# ICCD 2025

# SecNPU: Securing LLM inference on NPU

**Xuanyao Peng[1,2]**, Yinghao Yang[1], Shangjie Pan[1,3], Junjie Huang[2], Yujun Liang[2], Hang Lu[1,3], Fengwei Zhang[2], Xiaowei Li[1,3]

[1]SKLP, Institute of Computing Technology, CAS

[2]Department of Computer Science and Engineering, SUSTech

[3]Zhongguancun Laboratory

# Summary

**The contributions of this work:**

1. Propose a CPU-decoupled TEE architecture for LLM inference – **SecNPU.**
2. Propose an near-zero-overhead secure startup mechanism for LLMs.
3. Implement the prototype based on RTL design and evaluate its performance using a cycle-accurate NPU simulator.

**Benefits:**

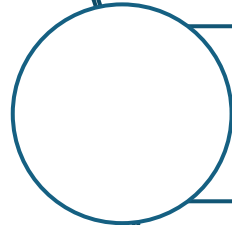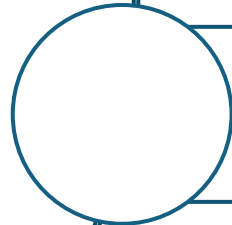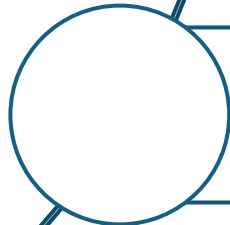| ① Broad compatibility | Use unified security metadata and apply to various kinds of CPU. |
| ② High performance | 1.6x speedup for LLM startup and 1.5x speedup for LLM decoding. |
| ③ Strict security guarantee | Protect from malicious OS and hardware attacks. |

# OUTLINE

**Threat Model & Motivation**

Methodology

Evaluation

Recap

# Threat Model

**Security Threats** in **CPU-NPU** Heterogeneous Systems:

**User's Privacy:**

- **Confidential user prompts**
- **Private user data**

**Model's Parameters:**

- **Data poisoning attacks**
- **Theft of model weights**



Image sources: hiddenlayer.com; forbes.com

# Threat Model

**Security Threats** in **CPU-NPU** Heterogeneous Systems:
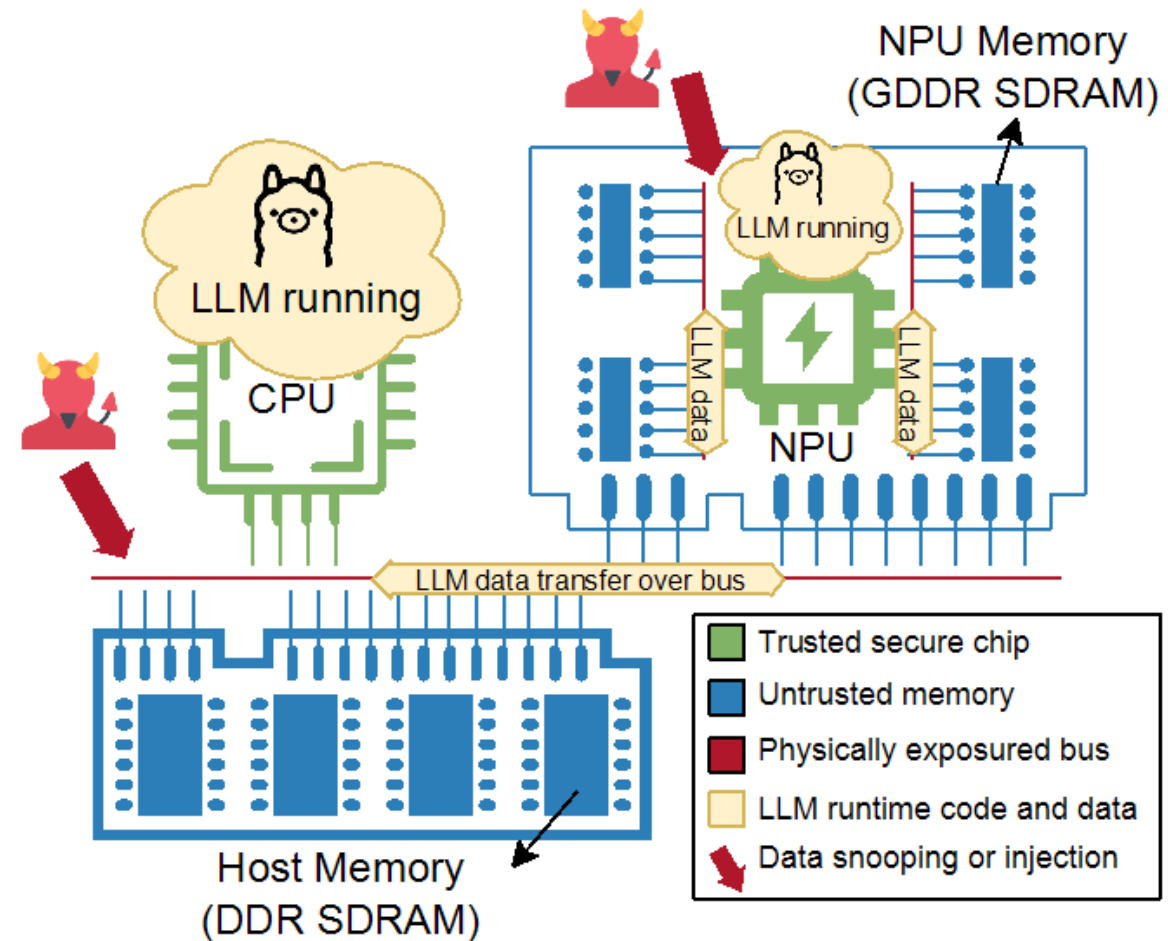
**User's Privacy:**

- **Confidential prompts**

- **Private inputs**

**Model's Parameters:**

- **Data poisoning**

- **Steal confidential weights**

**Inputs and model parameters** can be transmitted to the NPU by a malicious OS, exposing the data **on the physical bus**.
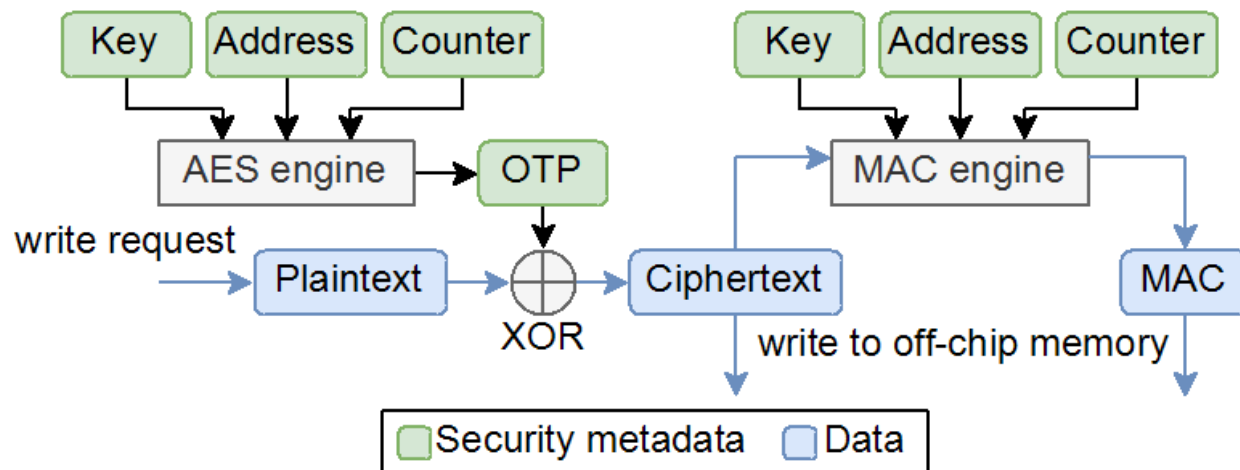


NPU Memory (GDDR SDRAM)

LLM running

LLM running

CPU

LLM data

LLM data

NPU

LLM data transfer over bus

Host Memory (DDR SDRAM)

Trusted secure chip
Untrusted memory
Physically exposed bus
LLM runtime code and data
Data snooping or injection

# Motivation

**The security mechanisms of traditional TEE (Trusted Execution Environment):**

- **Encrypt plaintext using AES-GCM**

- **Protect ciphertext's integrity using MAC**
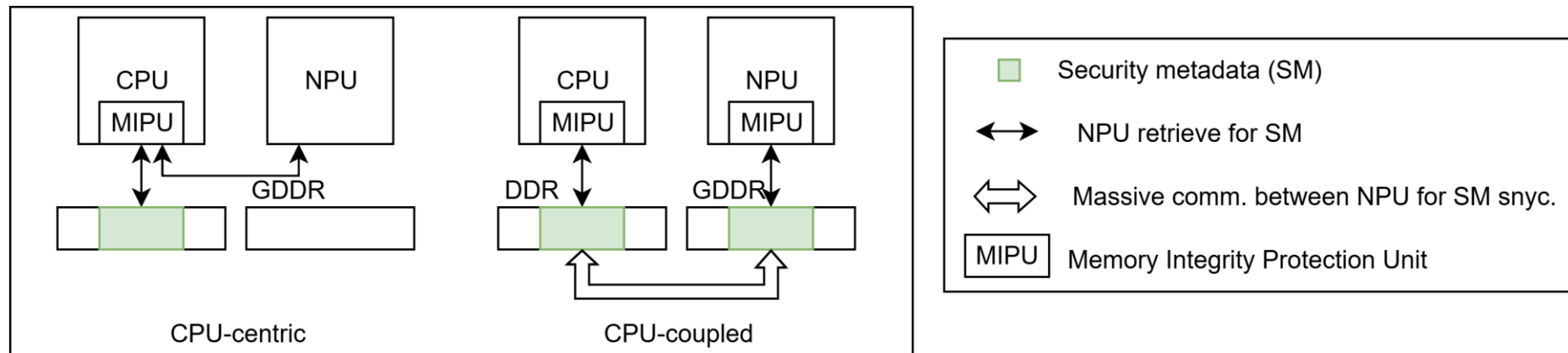
    **(Message Authentication Code)**



**3 types of security metadata:**

- **Private key stored within the Root of Trust (RoT) on chip**

- **Physical address of data**

- **Counter to ensure data freshness**

# Motivation

The traditional CPU-NPU TEE can be **classified in two categories**:

- **CPU-centric: All security functions (AES-GCM/MAC) are handled by the CPU (e.g. TNPU HPCA'22)**

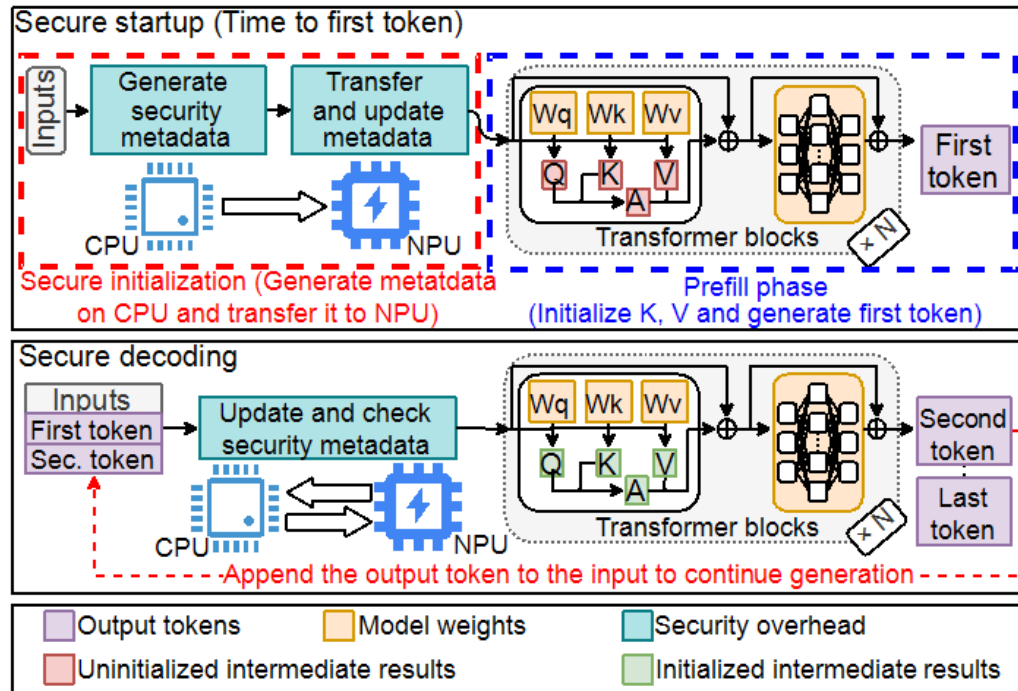- **CPU-coupled: Part of security functions (MAC) are delegated to the CPU (e.g. TensorTEE ASPLOS'24)**



**Both face:**

- **Slow startup due to security metadata initialization and transmission**
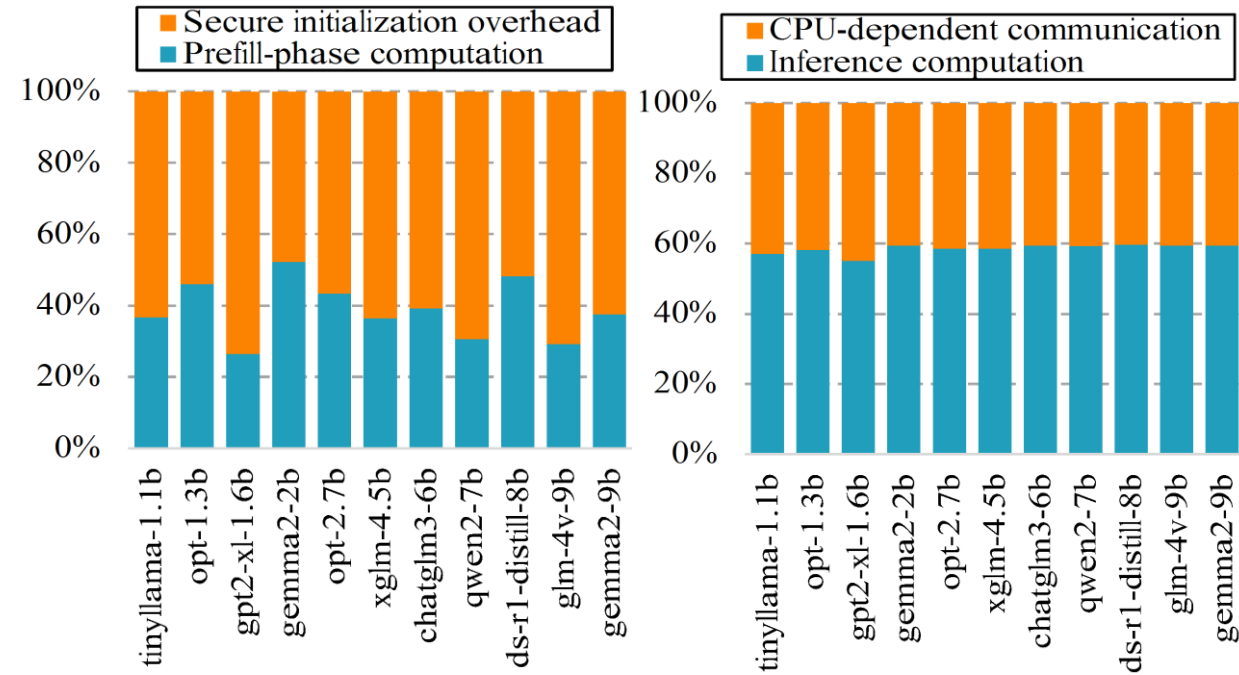- **High communication overhead during LLM inference**

# Motivation

## The security overhead of CPU-centric/coupled TEE is significant



**The security overhead introduced:**

- **Secure startup**
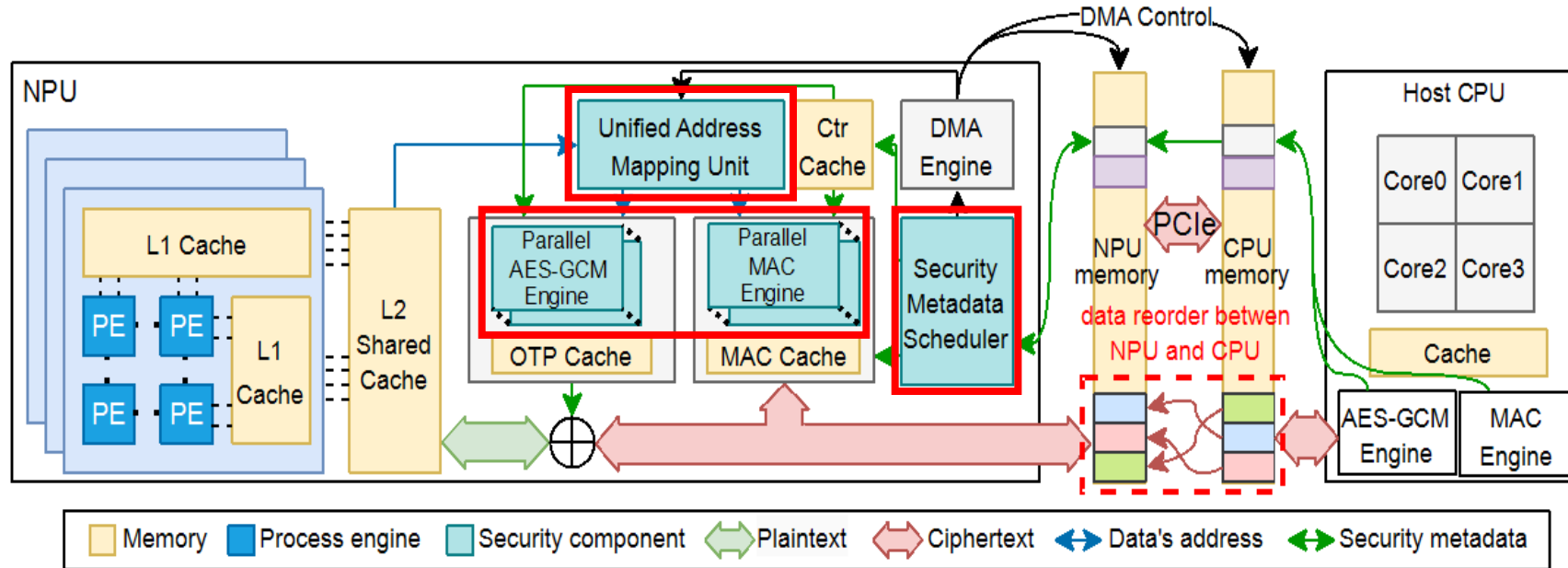- **Secure decoding**

**Startup overhead: 60% (average)**

**Inference overhead: 40% (average)**

# OUTLINE

Threat Model & Motivation

**Methodology**

Evaluation

Recap

**The overview architecture of SecNPU: CPU-decoupled TEE**



**3 key security components introduced:**

- **Unified Address Mapping Unit**: Handles data remapping after transfer from the CPU
- **Parallel AES-GCM/MAC Engine**: Accelerates NPU encryption & integrity verification
- **Security Metadata Scheduler**: Mitigates security overhead during startup

**Towards Unified Security Metadata and Near-Zero-Overhead Secure Startup!**

# Methodology

**Unified Security Metadata: Unified Physical Address**



**Data is reordered on NPU-side to improve prefetch accuracy.**

**Security metadata generated for old addresses becomes invalid!**

**How to maintain security metadata?**
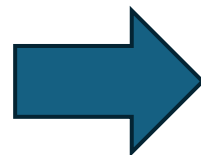
**Unified Security Metadata: Unified** **Physical Address**



**Design a dynamic mapping table:**
- **Original order**
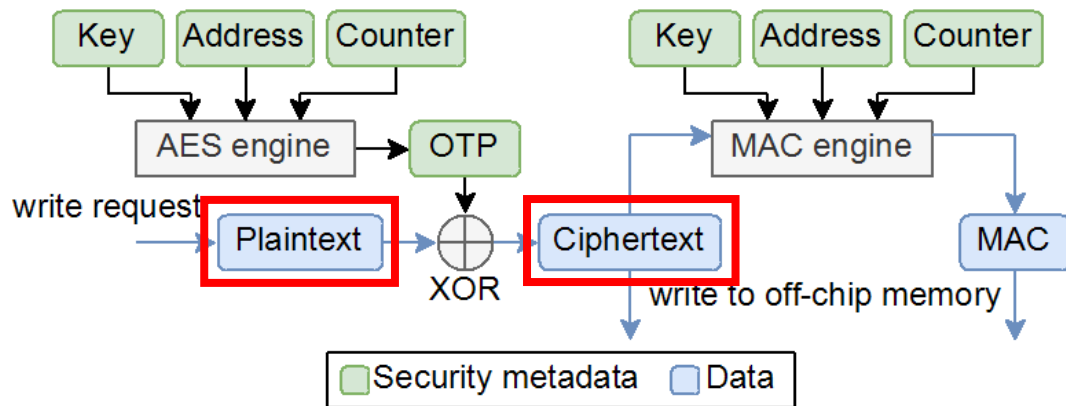- **Reorder granularity**
- **Original base physical address**

→ **Original physical address**

## Unified Security Metadata: Unifying Memory Protection Granularity

**Access Granularity Mismatch:**

- **CPU: 64-byte**
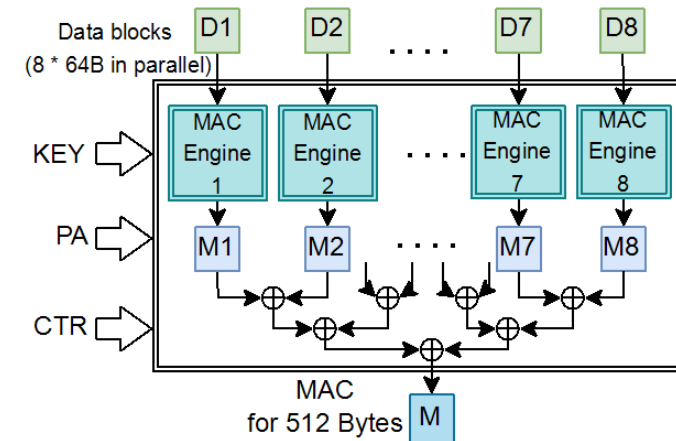
- **NPU: Larger blocks with DMA (Vendor-dependent)**



**CPU:**

**NPU:**

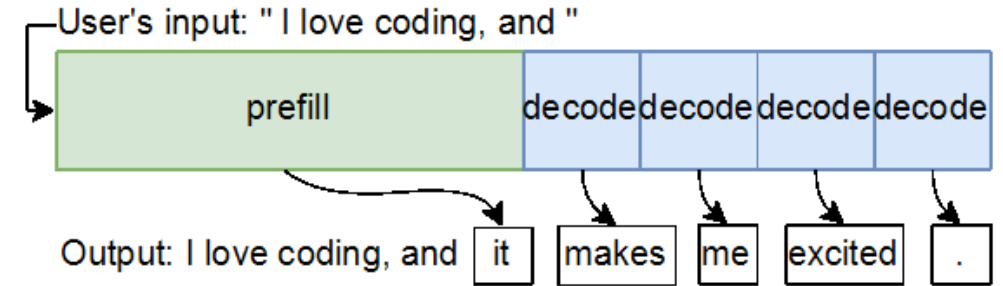**The CPU and NPU use different plaintext/ciphertext block sizes.**

**Unified to NPU-side granularity (CPU performs software-based metadata alignment to match NPU)**
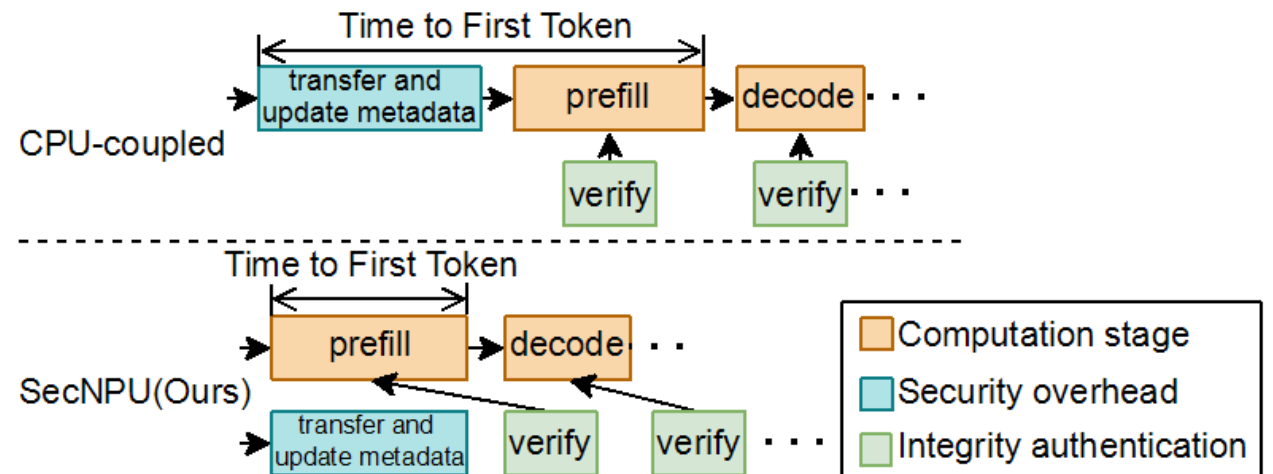
**Near-Zero-Overhead Startup: LLM-Oriented Optimization**

**The two stages of LLM inference**

- **Prefill: Compute-intensive stage**
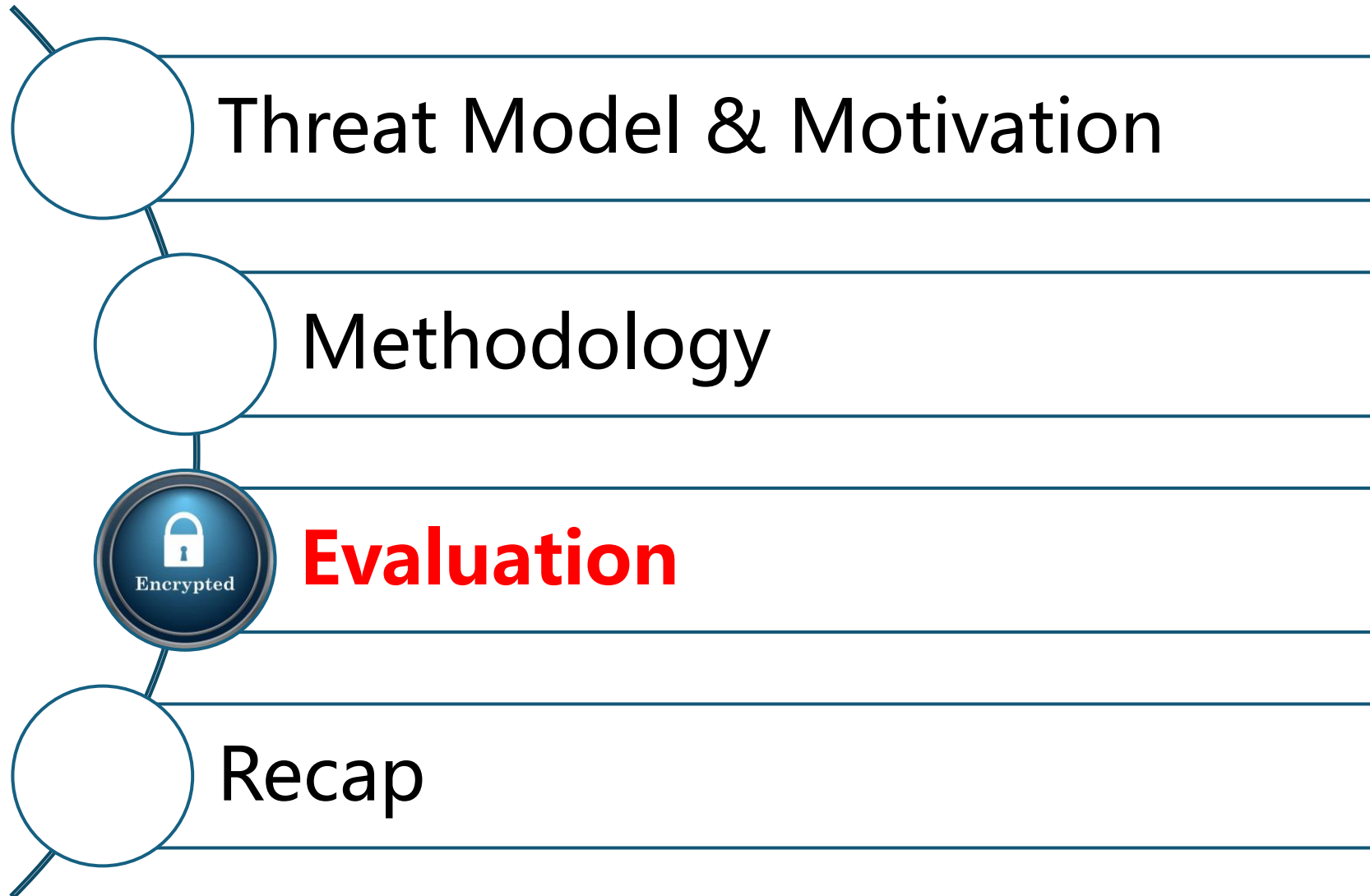
- **Decode: Memory-intensive stage**



This results in **the prefill phase occupying significantly less memory bandwidth** than the decode phase.

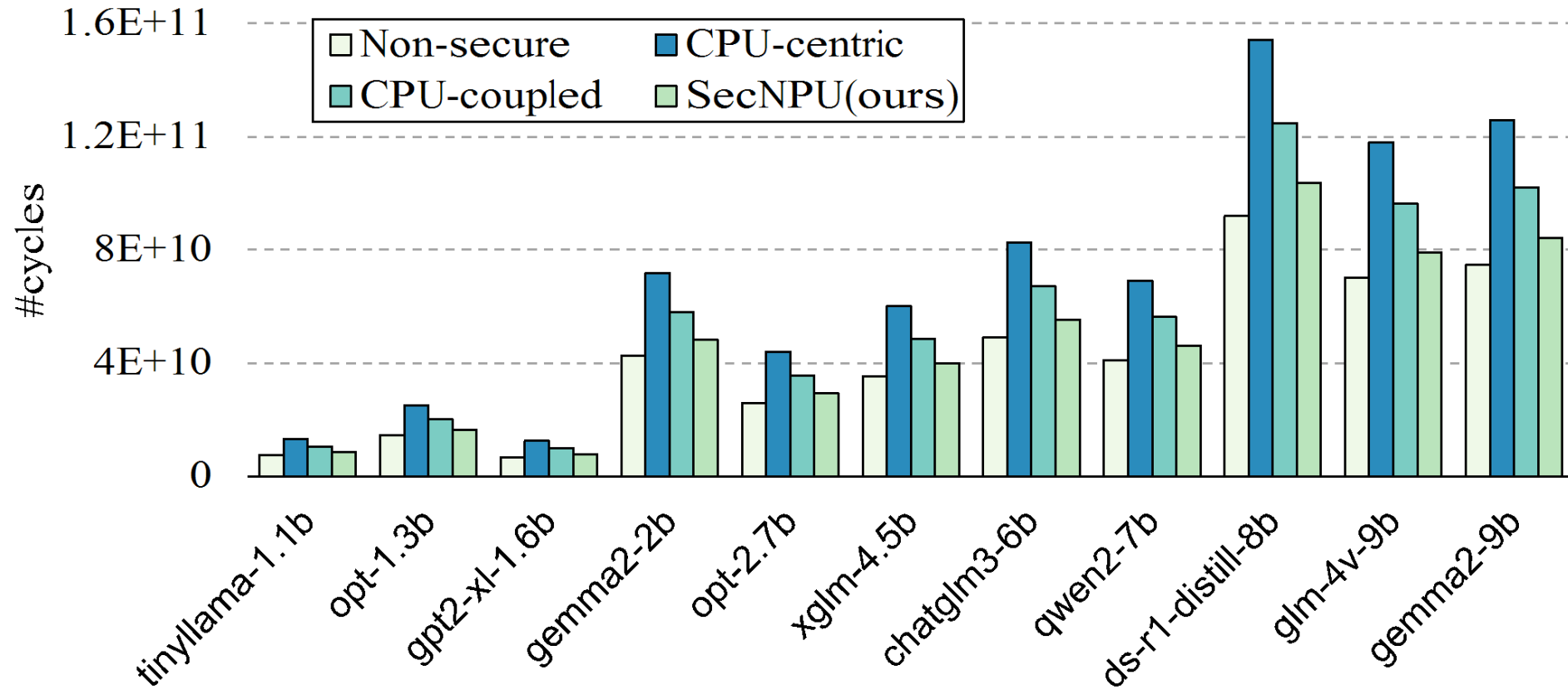Transfer security metadata during the prefill stage to eliminate overhead!

# OUTLINE

Threat Model & Motivation

Methodology

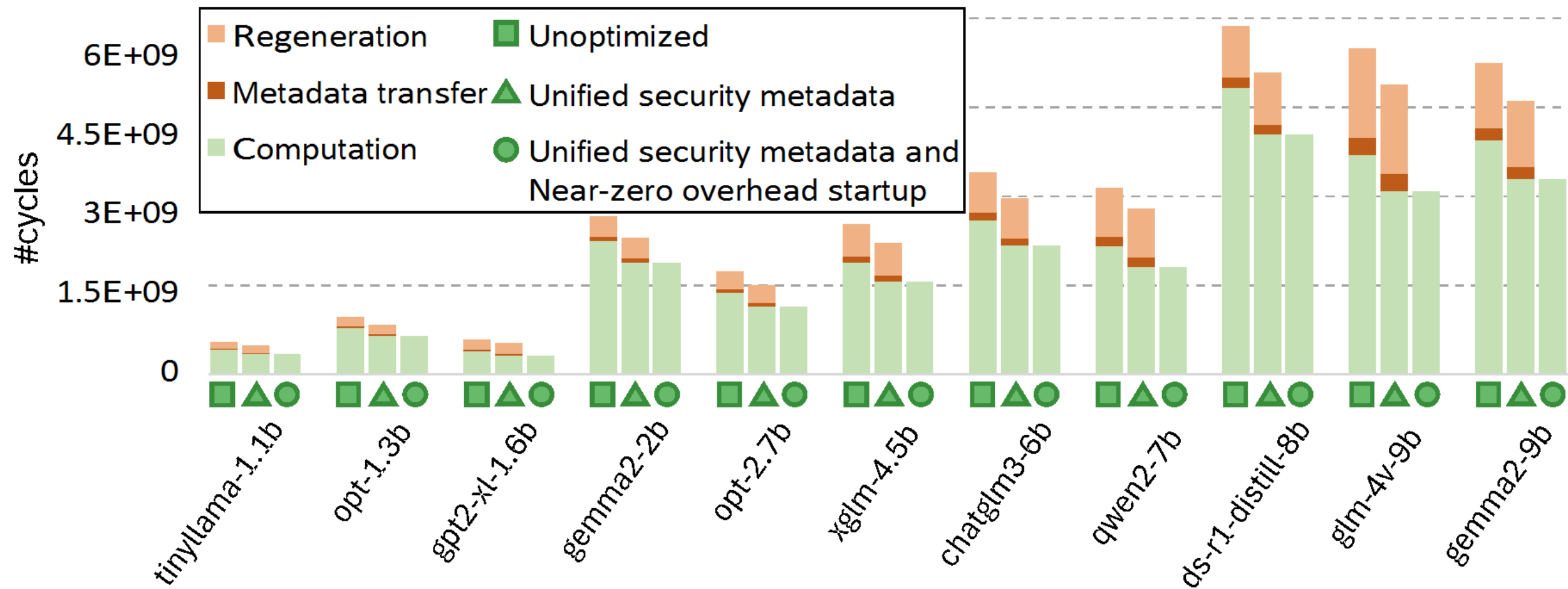**Evaluation**

Recap

# Evaluation

CPU-centric：TNPU
CPU-coupled：TensorTEE



The index measures the total cycles required; a lower value is better.
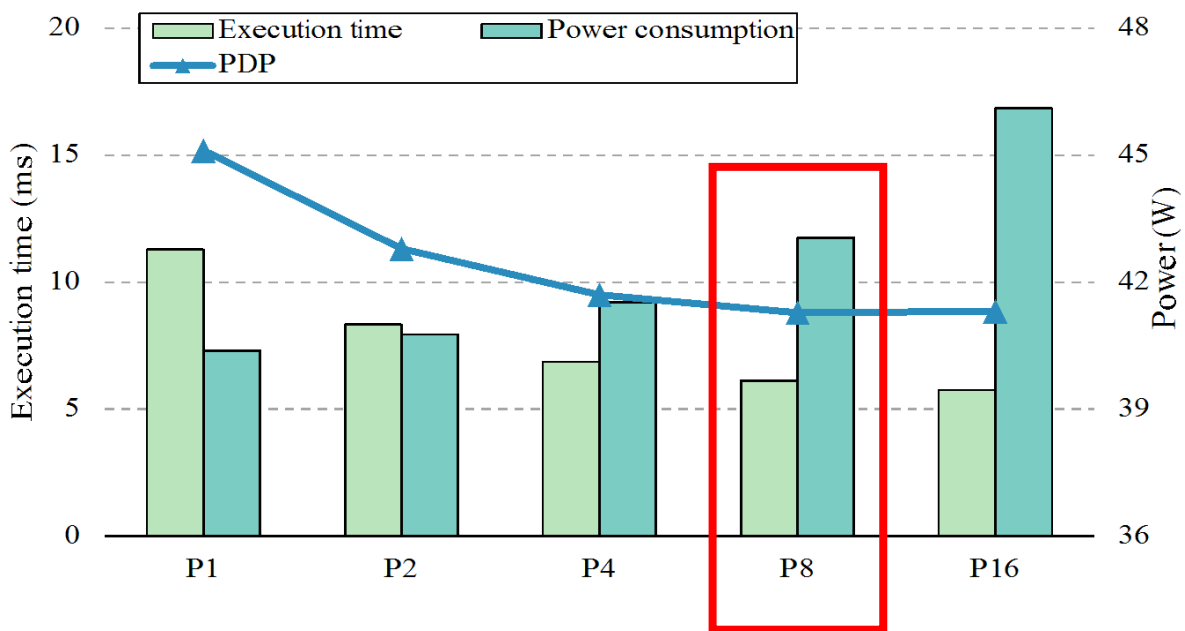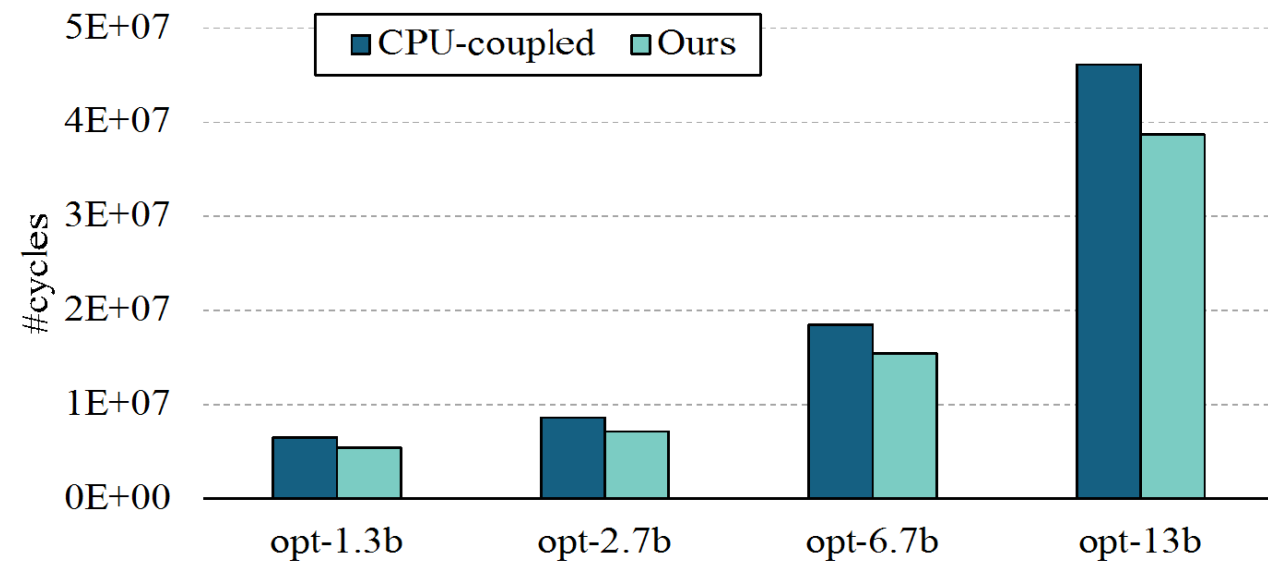
# Evaluation

## Ablation Study

# Evaluation

## Design Space Exploration



## Sensitivity Analysis of Multi-size LLMs

# Recap

**The contributions of this work:**

1. **Our work, SecNPU, proposes a CPU-decoupled TEE architecture:**
   - **unified security metadata**
   - **near-zero-overhead startup**
2. **Our prototype demonstrates speedups of up to 1.6x, all while providing robust security guarantees against both OS and hardware attacks**