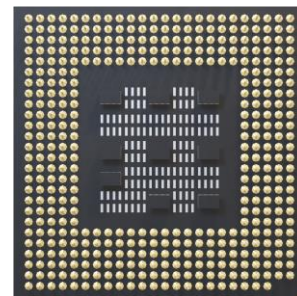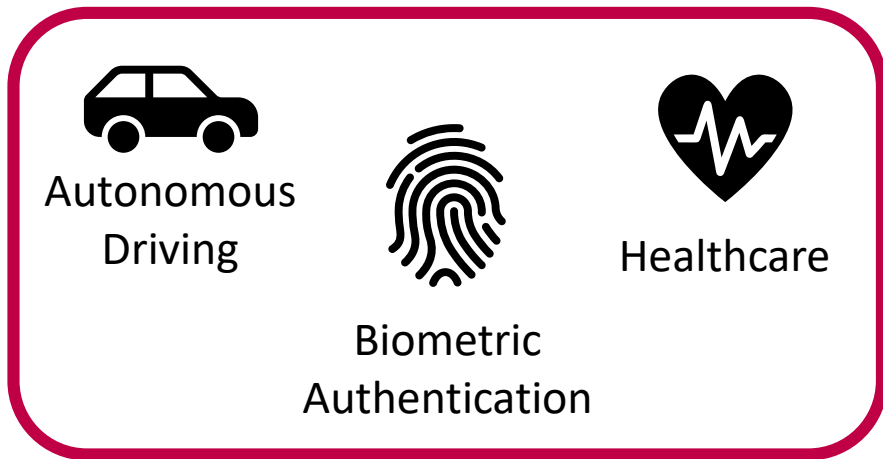# SecureLoop:
# Design Space Exploration of Secure DNN Accelerators

**Kyungmi Lee**, Mengjia Yan, Joel S. Emer*, Anantha P. Chandrakasan

MIT, *MIT/NVIDIA
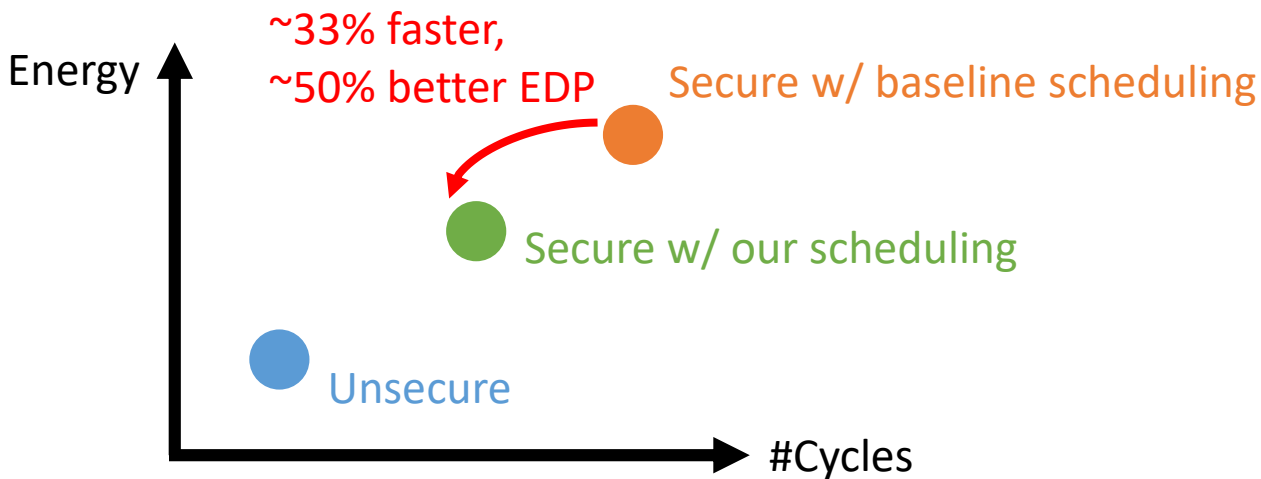
# ML Needs Both Security and Performance



**Applications require security**

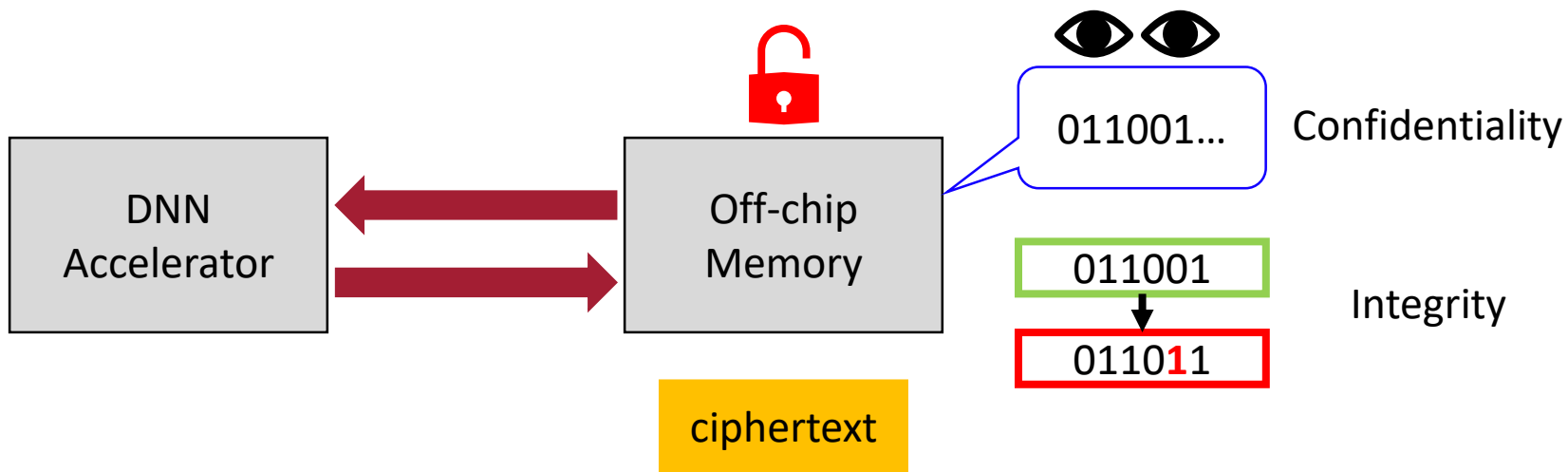**Accelerator design requires performance, area, energy**

Design space exploration for "secure" DNN accelerators

# SecureLoop

A tool for **design space exploration of secure DNN accelerators** equipped with cryptographic engines



~33% faster,
~50% better EDP

Energy

Secure w/ baseline scheduling

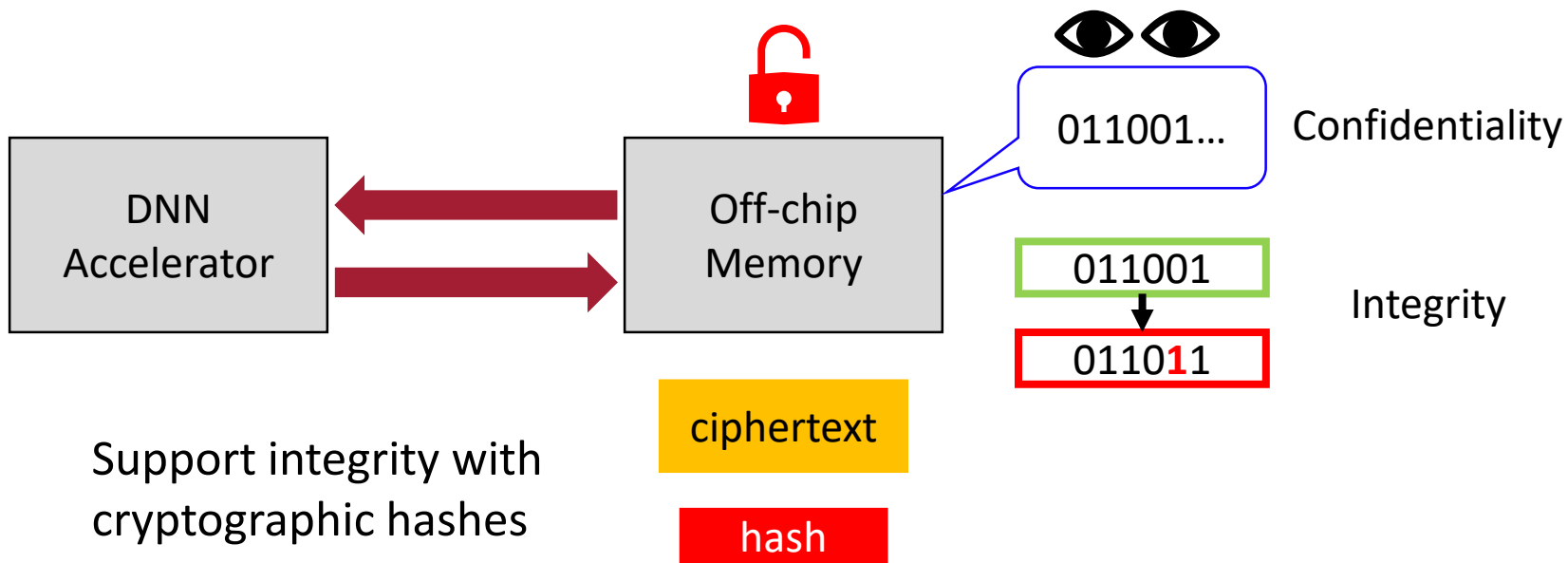Secure w/ our scheduling

Unsecure

#Cycles
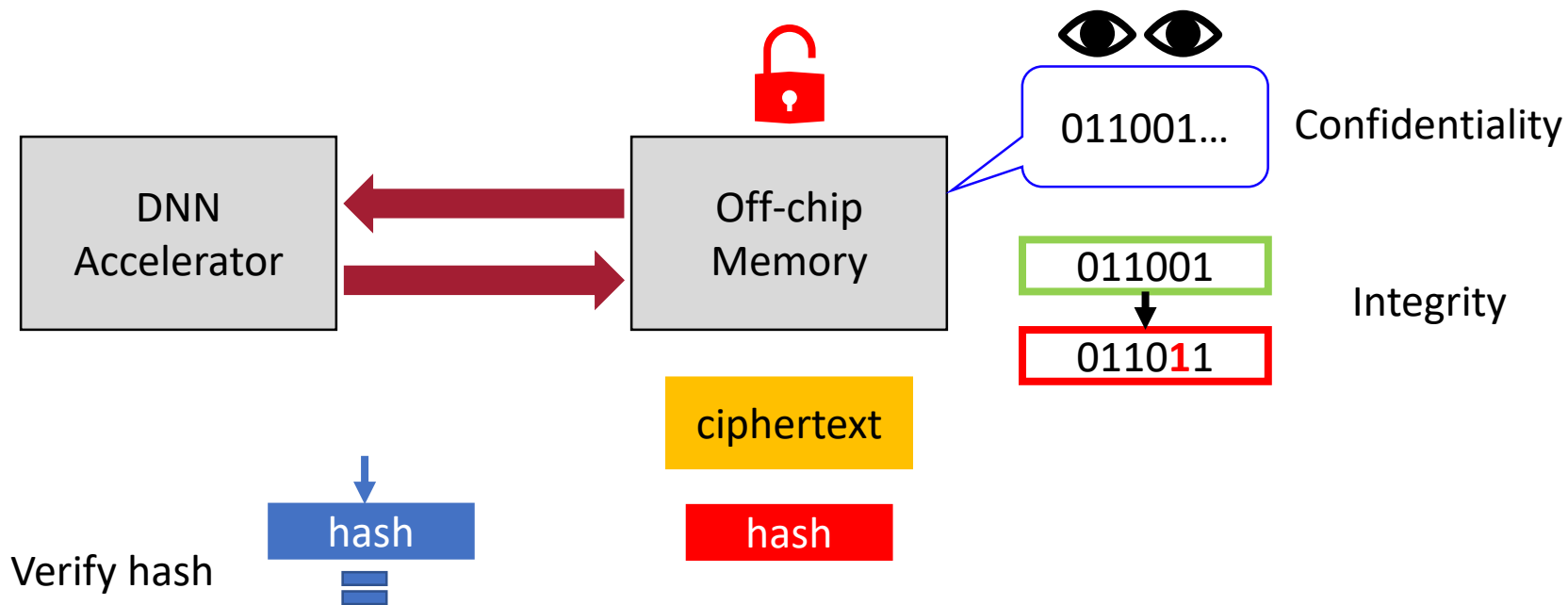
# Background: TEE and DNN Accelerators



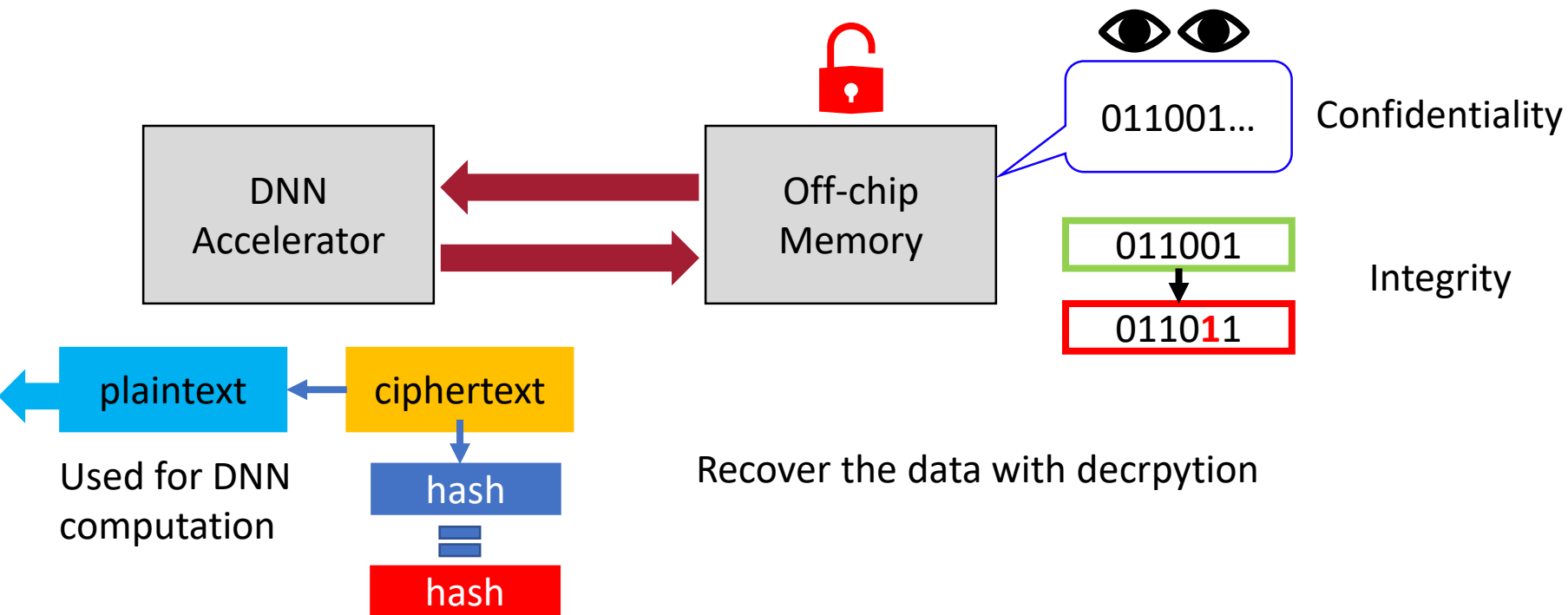Support confidentiality with cryptographic encryption

# Background: TEE and DNN Accelerators

# Background: TEE and DNN Accelerators

# Background: TEE and DNN Accelerators

# Background: TEE and DNN Accelerators



- Accelerator scheduling has to be coordinated with cryptographic operations

# What if Tile != Authentication Block

Tensor Data



hash

Tile
Auth
Block

Redundant reads

Tile

DNN Accelerator ⬅️➡️ Off-chip Memory

Can't get a hash only with this ▨▨ ta

Fetch extra "redundant" data for hash

additional off-chip traffic

**Baseline** Tile

**Secure Accelerator** Tile hash ▨▨

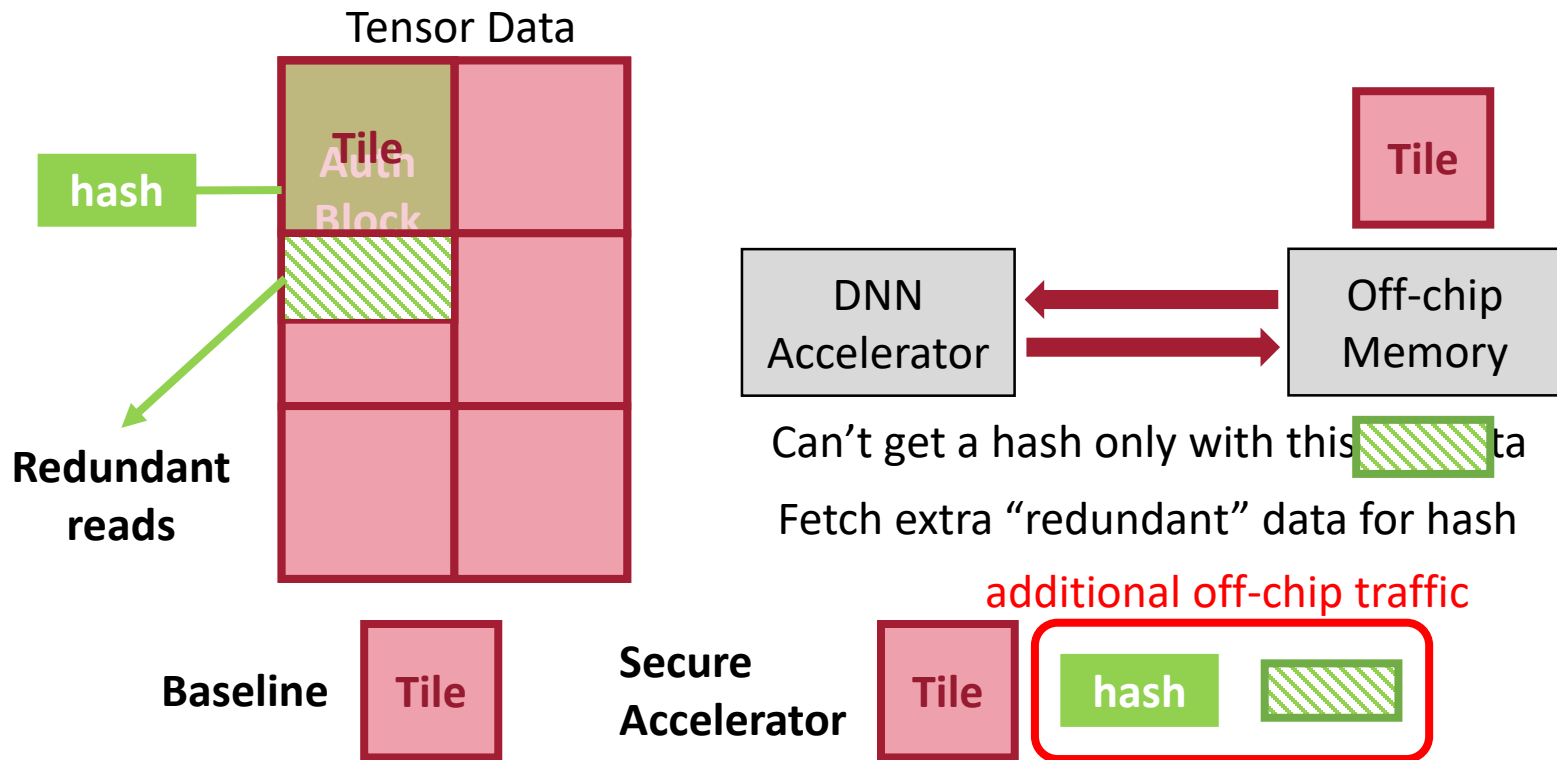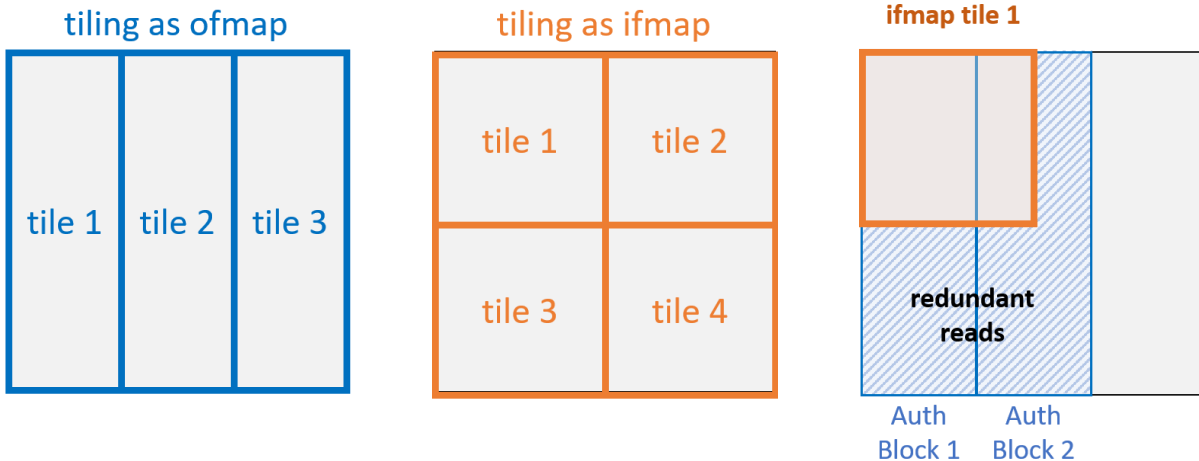# Tile-as-an-AuthBlock is not optimal

tiling as ofmap

| tile 1 | tile 2 | tile 3 |
|---|---|---|

tiling as ifmap

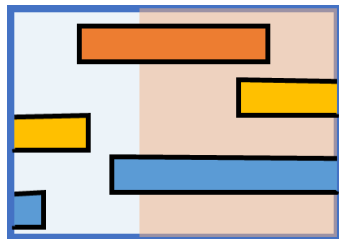| tile 1 | tile 2 |
|---|---|
| tile 3 | tile 4 |

ifmap tile 1

redundant reads
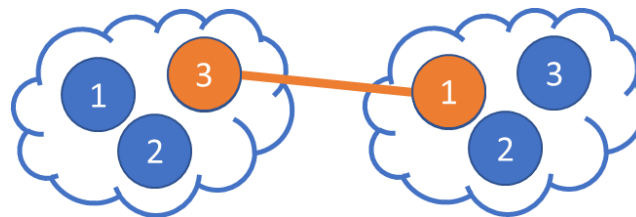
Auth Block 1    Auth Block 2

- Option 1: Tile-as-an-AuthBlock based on the output tiling ☹
- Option 2: Rehash between layers ☹

# Summary of our techniques

Analytical approach to identify the optimal AuthBlock assignment

Cross-layer fine tuning from the loopnest schedule
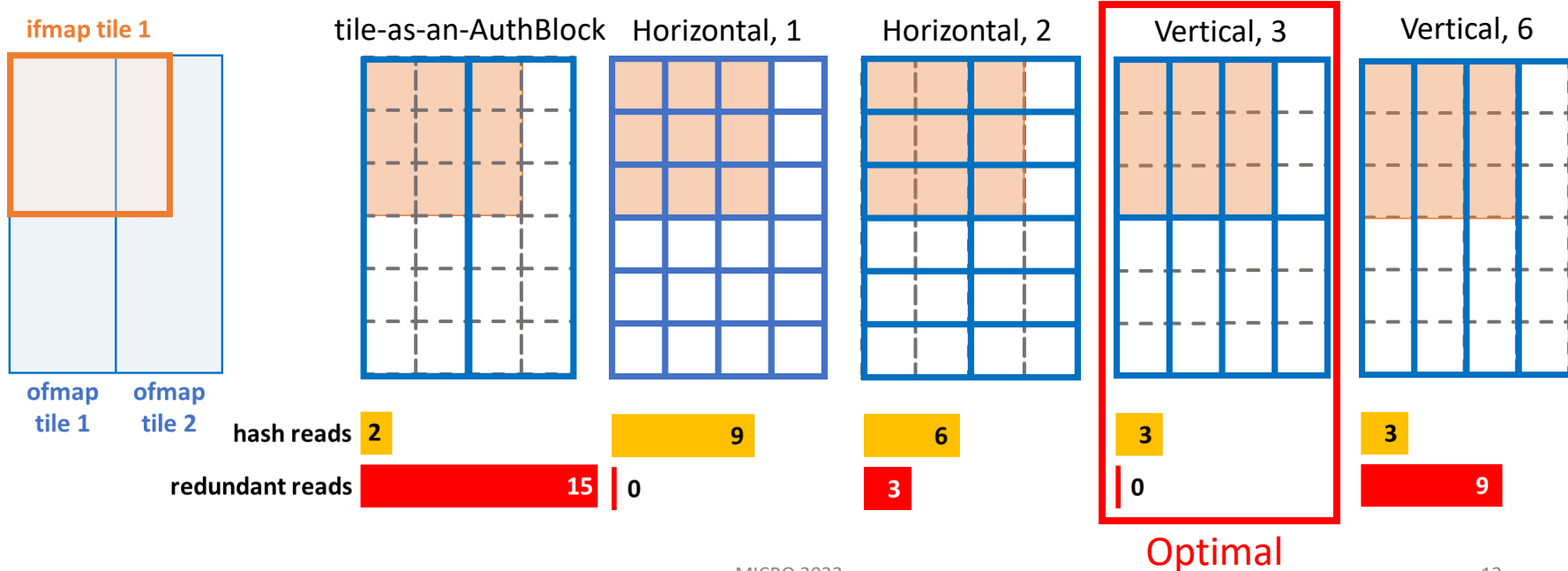
# Search space of AuthBlocks is complex

- Find the AuthBlock assignment that minimizes the additional off-chip traffic
- Both size and orientation of AuthBlocks matter



ifmap tile 1

ofmap tile 1    ofmap tile 2

tile-as-an-AuthBlock    Horizontal, 1    Horizontal, 2    Vertical, 3    Vertical, 6

hash reads: 2 | 9 | 6 | 3 | 3

redundant reads: 15 | 0 | 3 | 0 | 9

Optimal

# Search space of AuthBlocks is complex

- Find the AuthBlock assignment that minimizes the additional off-chip traffic
- Both size and orientation of AuthBlocks matter

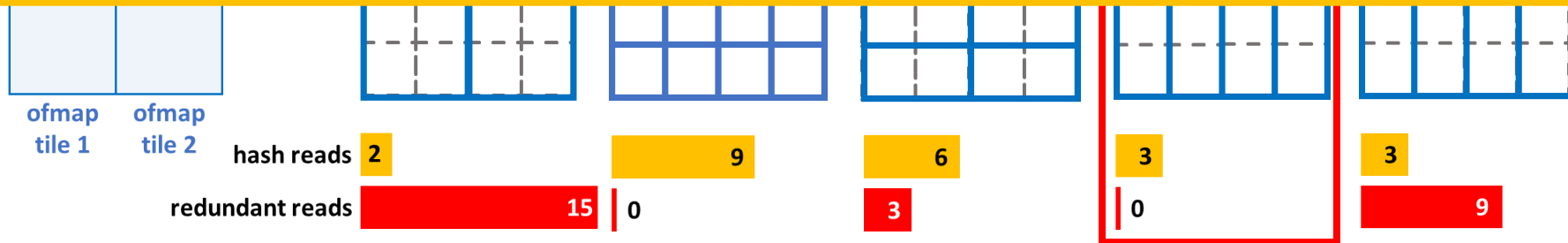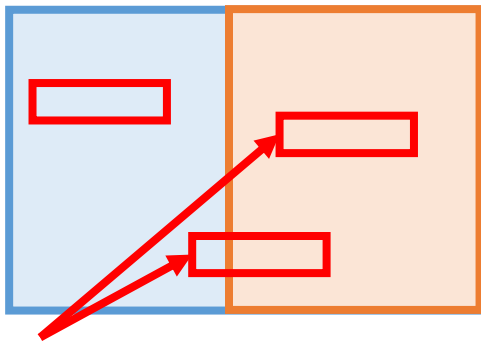ifmap tile 1    tile-as-an-AuthBlock    Horizontal, 1    Horizontal, 2    Vertical, 3    Vertical, 6

Exhaustive search using cycle-accurate simulation is time consuming

ofmap tile 1    ofmap tile 2

hash reads    2    9    6    3    3

redundant reads    15    0    3    0    9

Optimal

# Analytical approach to AuthBlock assignment

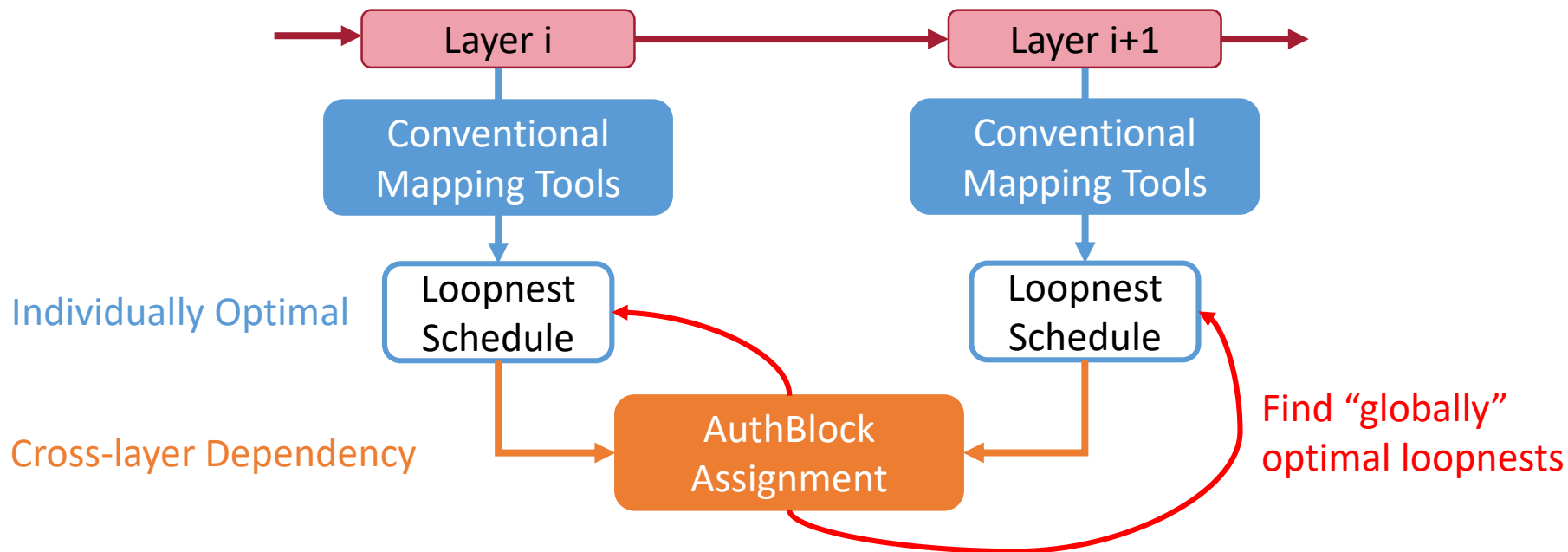- Counting how many AuthBlocks intersect a tile

Please refer to the paper for details

**Linear congruence problem**

$$u \times k \equiv \min(w_i - w_j - u + 1, 0)$$
$$, \ldots, w_i - 1 \bmod w_i$$

Count these AuthBlocks

# Cross-layer dependency from the loopnest level



**Loopnest schedules optimal for individual layer $\not\to$ globally optimal**

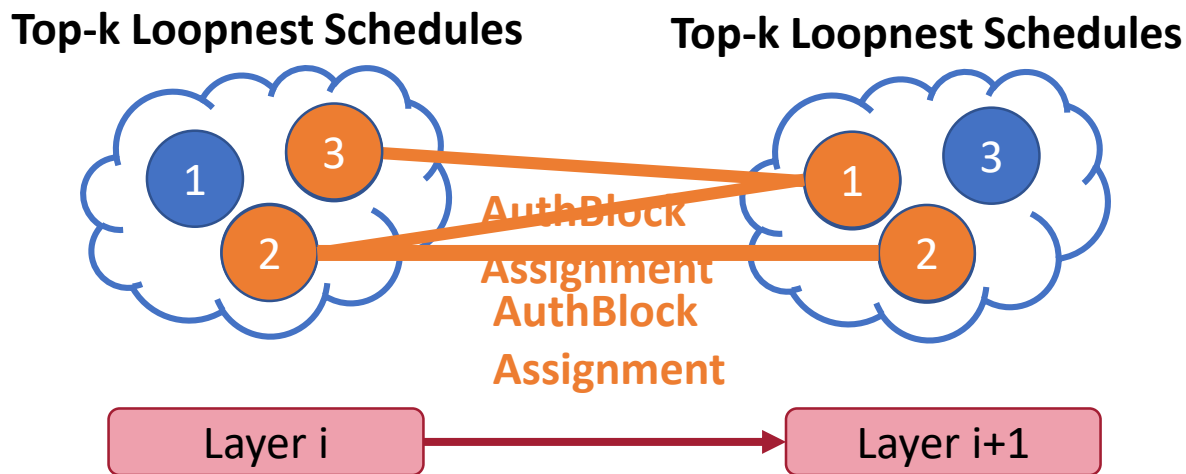# Cross-layer dependency from the loopnest level



**Loopnest schedules optimal for individual layer** $\not\to$ **globally optimal**

# Heuristic approach to joint optimization

- Simulated annealing to find the approximate solution



**Top-k Loopnest Schedules**

**Top-k Loopnest Schedules**

AuthBlock
Assignment
AuthBlock
Assignment

Layer i

Layer i+1

# SecureLoop: Scheduling Search Engine



**Step 1:**
**Augmenting Loopnest Schedulers**

DNN Workload

DNN Accel Arch.

CryptEngine Arch.

DNN Accelerator

Crypt Engine

DRAM

**Effective DRAM Model**

Loopnest Schedulers

**(top-k) Loopnest Schedules for Each Layer in a DNN**

# SecureLoop: Scheduling Search Engine



**Step 1:**
**Augmenting Loopnest Schedulers**

DNN Workload

DNN Accel Arch.

CryptEngine Arch.

DNN Accelerator

Crypt Engine

DRAM

**Effective DRAM Model**

Loopnest Schedulers

**(top-k) Loopnest Schedules for Each Layer in a DNN**

**Step 3:**
**Cross-layer Fine Tuning**

Layer i

Layer i+1

1 2 3

1 2 3
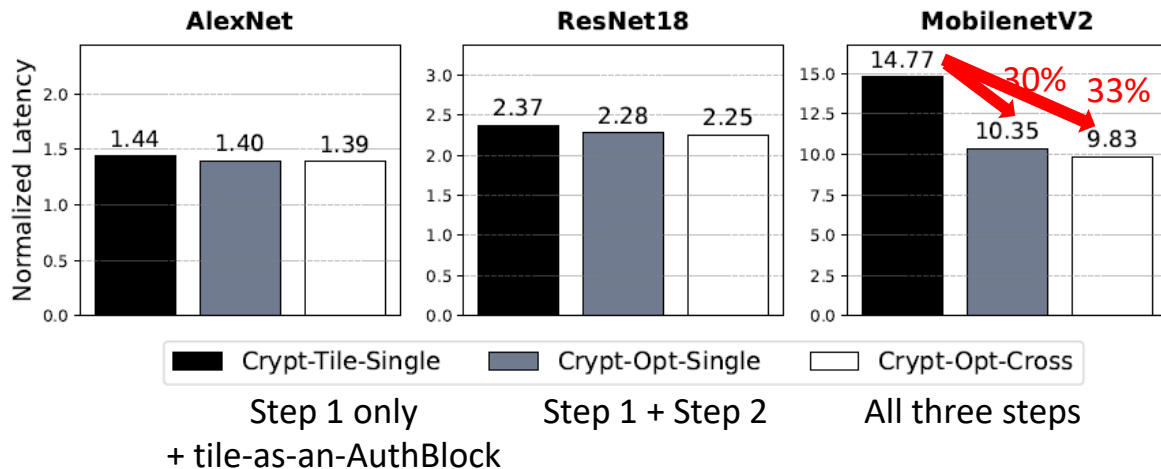
**Step 2:**
**Authentication Block Assignment**

# Comparing scheduling algorithms

**Setup**

- Eyeriss-like + AES-GCM
- Mostly conv workloads
- Different scheduling algo.

*shallow,
unavoidable rehashing*

*deeper, ↑ opportunity*



Step 1 only
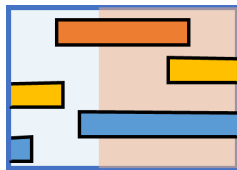+ tile-as-an-AuthBlock

Step 1 + Step 2

All three steps

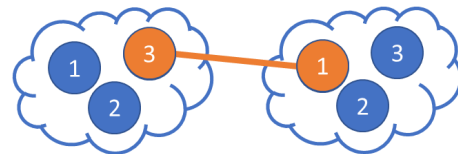**Summary:** ~33% faster, ~50% better in EDP compared to the "tile-as-an-AuthBlock"

# SecureLoop

- Scheduling algorithm considering cryptographic operations

Analytical approach to identify the optimal AuthBlock assignment



Cross-layer fine tuning from the loopnest schedule



- Design space exploration for secure accelerators

more in paper, including sweeps for different design choices and trade-off curves..