

# DAGguise

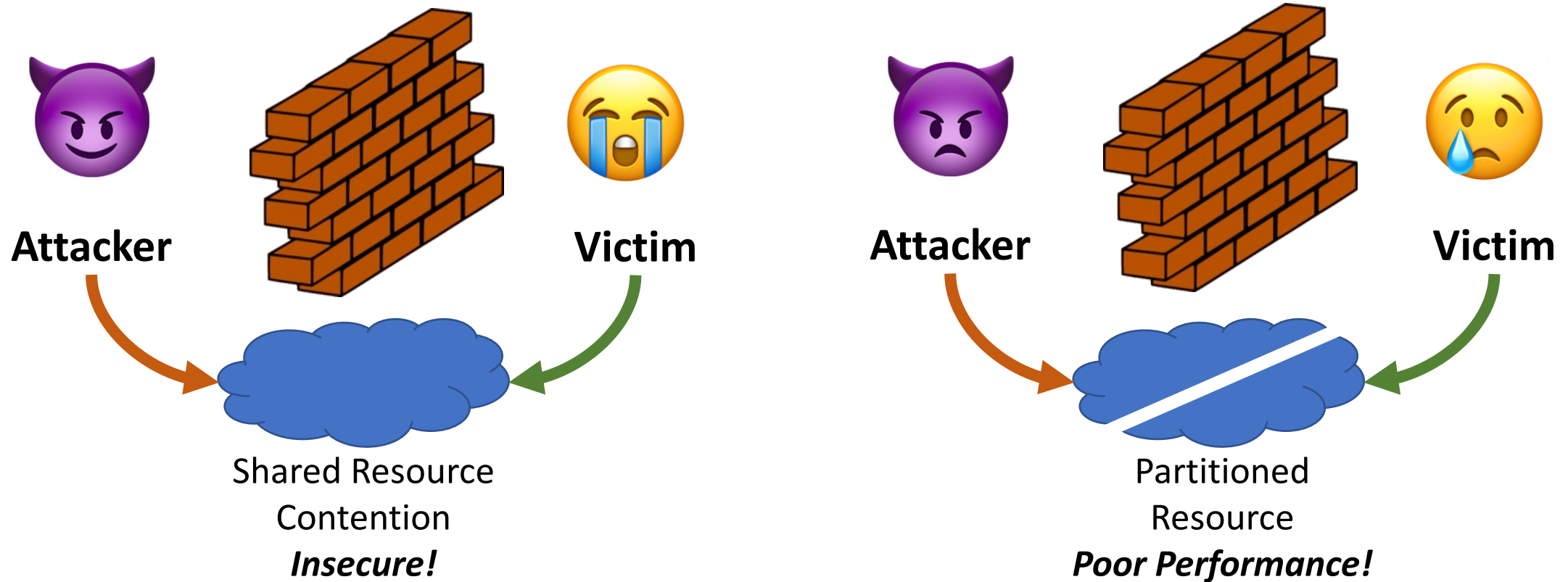
## Mitigating Memory Controller Side Channels

*Peter W. Deutsch\*, Yuheng Yang\*, Thomas Bourgeat,  
Jules Drean, Joel Emer, and Mengjia Yan*

ASPLOS 2022 (Session 3A)

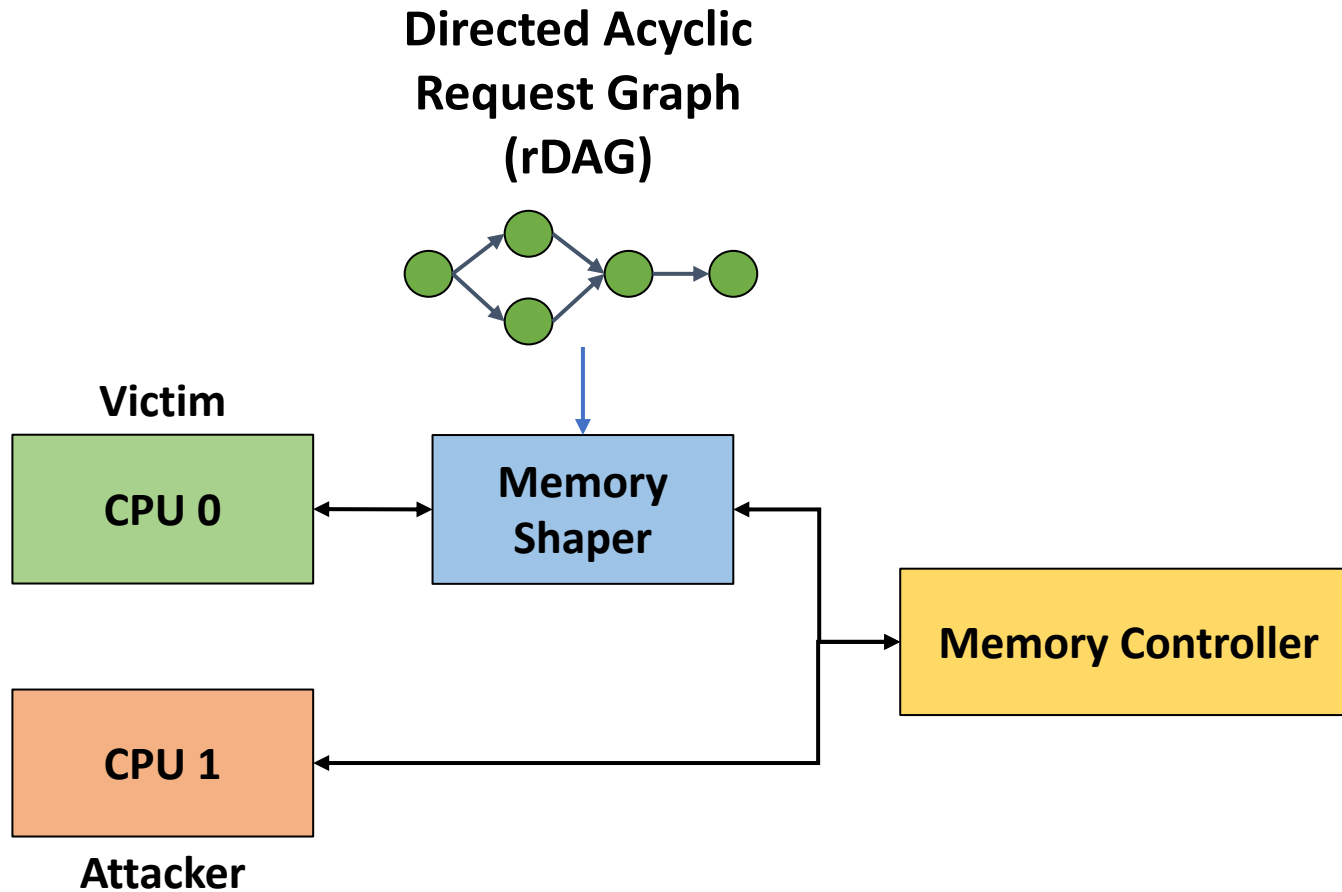


# Microarchitectural Side-Channels



**Key Defense Tradeoff: Security vs. Performance**

# DAGguise Key Idea



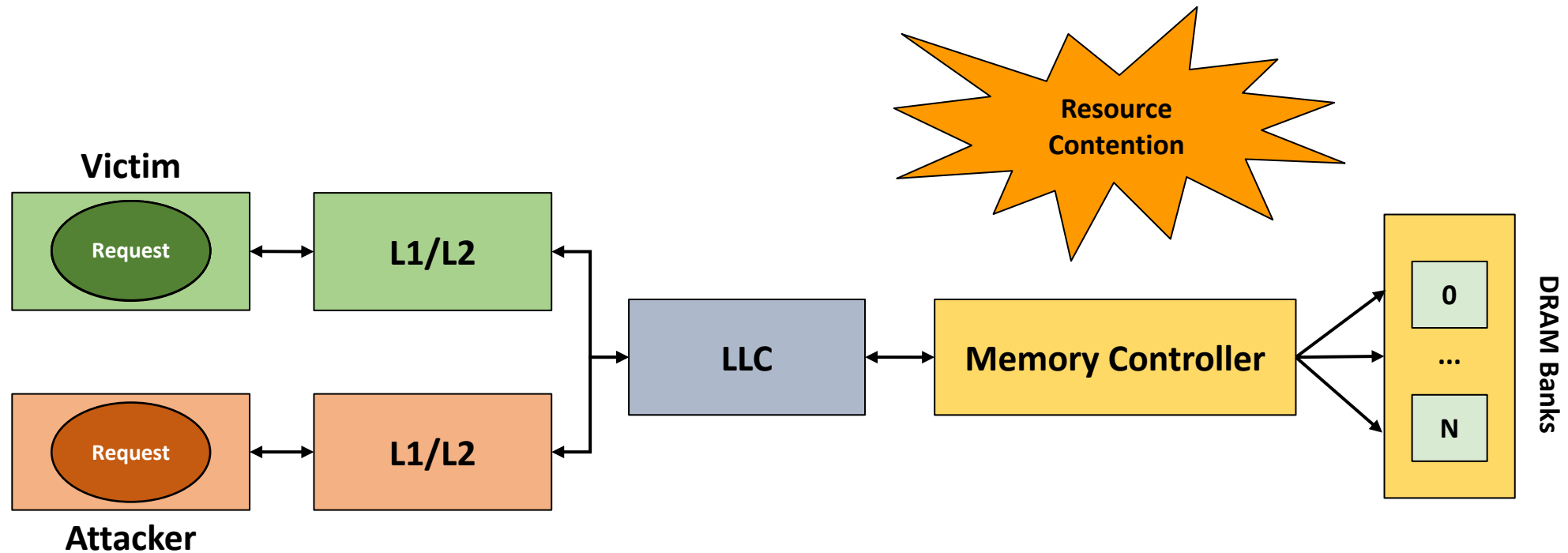
**DAGguise achieves:**

- ✓ Formally-Verified Security
- and*
- ✓ Good Performance

# Outline

- Memory Controller + Scheduler-based Side Channels
- Existing Approaches
  - Static Partitioning
  - Traffic Shaping
- DAGguise
  - Directed Acyclic Request Graphs (rDAGs)
- Security + Performance Evaluation
- Generalizability

# Memory Controller Side Channels



This is a class of “*scheduler-based*” side channels!

# Scheduler-Based Side Channels

This is the extended version of a paper that appears in USENIX Security 2021

## Lord of the Ring(s): Side Channel Attacks on the CPU On-Chip Ring Interconnect Are Practical

Riccardo Paccagnella Licheng Luo Christopher W. Fletcher  
University of Illinois at Urbana-Champaign



## DRAMA: Exploiting DRAM Addressing for Cross-CPU Attacks

Peter Pessl, Daniel Gruss, Clémentine Maurice, Michael Schwarz and Stefan Mangard  
Graz University of Technology, Austria

## Bandwidth Utilization Side-Channel on ML Inference Accelerators

Sarbartha Banerjee The University of Texas at Austin sarbartha@utexas.edu  
Shijia Wei The University of Texas at Austin shijiawei@utexas.edu  
Prakash Ramrakhiani ARM Research prakash.ramrakhiani@arm.com  
Mohit Tiwari The University of Texas at Austin tiwari@austin.utexas.edu

Abstract—Accelerators used for machine learning (ML) inference provide great performance but expose more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

Abstract—Accelerators used for machine learning (ML) inference provide great performance but expose more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

Abstract—Accelerators used for machine learning (ML) inference provide great performance but expose more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

Abstract—Accelerators used for machine learning (ML) inference provide great performance but expose more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

Abstract—Accelerators used for machine learning (ML) inference provide great performance but expose more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

Abstract—Accelerators used for machine learning (ML) inference provide great performance but expose more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

Abstract—Accelerators used for machine learning (ML) inference provide great performance but expose more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

Abstract—Accelerators used for machine learning (ML) inference provide great performance but expose more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

Abstract—Accelerators used for machine learning (ML) inference provide great performance but expose more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

Abstract—Accelerators used for machine learning (ML) inference provide great performance but expose more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

Abstract—Accelerators used for machine learning (ML) inference provide great performance but expose more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

Abstract—Accelerators used for machine learning (ML) inference provide great performance but expose more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

Abstract—Accelerators used for machine learning (ML) inference provide great performance but expose more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

In cloud computing, servers are often co-located on the same physical hardware. Thus, preventing information leakage is crucial. While the CPU cache is a well-known side channel, shared hardware components like the memory bus, shared memory, and cache coherence protocols can also leak information. Furthermore, CPU cache is not the only side channel. In this setting, a slow cross-CPU communication channel as well as cache contention can be used to build these attacks, address mappings.

Abstract—SMT architectures are attractive targets for side-channel enabled attackers, with their inherently broader attack surface that exposes more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

Abstract—SMT architectures are attractive targets for side-channel enabled attackers, with their inherently broader attack surface that exposes more periphery than their host CPUs. This side-channel leakage source is exposed to an on-chip interconnect. We demonstrate our attack on EdDSA and RSA by measuring the precise timing of key stream generation.

## SMoTherSpectre: Exploiting Speculative Execution through Port Contention

Atri Bhattacharyya\* EPFL  
Alessandro Sorniotti† IBM Research - Zurich  
Alexandra Sandulescu† IBM Research - Zurich  
Babak Falsafi\* EPFL  
Matthias Neugschwandtner\* IBM Research - Zurich  
Mathias Payer\* EPFL  
Amil Kurmus† IBM Research - Zurich

Session 10D: VulnDet 2 + Side Channels 2

CCS'18, October 15-19, 2018, Toronto, ON, Canada

## Rendered Insecure: GPU Side Channel Attacks are Practical

Hoda Naghibijouybari University of California, Riverside hnagh001@ucr.edu  
Zhiyun Qian University of California, Riverside zhiyunq@cs.ucr.edu  
Ajaya Neupane University of California, Riverside ajaya@ucr.edu  
Nael Abu-Ghazaleh University of California, Riverside nael@cs.ucr.edu

### ABSTRACT

Graphics Processing Units (GPUs) are commonly integrated with computing devices to enhance the performance and capabilities of graphical workloads. In addition, they are increasingly being integrated in data centers and clouds such that they can be used to accelerate data intensive workloads. Under a number of scenarios the GPU can be shared between multiple applications at a fine granularity allowing a spy application to monitor side channels and attempt to infer the behavior of the victim. For example, OpenGL and WebGL send workloads to the GPU at the granularity of a frame, allowing an attacker to interleave the use of the GPU to measure the side-effects of the victim computation through performance counters or other resource tracking APIs. We demonstrate

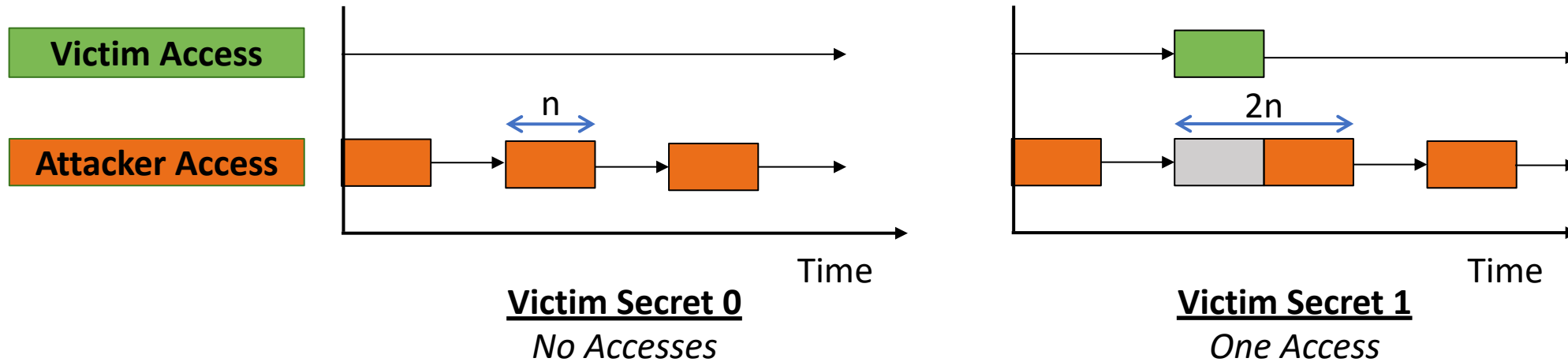
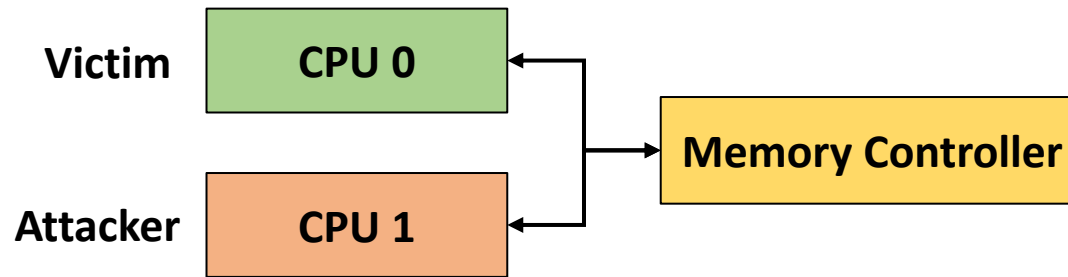
### 1 INTRODUCTION

Graphics Processing Units (GPUs) are integral components to most modern computing devices, used to optimize the performance of today's graphics and multi-media heavy workloads. They are also increasingly integrated on computing servers to accelerate a range of applications from domains including security, computer vision, computational finance, bio-informatics and many others [52]. Both these classes of applications can operate on sensitive data [25, 31, 57] which can be compromised by security vulnerabilities in the GPU stack. Although the security of GPUs is only starting to be explored, several vulnerabilities have already been demonstrated [46, 49, 55, 58, 63, 71, 74]. Most related to this paper, Luo et al. demonstrated a

arXiv:2110.07157v1 [cs.CR] 14 Oct 2021

arXiv:1903.01843v3 [cs.CR] 26 Sep 2019

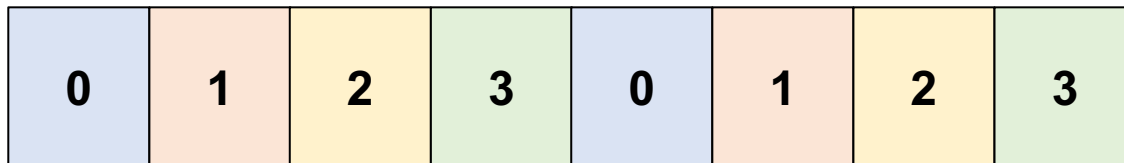
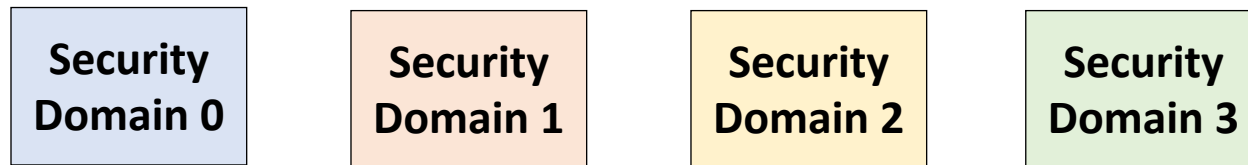
# Timing Attack Example



The attacker uses its *own* latencies to leak information!

# Static Partitioning in Time

Use a Round Robin, No-Skip Arbitration Policy



→  
Slot Allocation Timeline

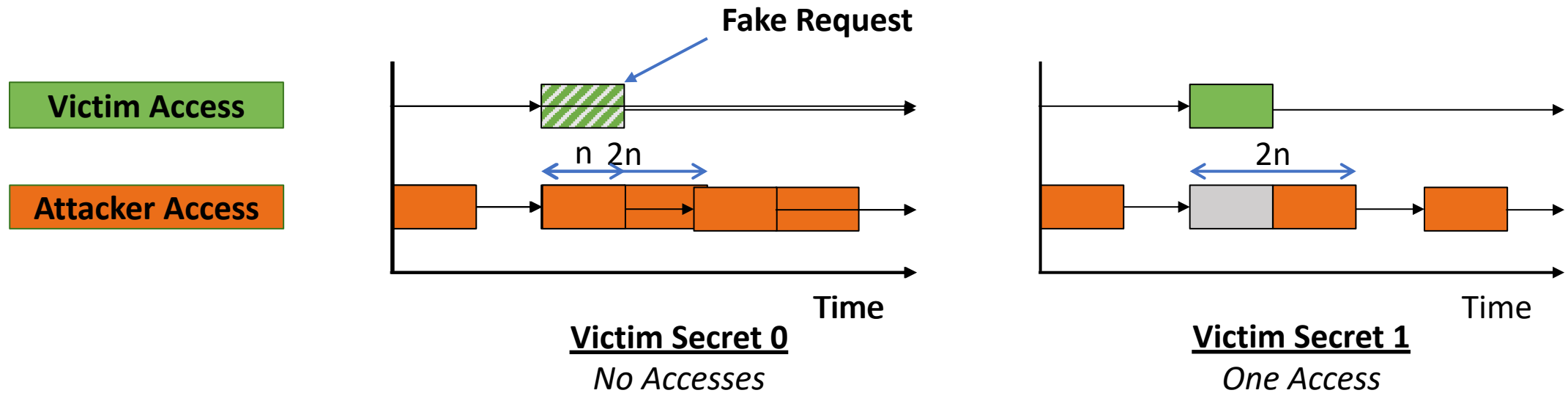
✓ Secure  
Static partitioning, no leakage

✗ Bad Performance  
Poor bandwidth utilization!



# Traffic Shaping

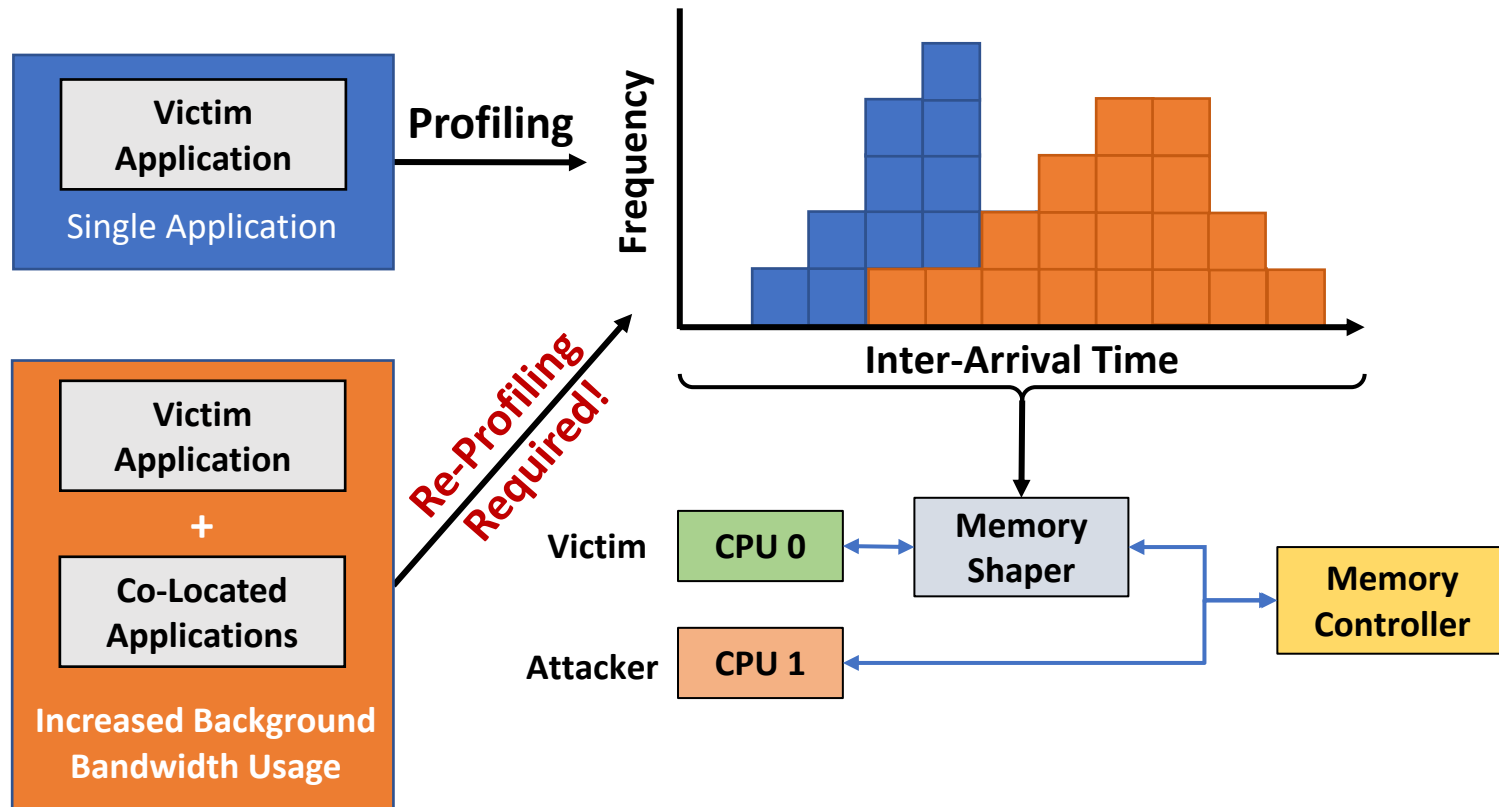
**Shaping Strategy:** Delay victim's existing requests and add fake requests



How do we do this for real applications without significant costs?

# Camouflage's Traffic Shaping Strategy

Shape memory requests to a secret-independent *timing distribution*



## ✓ Good Performance

Dynamic sharing of the memory controller

## ✗ Insecure

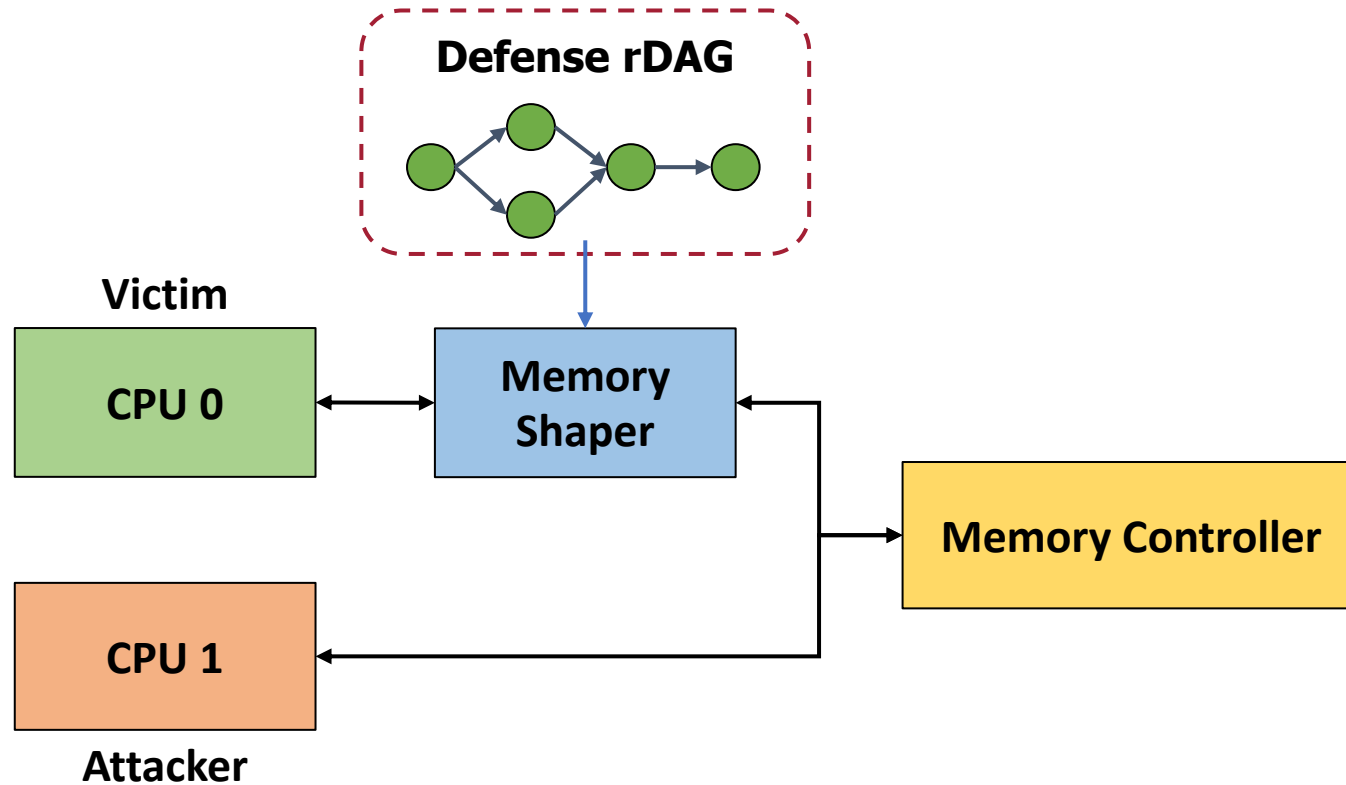
Ordering or bank information can reveal the secret

## ✗ Expensive Profiling

Ideal shaping distribution depends on co-running applications

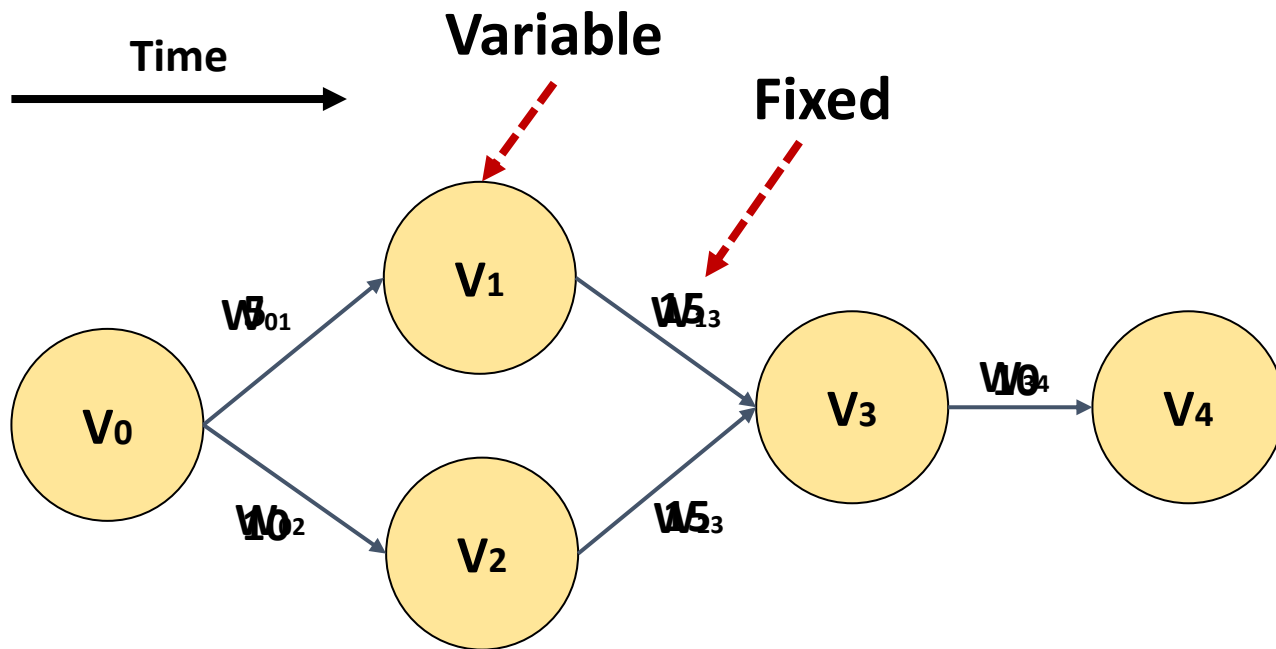
# DAGguise's Traffic Shaping Strategy

Shape memory requests to a secret-independent  
*Directed Acyclic Request Graph (rDAG)*



- ✓ Secure
- ✓ Good Performance
- ✓ Profile Victim Alone

# Directed Acyclic Request Graphs

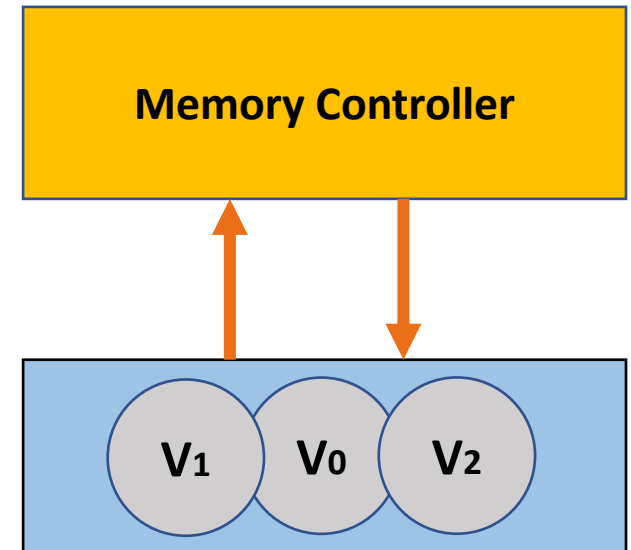


## Vertices

Memory requests with *variable* latency

## Edges

Dependencies between memory requests with *fixed* latency



# Why shape requests to an rDAG?

## ✓ Security

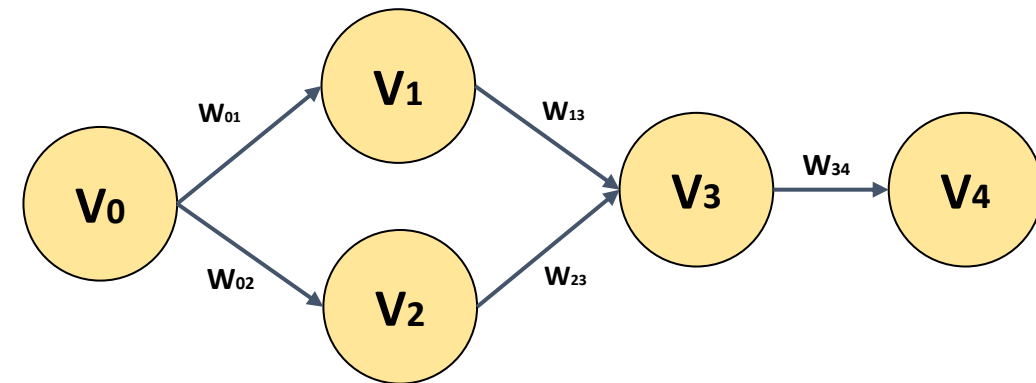
- Shaping to a secret-independent defense rDAG makes victim request patterns *indistinguishable*
- Defense rDAGs are public and are the only thing an attacker can recover

## ✓ Performance

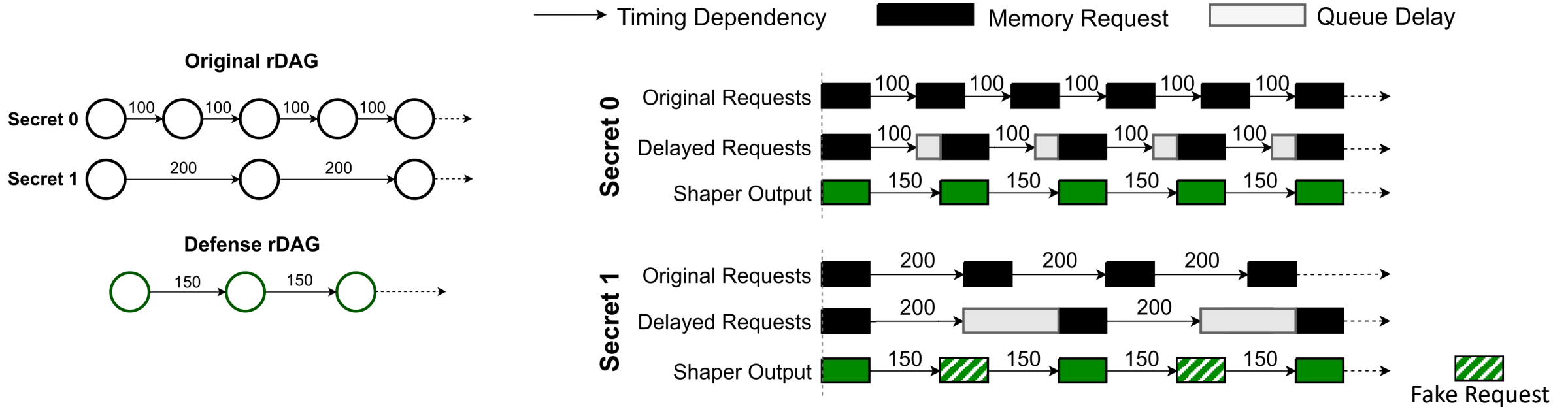
- Allows for *dynamic* sharing of memory resources in the memory controller

## ✓ Profiling Cost

- Does not require knowledge of co-located applications

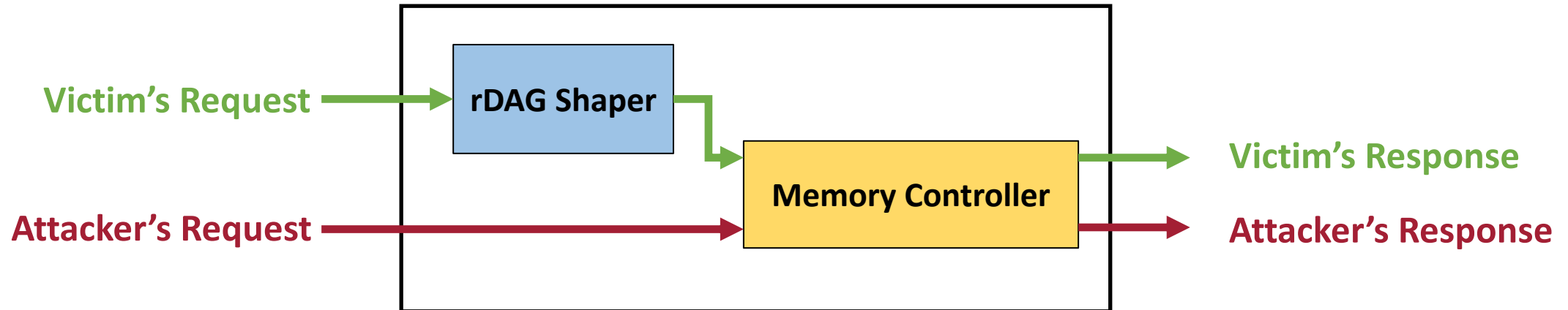


# Simple Shaping Example



**The shaper output is always the same, *no matter the secret!***

# Indistinguishability Property



The attacker's observations should be *independent* from victim's request pattern

# Indistinguishability Property

- Attacker's observation is independent from victim's request pattern
  - Given an attacker's request pattern, the attacker has an identical observation when contending with **ANY** victim's request pattern
  - This holds for **ANY** attacker's request pattern

Attacker's Observations when Contending with Victim

Victim Request Patterns Attacker Request Patterns	A	B	C	... ..
	Attacker's Response Pattern X			
X				



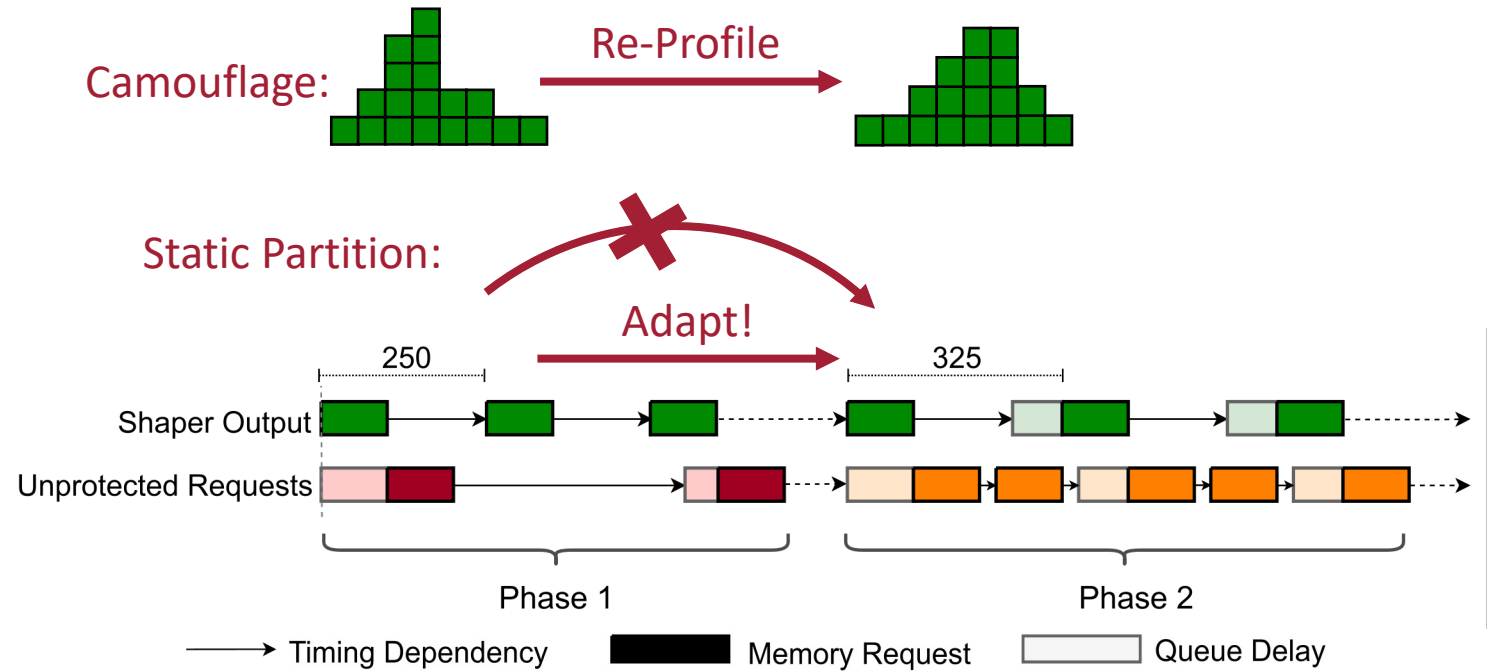
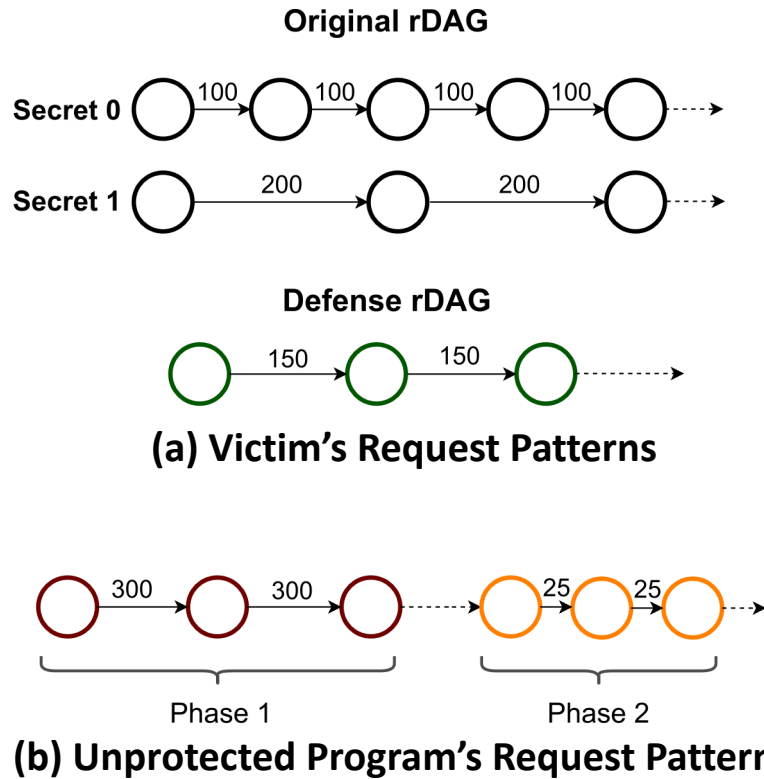
# Formalization & Verification

- Formalize the indistinguishability property using state transitions

$$\begin{aligned} P(S_0, n) := & \forall \text{Req}_{Tx}, \text{Req}'_{Tx}, \forall \text{Req}_{Rx} \\ & \text{if } S_0 \xrightarrow[\text{Req}_{Tx}, \text{Req}_{Rx}]{\text{Resp}_{Tx}, \text{Resp}_{Rx}} S_n \text{ and } S_0 \xrightarrow[\text{Req}'_{Tx}, \text{Req}_{Rx}]{\text{Resp}'_{Tx}, \text{Resp}'_{Rx}} S'_n \\ & \text{then } \text{Resp}_{Rx} = \text{Resp}'_{Rx} \end{aligned}$$

- Verification with *Rosette*:
  - First k cycles: symbolic execution
  - Arbitrary cycles: k-induction

# rDAG Adaptivity

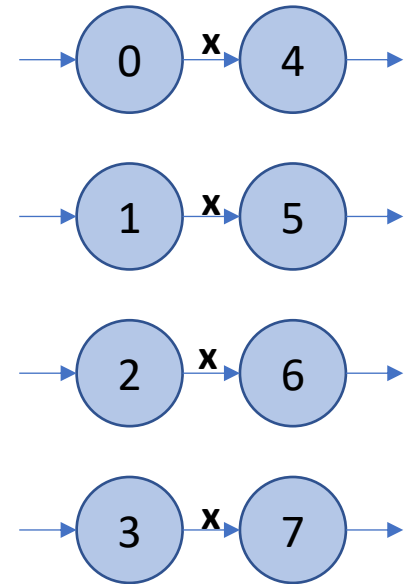


**(c) Contention between Victim and Unprotected Program on Memory Controller**

**rDAG's adaptivity allows for better bandwidth utilization!**

# Offline Profiling Step

- Not for security, **any** secret-independent rDAG ensures security
- Low profiling cost
  - Victim is profiled **alone**
  - Reduce search space by finding parameters for an rDAG *template*



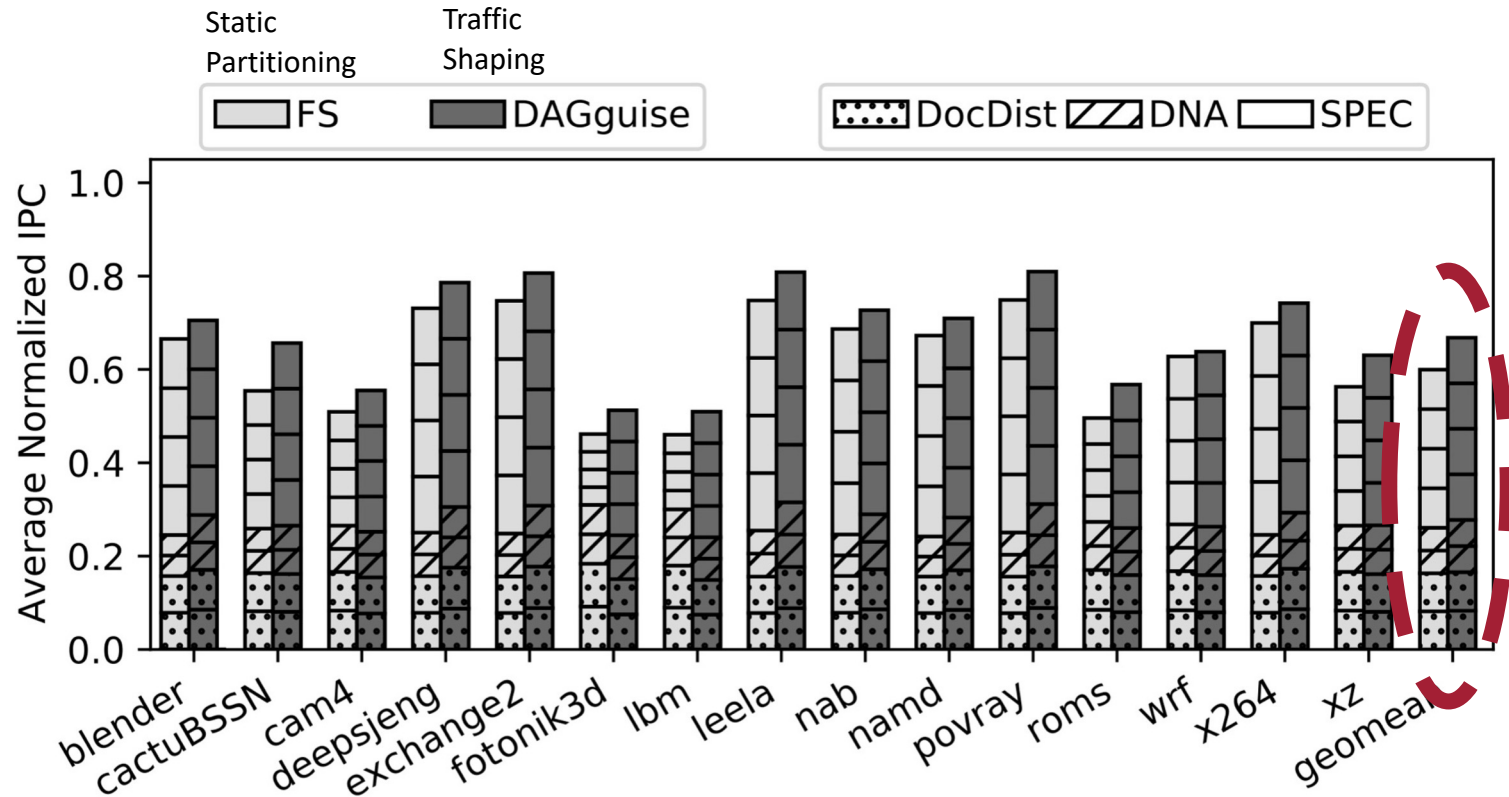
**4-Parallel rDAG Template**

# Experimental Setup

- **Simulator:** gem5 and DRAMSim2
- **Architectural Specifications:**
  - 2 and 8 out-of-order CPU cores
  - 32KB L1i/d, 256kB L2, 1MB/core L3
- **Evaluated Configurations:**
  - DAGguise
  - Fixed Service (Bank Triple Alternation)
  - Baseline
- **Evaluated Applications:**
  - Unprotected SPEC benchmark(s) co-running alongside DAGguise protected application(s)



# Experimental Results

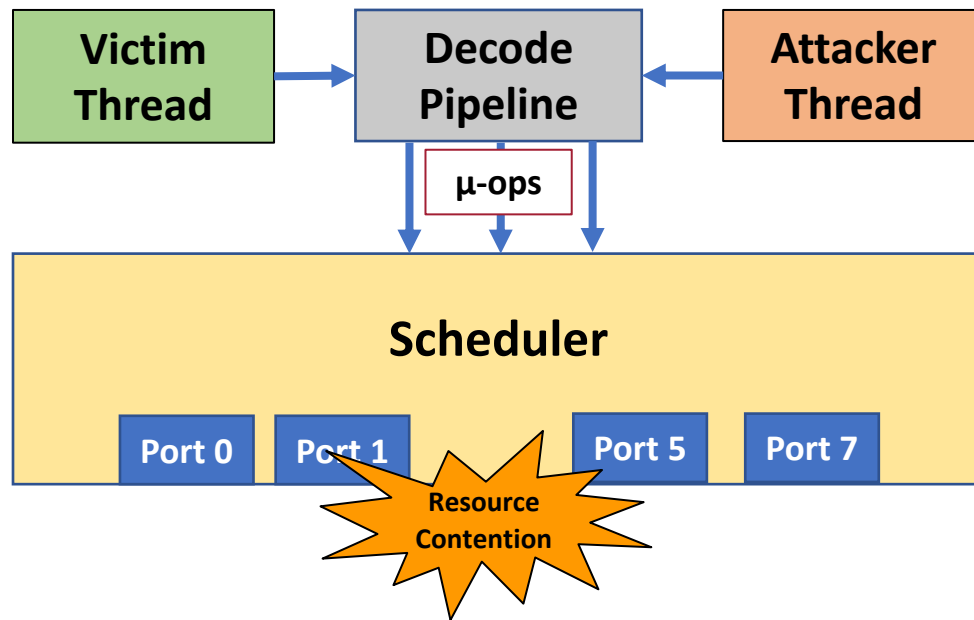


**DAGguise's improves performance for both protected *and* unprotected applications!**

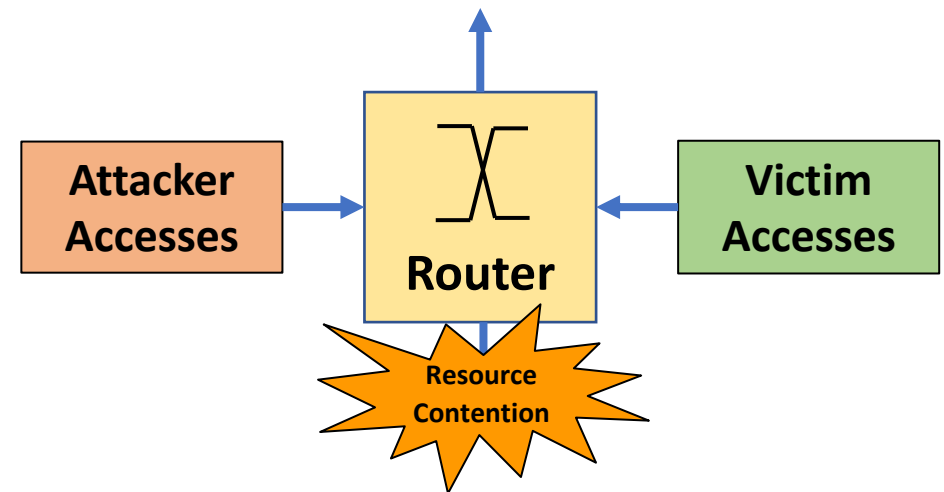
**DAGguise achieves a 12% performance improvement over Fixed Service in an 8-CPU system**

# DAGguise Generalization

## SMT Contention



## Network on Chip Contention

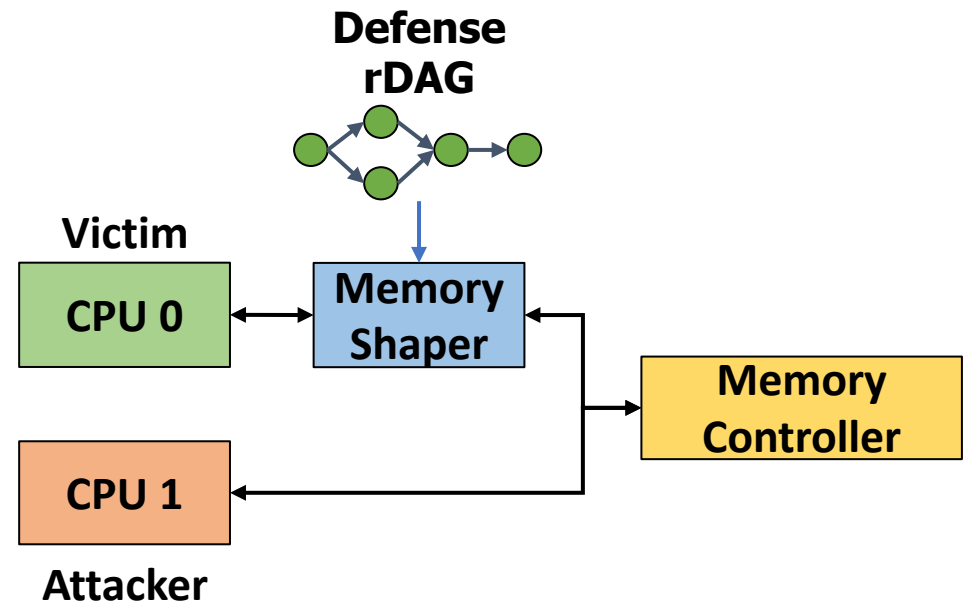


# More in the Paper

- Implementation details of DAGguise shaper
- Formal security verification using symbolic execution and k-induction
- Detailed rDAG offline profiling process
- More performance and area overhead evaluation
- Generalizations to other scheduler-based side channels (e.g. port contention)

# Conclusion

- DAGguise
  - A memory traffic shaper which:
    - Completely eliminates data leakage
    - Allows for dynamic contention
    - Requires only simple profiling
- rDAGs
  - A *general and adaptive* request representation
- A formal model of correctness using *Rosette*
- A generalized scheduler-based attack mitigation framework





# DAGguise

## Mitigating Memory Controller Side Channels

**Peter W. Deutsch**

[pwd@mit.edu](mailto:pwd@mit.edu)

**Yuheng Yang**

[yuhengy@mit.edu](mailto:yuhengy@mit.edu)

Thomas Bourgeat

[bthom@mit.edu](mailto:bthom@mit.edu)

Jules Drean

[drean@mit.edu](mailto:drean@mit.edu)

Joel S. Emer

[jsemer@mit.edu](mailto:jsemer@mit.edu)

Mengjia Yan

[mengjiay@mit.edu](mailto:mengjiay@mit.edu)

