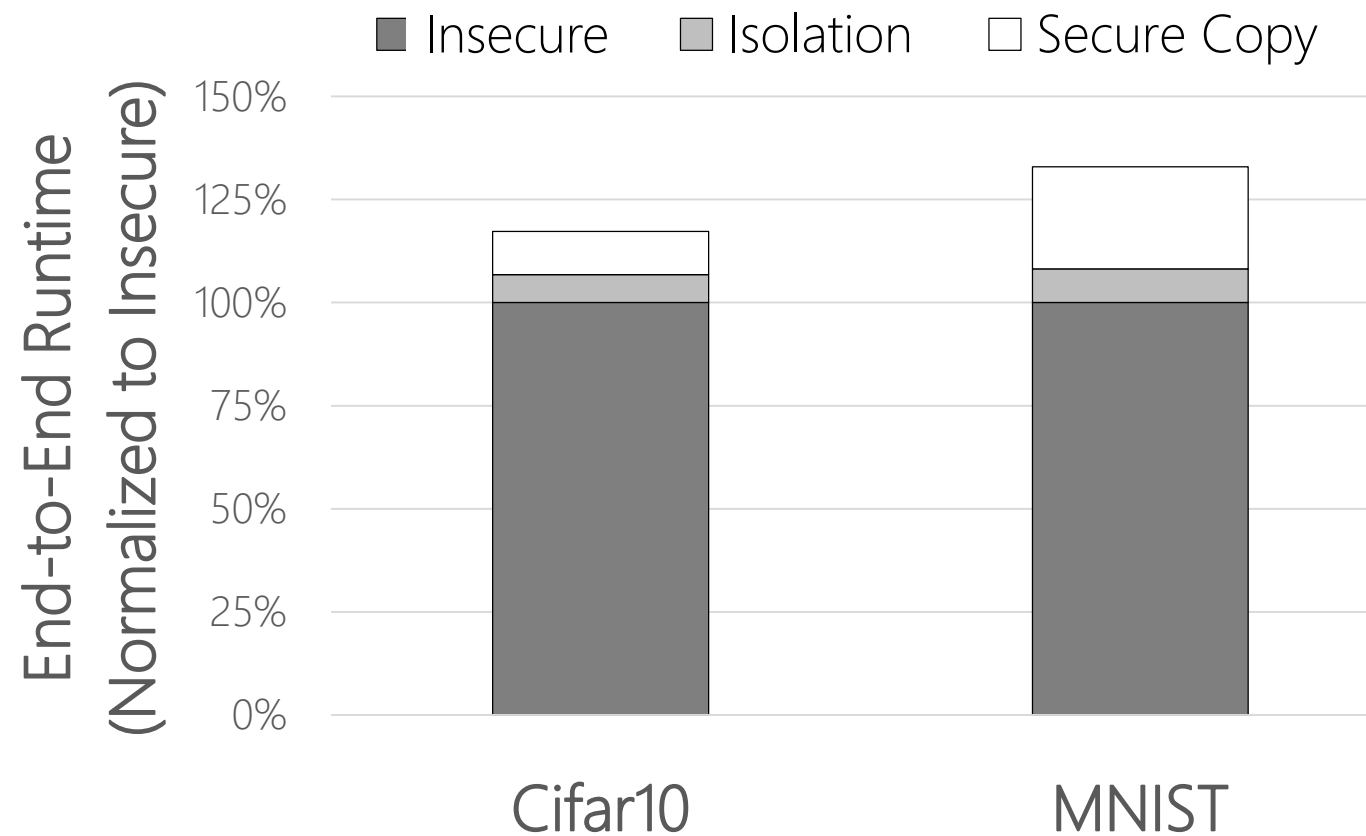
# Implementation

NVIDIA GTX Titan Black

- 2880 CUDA cores, 6GB of memory, peak performance 5.6 TFLOPS

Prototype

- GPU runtime: *secure task submission* and *secure memory management*
- Device driver: *address-space mgmt. command submission*
- Hardware primitives: emulation of new commands and crypto in device driver

Benchmarks: Cifar10-CNN and MNIST-autoencoder

# Implications on System Performance



Legend: ■ Insecure  ■ Isolation  □ Secure Copy

Y-axis: End-to-End Runtime (Normalized to Insecure) — 0%, 25%, 50%, 75%, 100%, 125%, 150%

X-axis: Cifar10, MNIST

Isolation
- Secure context management
- Secure command submission

Secure copy
- Host-side authenticated encryption
- GPU-side authenticated encryption

*Overhead correlates with ratio between computation and I/O*

# Concluding Remarks

Cloud trends in collision

- Confidentiality and hardware acceleration
- But, confidential computing restricted to CPUs

Graviton: Trusted Execution Environments on GPUs

- Low hardware complexity
- Low performance overheads
- Hardware complexity hidden by GPU programming model