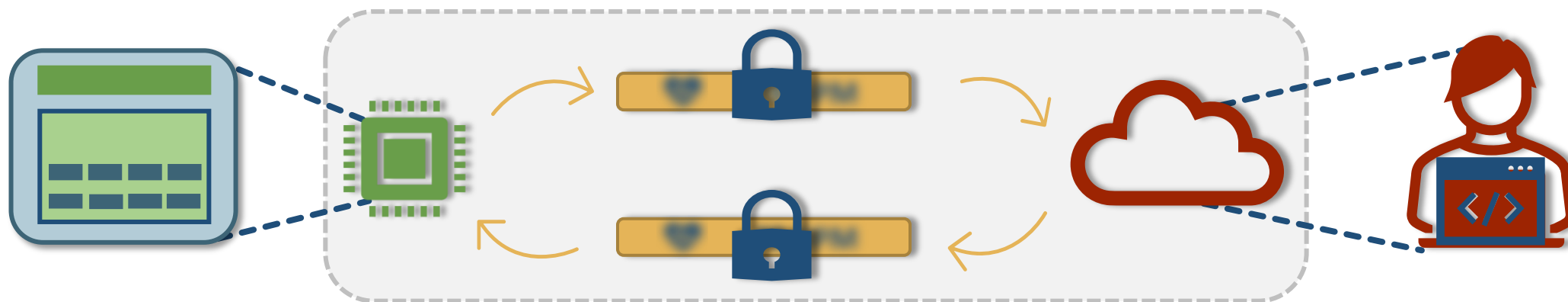# Client-Optimized Algorithms & Acceleration for Encrypted Compute Offloading

**ASPLOS '22 | Lausanne, Switzerland | March 3, 2022**

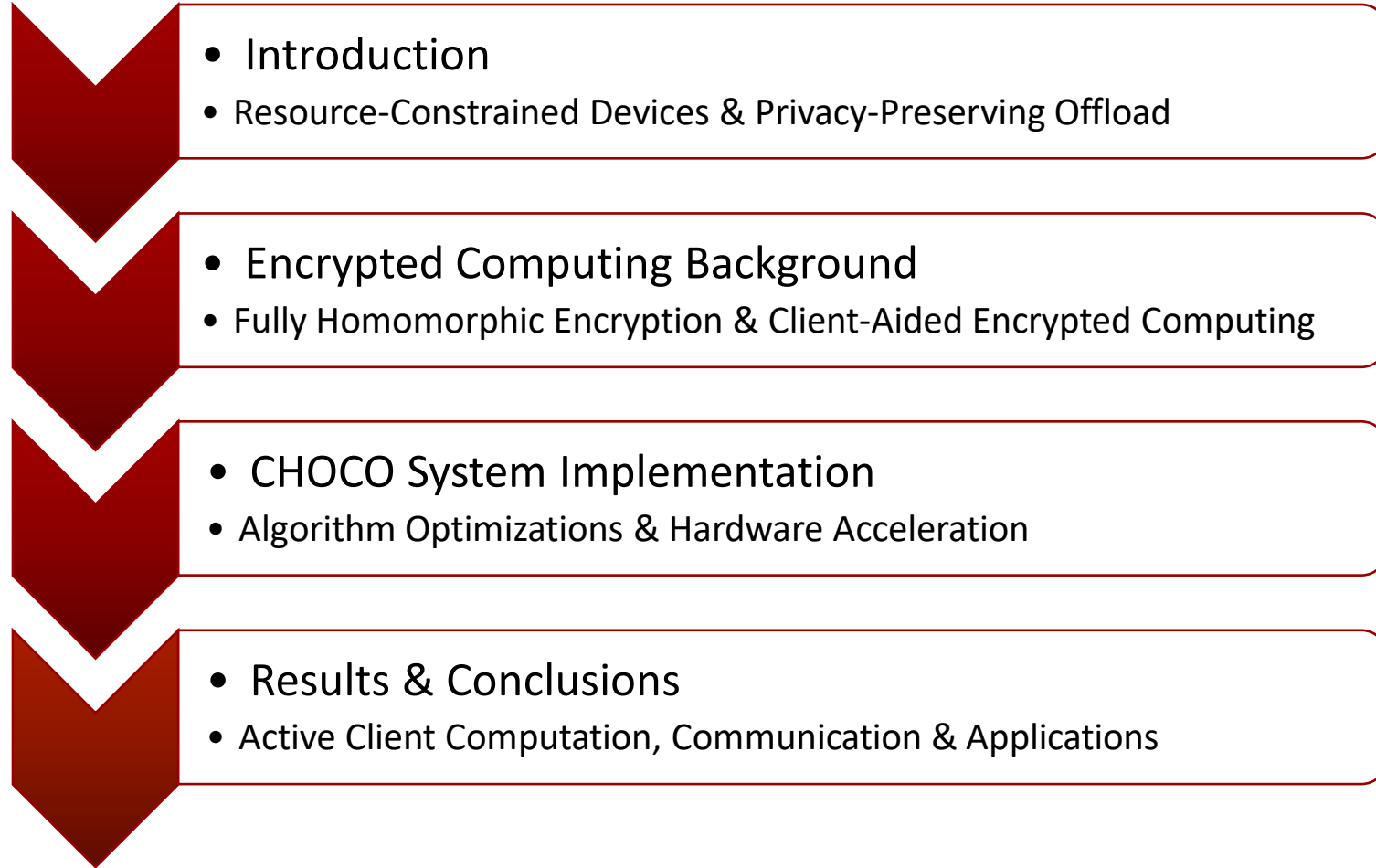McKenzie van der Hagen – mckenziv@andrew.cmu.edu

Brandon Lucia – blucia@andrew.cmu.edu

Electrical & Computer ENGINEERING

# **C**lient-aided **H**E for **O**paque **C**ompute **O**ffloading



- CHOCO enables privacy-preserving computation for **resource-constrained devices**

- CHOCO utilizes **Homomorphic Encryption (HE)** and Client-Aided Encrypted Computing

- CHOCO introduces client-optimized **encrypted algorithms & hardware acceleration**

- CHOCO makes client responsibility **competitive with local compute**
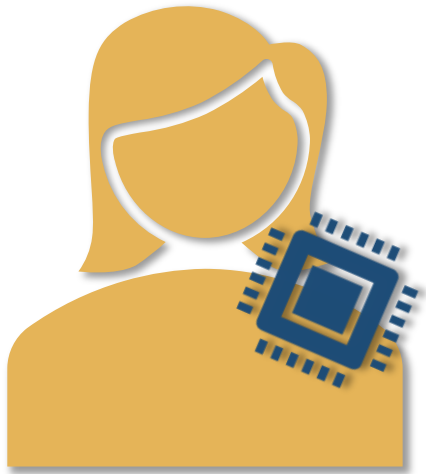
- CHOCO benefits **generalize to diverse applications**

Electrical **&** Computer
ENGINEERING

# Outline

- Introduction
  - Resource-Constrained Devices & Privacy-Preserving Offload

- Encrypted Computing Background
  - Fully Homomorphic Encryption & Client-Aided Encrypted Computing

- CHOCO System Implementation
  - Algorithm Optimizations & Hardware Acceleration

- Results & Conclusions
  - Active Client Computation, Communication & Applications
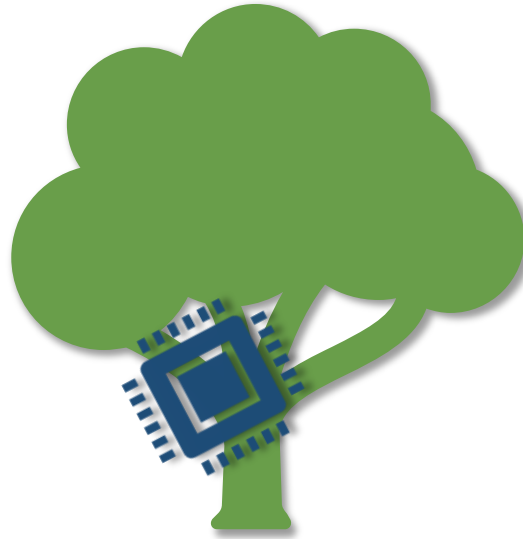
Electrical & Computer
ENGINEERING

# Outline

- **Introduction**
  - **Resource-Constrained Devices & Privacy-Preserving Offload**

- Encrypted Computing Background

- CHOCO System Implementation

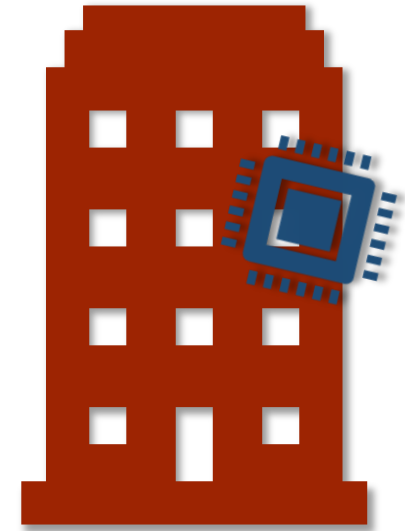- Results & Conclusions

Electrical & Computer
ENGINEERING

# Resource-Constrained Devices are Everywhere

Health Monitoring

Wildlife Monitoring

Infrastructure Monitoring

# Computational Demands Exhaust Sensor Resources



114 mg/dL

109 / 72

97% SpO2

78 BPM

Device Size

Data Quantity

Computation Complexity

Electrical & Computer ENGINEERING

# Privacy-Preserving Computation Offload



**Communication**
Transfer Data & Results

**Shared Offload Server**
Semi-Honest Central Compute

**Resource-Constrained Private Client**
Secure Data Collection

78 BPM

# Privacy-Preserving Computation Offload



**Communication**
Transfer Data & Results

**Shared Offload Server**
Semi-Honest Central Compute

**Resource-Constrained Private Client**
Secure Data Collection

**FHE: [Fully] Homomorphic Encryption**

# Outline

- **Introduction**
  - **Resource-Constrained Devices & Privacy-Preserving Offload**

- Encrypted Computing Background

- CHOCO System Implementation

- Results & Conclusions
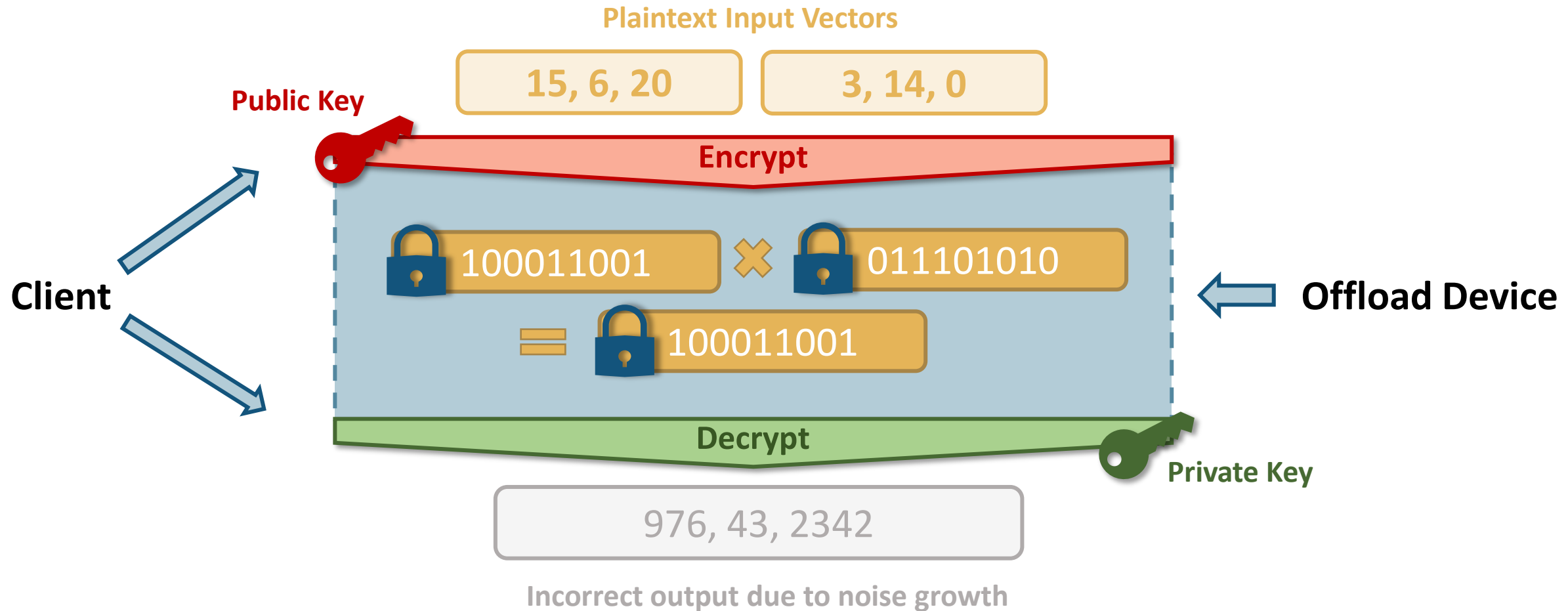
Electrical & Computer
ENGINEERING

# Outline

- Introduction

- **Encrypted Computing Background**
  - **Fully Homomorphic Encryption & Client-Aided Encrypted Computing**
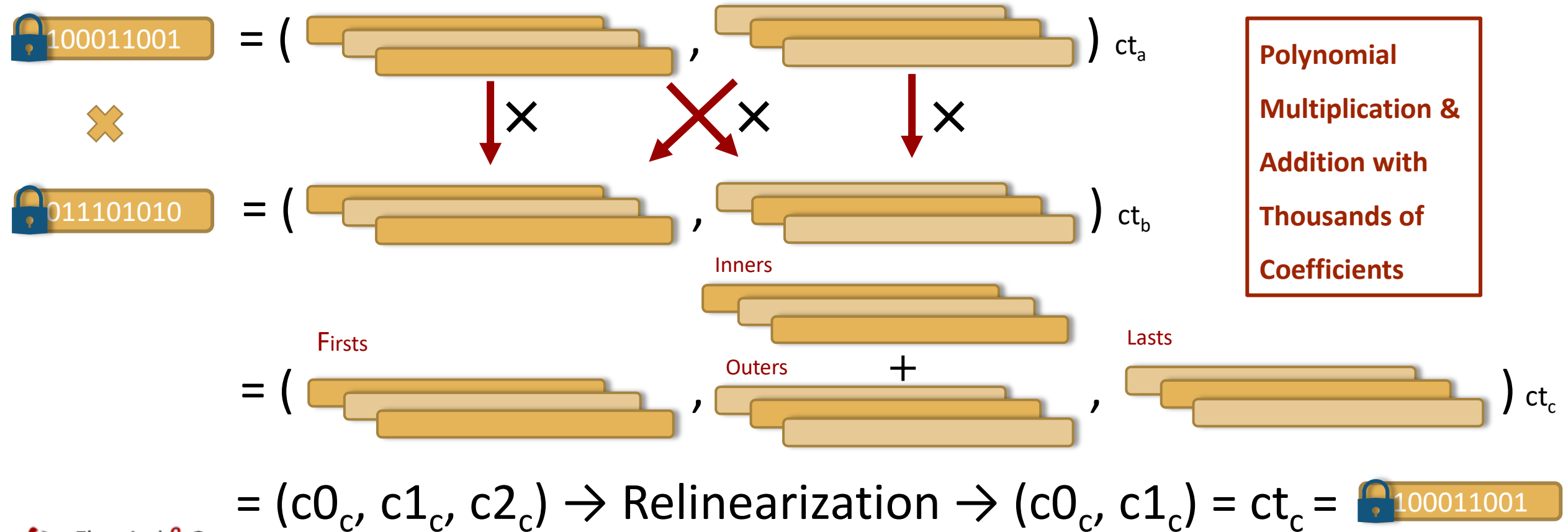
- CHOCO System Implementation

- Results & Conclusions

Electrical & Computer ENGINEERING

# Polynomial Operations
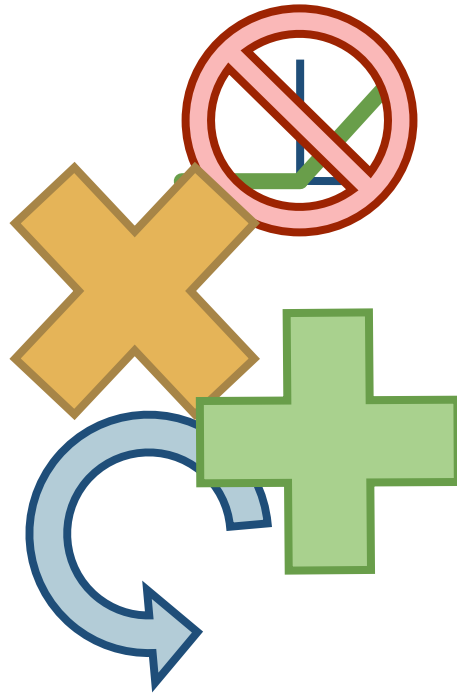
$$ct_a \times ct_b = (c0_a, c1_a) \times (c0_b, c1_b)$$

$$= ([c0_a \times c0_b]_q , [c0_a \times c1_b + c1_a \times c0_b ]_q , [c1_a \times c1_b]_q )$$



**Polynomial Multiplication & Addition with Thousands of Coefficients**

= (c0_c, c1_c, c2_c) → Relinearization → (c0_c, c1_c) = ct_c =

Electrical & Computer ENGINEERING

# HE Challenges & Limitations

High Computation Costs

Linear Operations

Noise Growth & Arithmetic Depth

Parameter Selection

# HE Challenges & Limitations

**Offloading Entire Encrypted Applications is Still Infeasible in Many Scenarios**
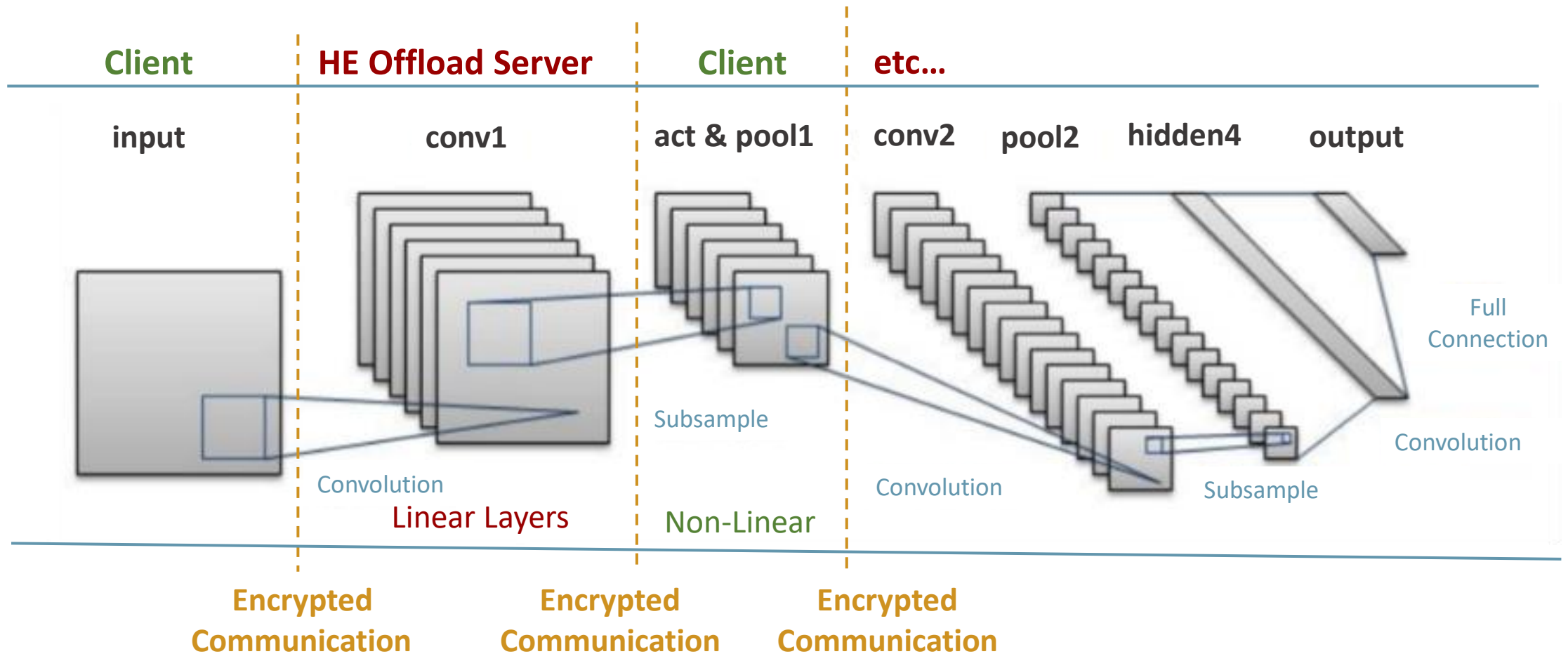
High Computation Costs       Linear Operations       Noise Growth       Parameter Selection

Electrical & Computer
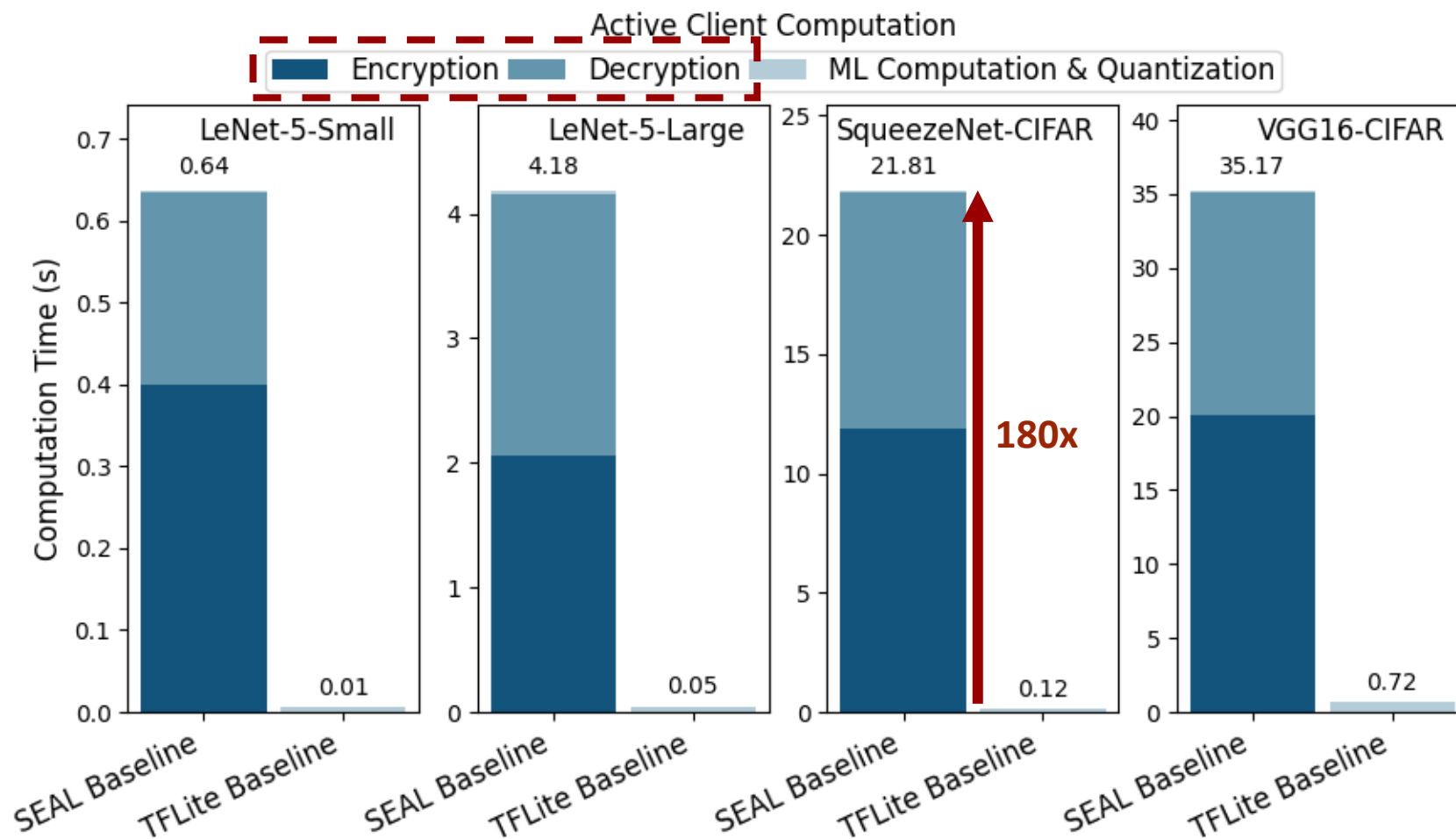ENGINEERING

14

# Client-Aided Encrypted Computing

# Client-Aided Encrypted Inference



Image: https://www.pyimagesearch.com/2016/08/01/lenet-convolutional-neural-network-in-python/

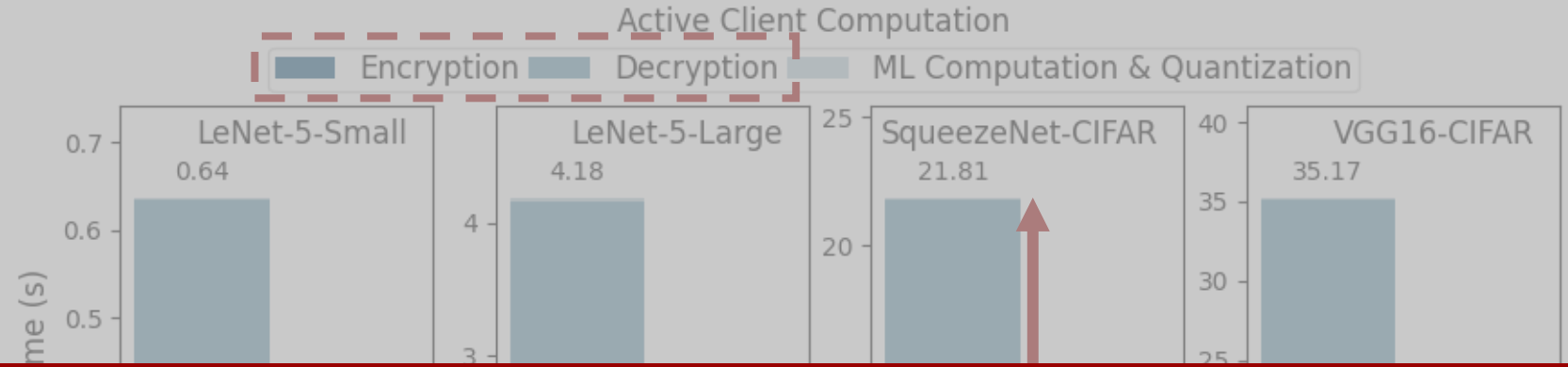- **Systematically limits arithmetic depth & regularly refreshes noise**

# Quantifying Client Responsibility

- ARM Cortex-A7 CPU Client

- **Up to 180x** overhead to offload compute

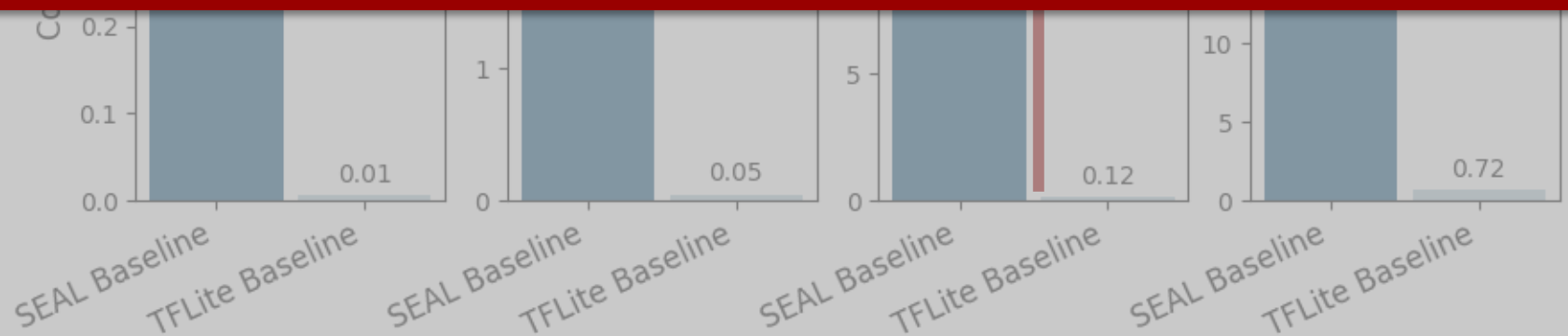- Dominated by Homomorphic Encryption (HE) operations

Electrical & Computer ENGINEERING

# Quantifying Client Responsibility

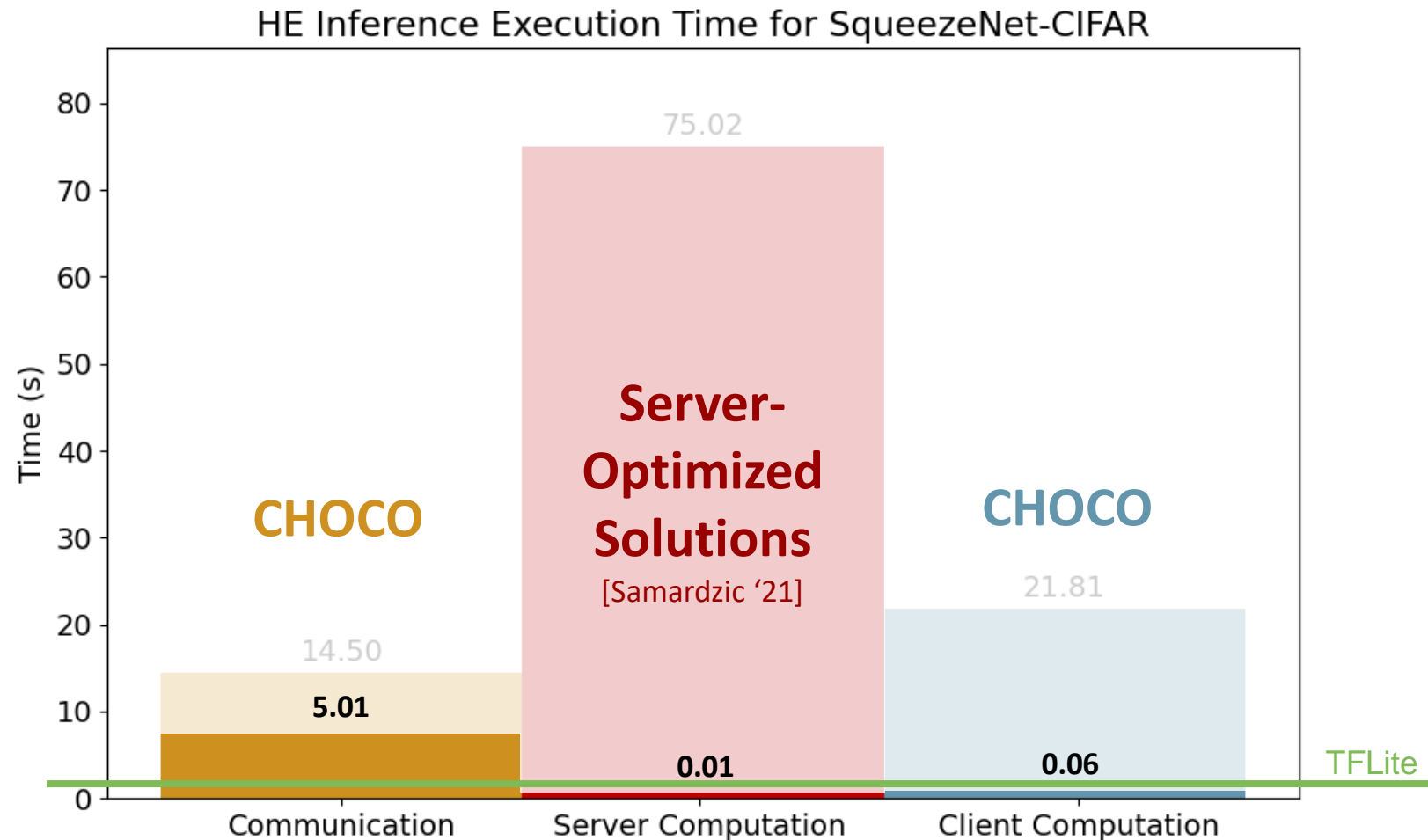- ARM Cortex-A7 CPU Client

- **Up to 180x** overhead to

  Encryption (HE) operations



**CHOCO Reduces Client-Side Computation by up to 341x through SW Algorithms & HW Acceleration**

# Complete Client-Aided System Improvements



HE Inference Execution Time for SqueezeNet-CIFAR

# Outline

- Introduction

- **Encrypted Computing Background**
  - **Fully Homomorphic Encryption & Client-Aided Encrypted Computing**

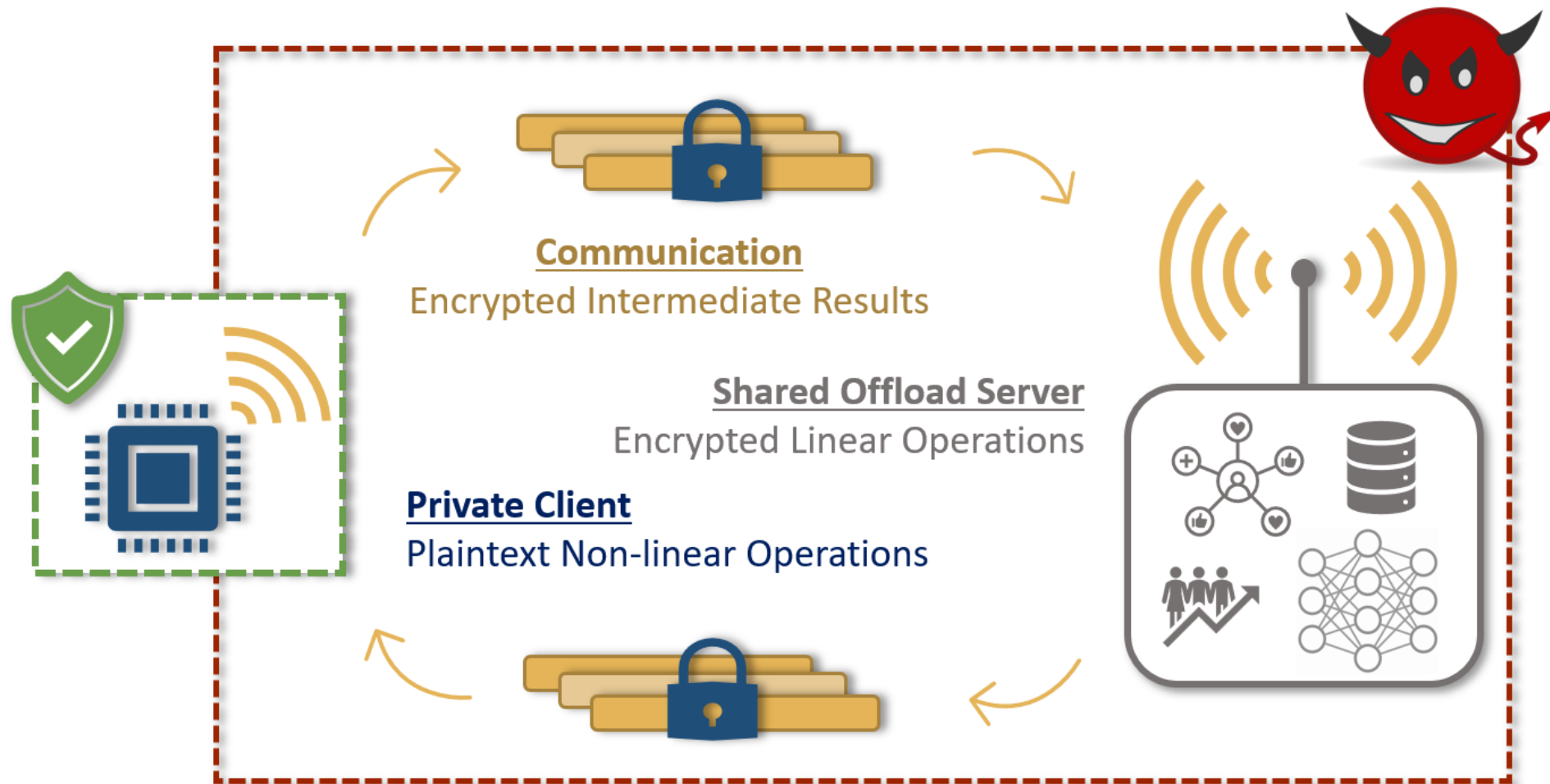- CHOCO System Implementation

- Results & Conclusions

Electrical & Computer
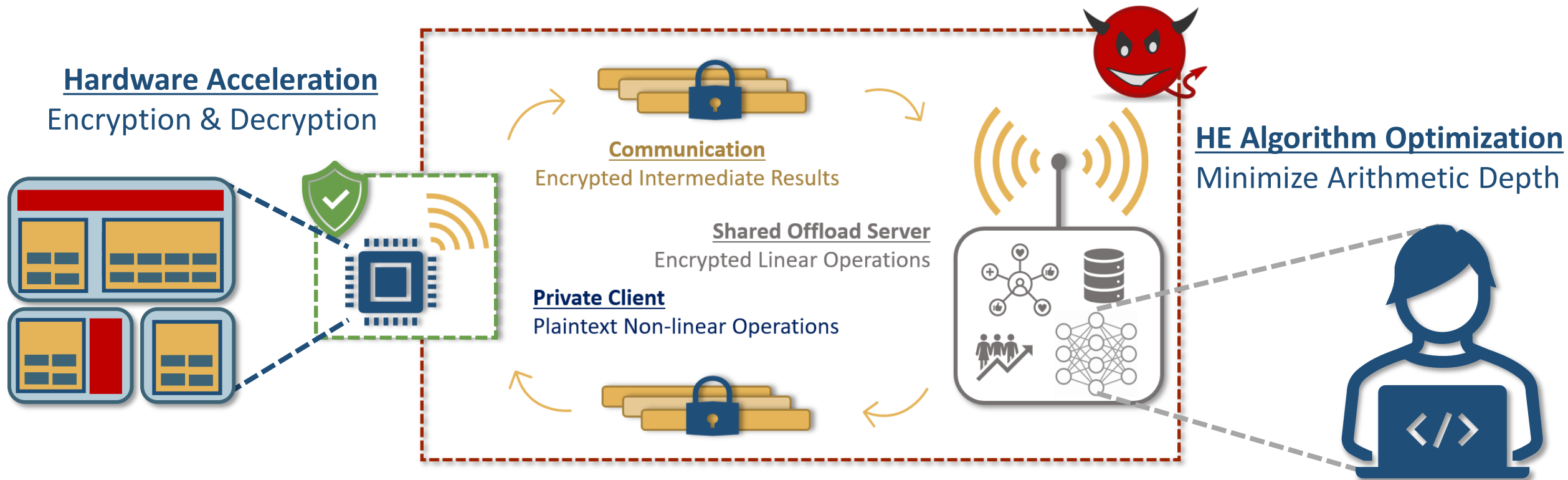ENGINEERING

# Outline

- Introduction

- Encrypted Computing Background

- **CHOCO System Implementation**
  - **Algorithm Optimizations & Hardware Acceleration**

- Results & Conclusions

Electrical & Computer
ENGINEERING

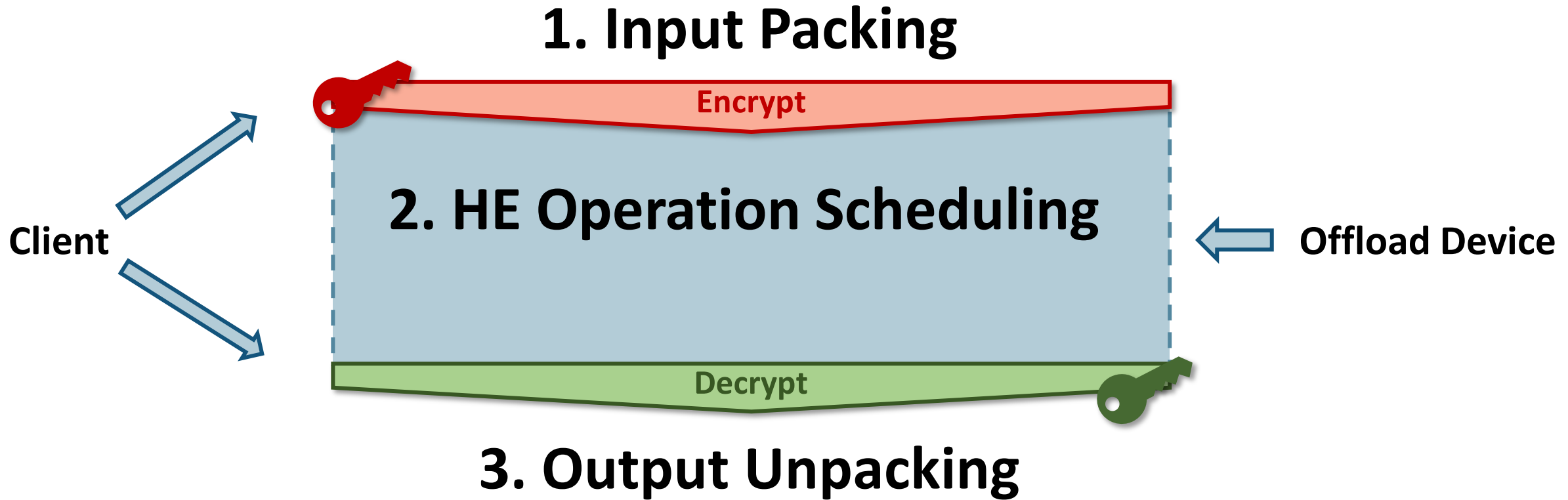# **C**lient-aided **H**E for **O**paque **C**ompute **O**ffloading



**Communication**
Encrypted Intermediate Results

**Shared Offload Server**
Encrypted Linear Operations

**Private Client**
Plaintext Non-linear Operations

22

# **C**lient-aided **H**E for **O**paque **C**ompute **O**ffloading

**Hardware Acceleration**
Encryption & Decryption

**HE Algorithm Optimization**
Minimize Arithmetic Depth

**Communication**
Encrypted Intermediate Results

**Shared Offload Server**
Encrypted Linear Operations

**Private Client**
Plaintext Non-linear Operations

# Encrypted Algorithm Optimization

# HE Algorithms

**1. Input Packing**

**Encrypt**

**2. HE Operation Scheduling**

**Client**

**Offload Device**

**Decrypt**

**3. Output Unpacking**

# Algorithm Optimizations Impact Client Costs



Minimize Arithmetic Depth → ↓ Noise Growth → ↓ Parameter Selections → ↓ Ciphertext Size → ↓ Computation & Mem / ↓ Communication

# Windowed Rotations

# Windowed Rotations



Ideally: Values wrap around within a window of interest

# Windowed Rotations



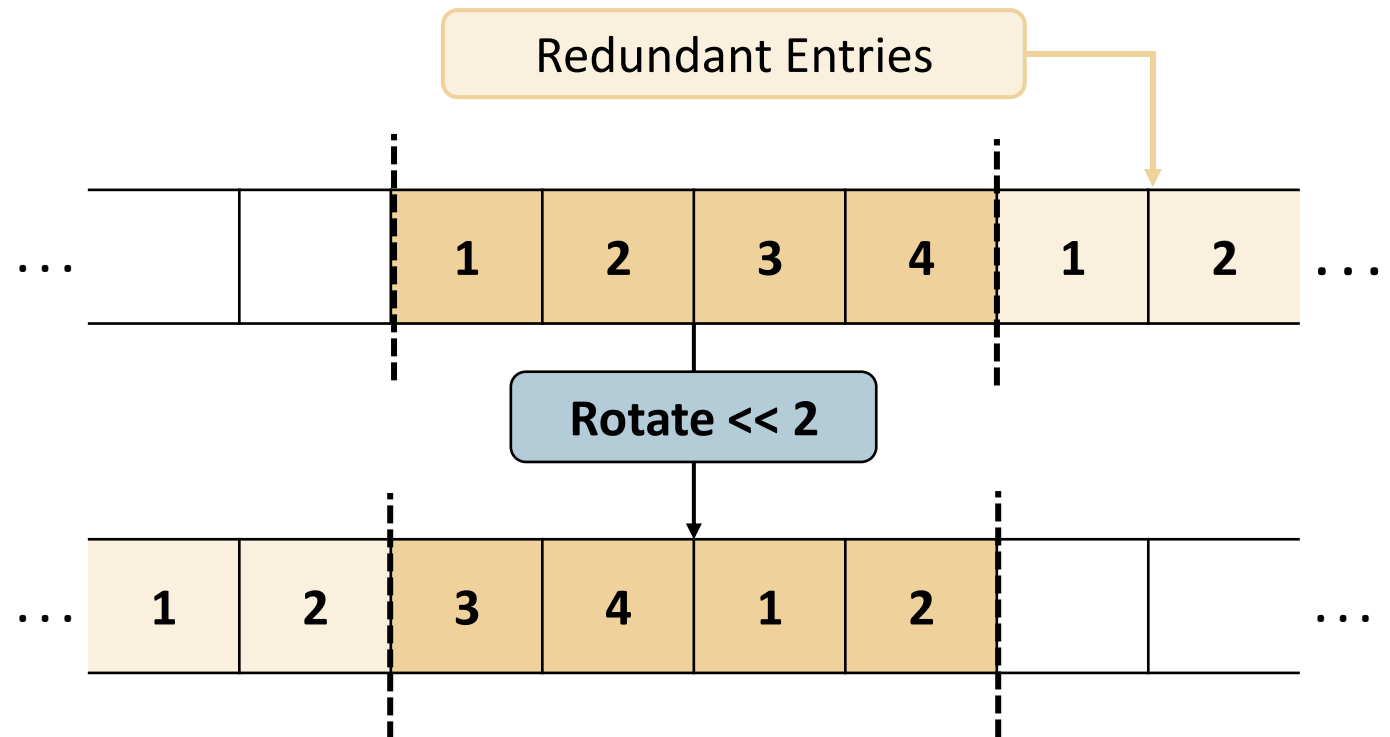Actually: Values wrap around the entire ciphertext vector

# Standard Permutations

- Expensive

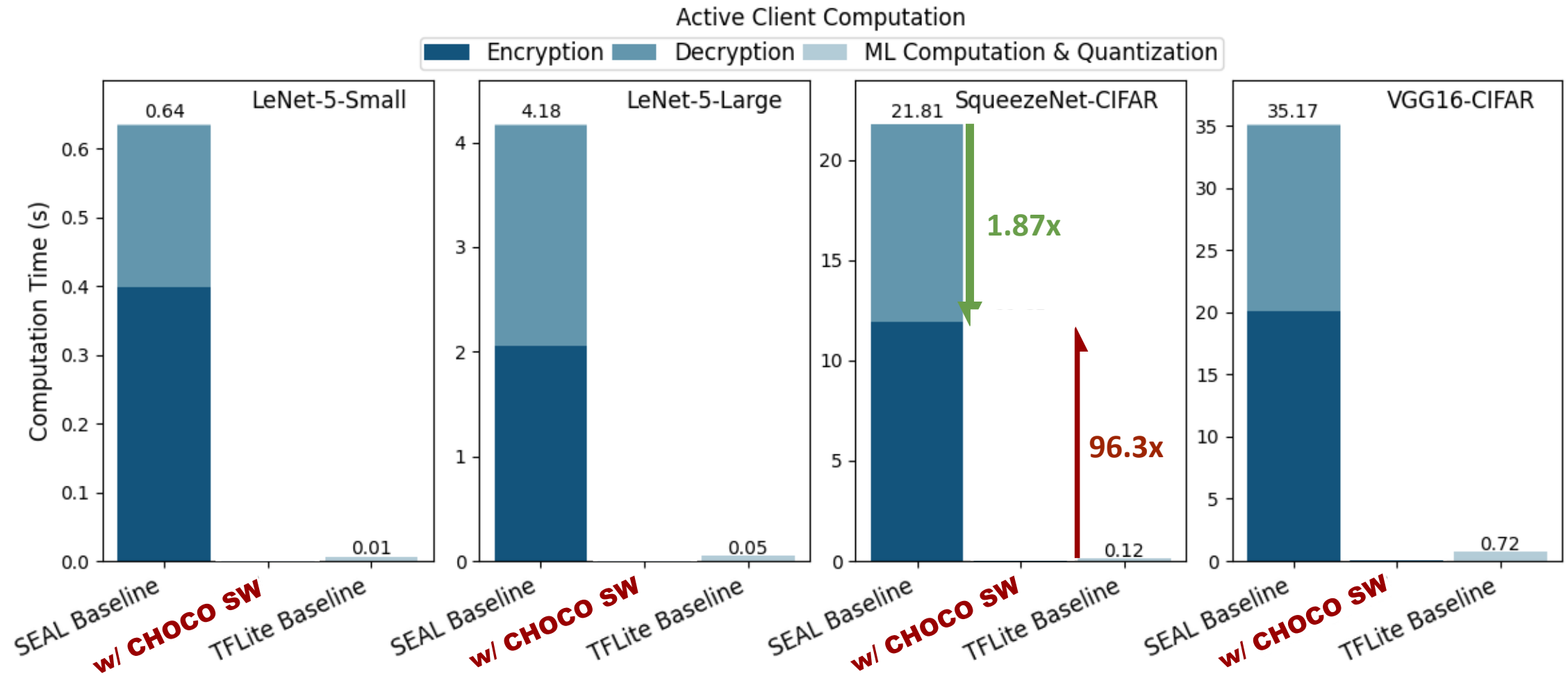- Computation & Noise Growth

- High Arithmetic Depth

# Rotational Redundancy

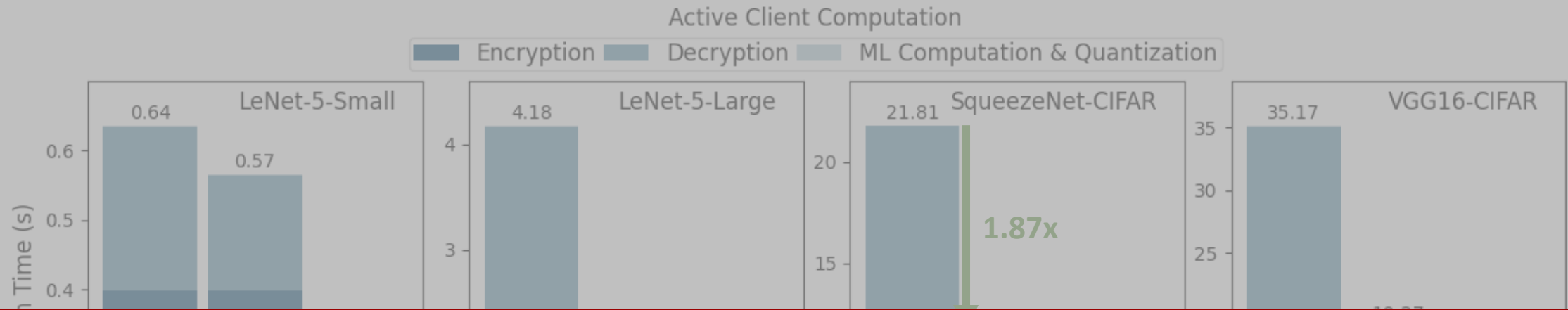- Novel Input Packing

- Single HE Rotation

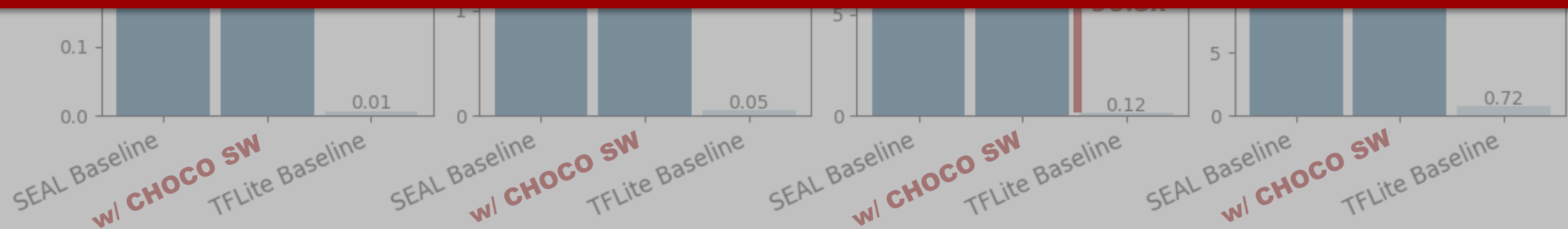- Low arithmetic depth

# CHOCO Algorithms Reduce Client Computation



Active Client Computation

- **50%** Smaller Ciphertexts
- **CHOCO SW** = SEAL baseline + Rotational Redundancy
- Average **1.7x** improvement over SEAL
- Average **62.5x** remaining overhead vs TFLite

# CHOCO Algorithms Reduce Client Computation



Active Client Computation

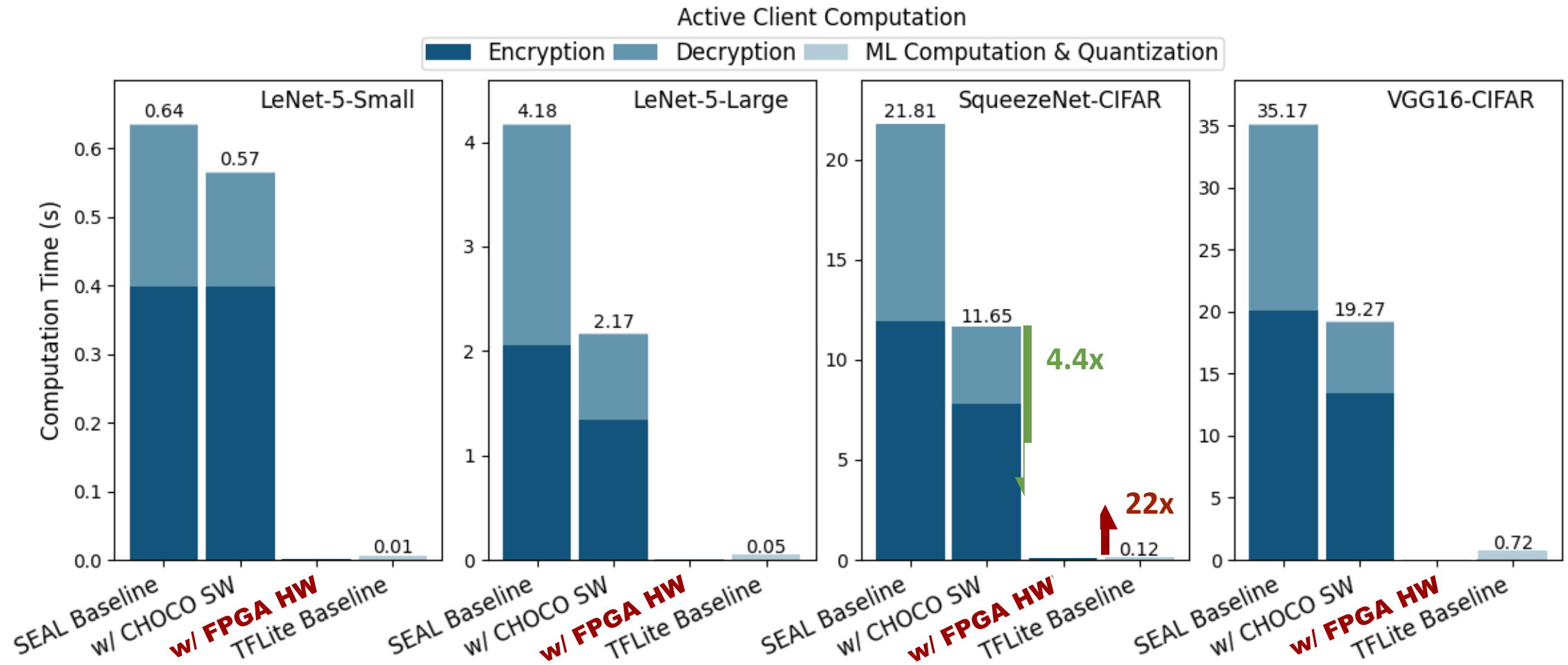Encryption ■ Decryption ■ ML Computation & Quantization

**CHOCO Algorithm Optimizations Provide a Critical but Insufficient Reduction in Client Computation**

- **CHOCO SW** = SEAL baseline + Rotational Redundancy
- **50%** Smaller Ciphertexts

- Average **1.7x** improvement over SEAL
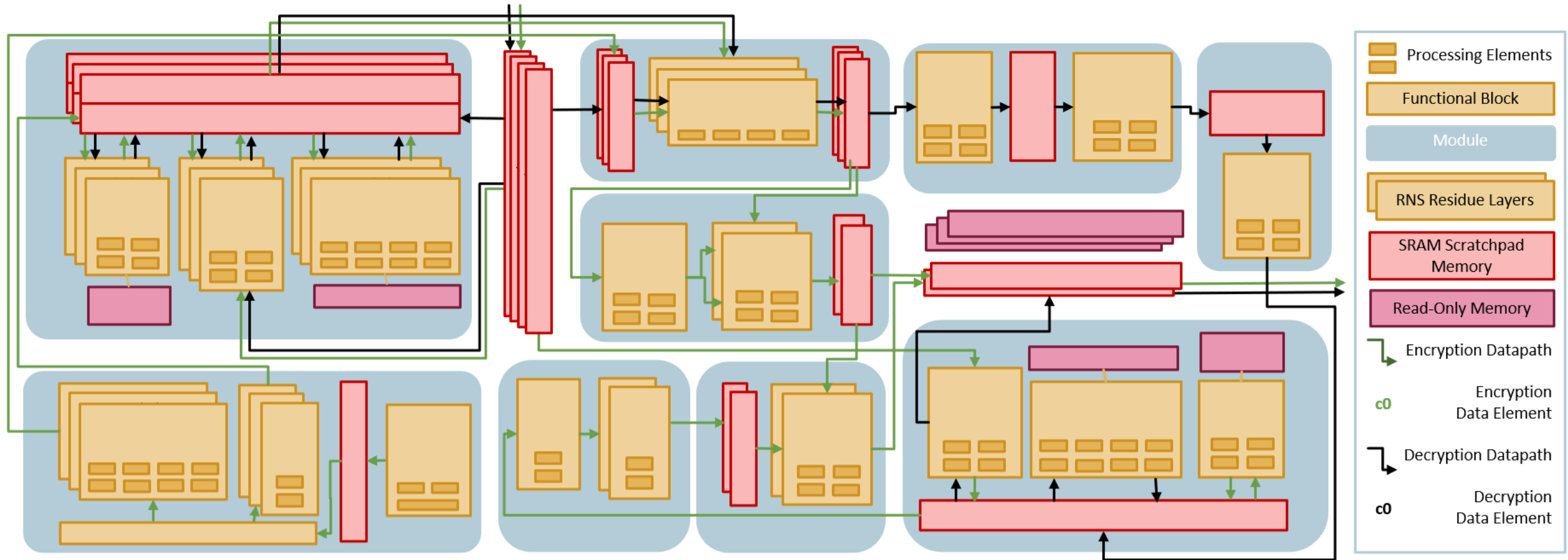- Average **62.5x** remaining overhead vs TFLite

# Hardware Acceleration

Electrical & Computer ENGINEERING

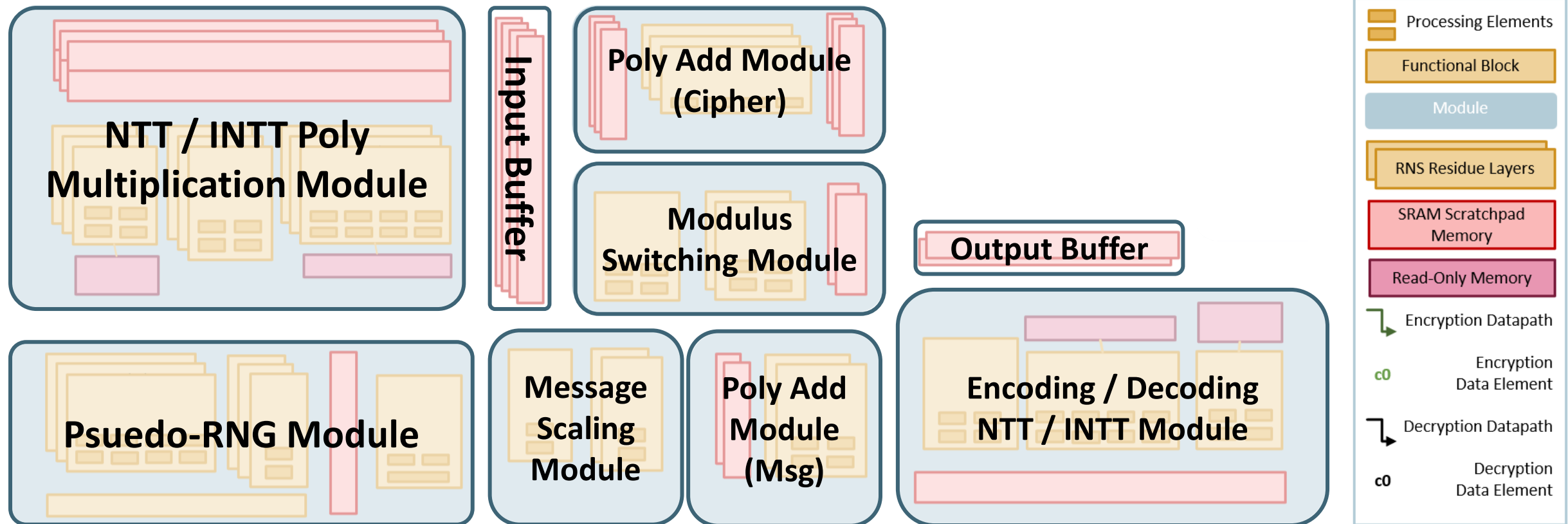# Existing FPGA Acceleration is Incomplete



Active Client Computation

- **FPGA HW** = CHOCO SW + Encryption/Decryption FPGA
- Average **14.5x** remaining overhead vs TFLite

A. Mert, E. Ozturk, and E. Savas. 2020. *Design and Implementation of Encryption/Decryption Architectures for BFV Homomorphic Encryption Scheme.* IEEE Transactions on Very Large Scale Integration (VLSI) Systems 28, 02 (feb 2020), 353–362.

Electrical & Computer ENGINEERING

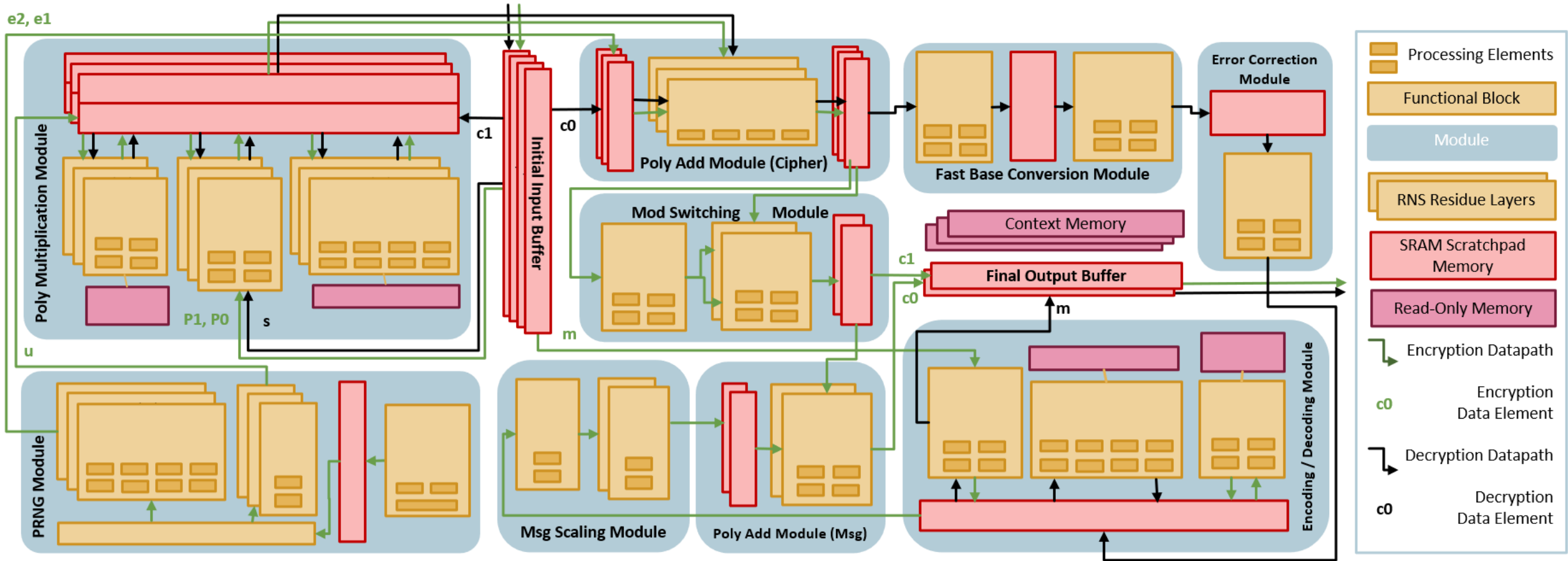# CHOCO – Through Accelerated Cryptographic Operations
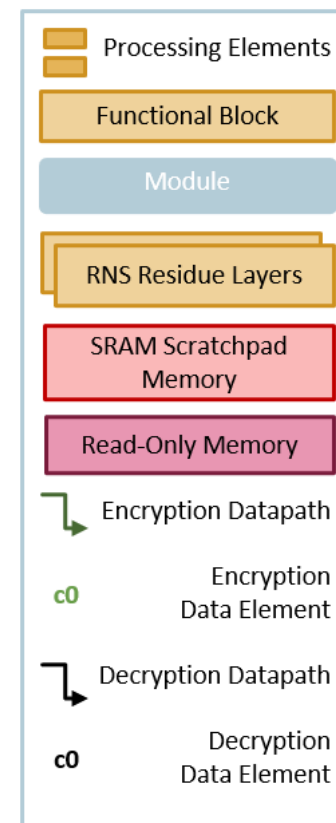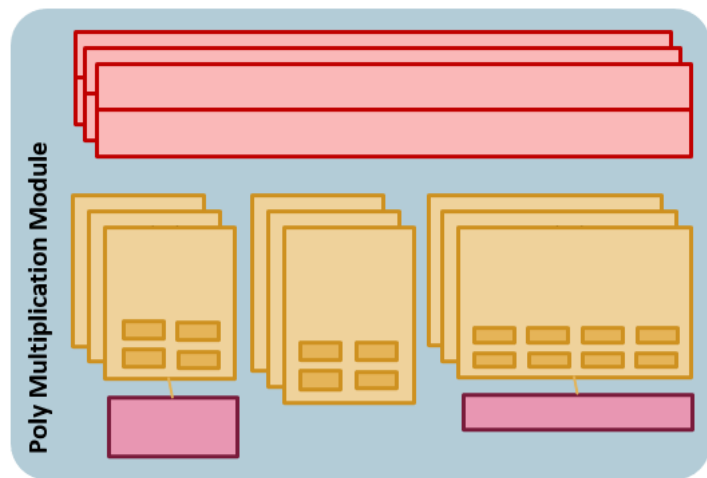
# CHOCO-TACO Encryption & Decryption Hardware



- $\textbf{Encrypt}(\text{pk}, m)$: For $m \in R_t$, let $\text{pk} = (p_0, p_1)$. Sample $u \xleftarrow{\$} R_2$, and $e_1, e_2 \leftarrow \chi$. Compute

$$\text{ct} = \left([\Delta m + p_0 u + e_1]_q, \; [p_1 u + e_2]_q\right).$$

# CHOCO-TACO Encryption & Decryption Hardware



- Encrypt(pk, $m$): For $m \in R_t$, let pk $= (p_0, p_1)$. Sample $u \xleftarrow{\$} R_2$, and $e_1, e_2 \leftarrow \chi$. Compute

$$\text{ct} = \left([\Delta m + p_0 u + e_1]_q, [p_1 u + e_2]_q\right).$$
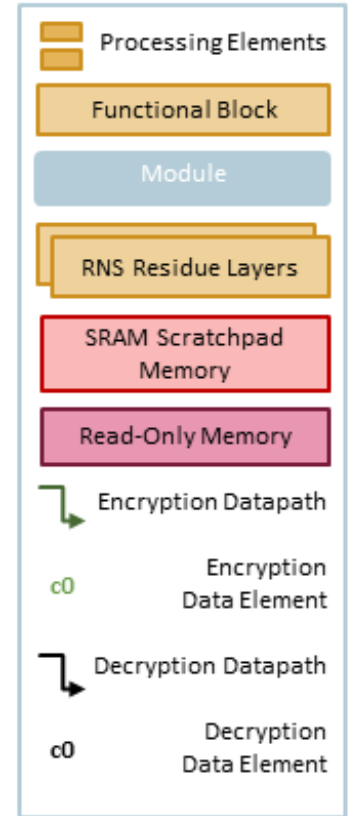
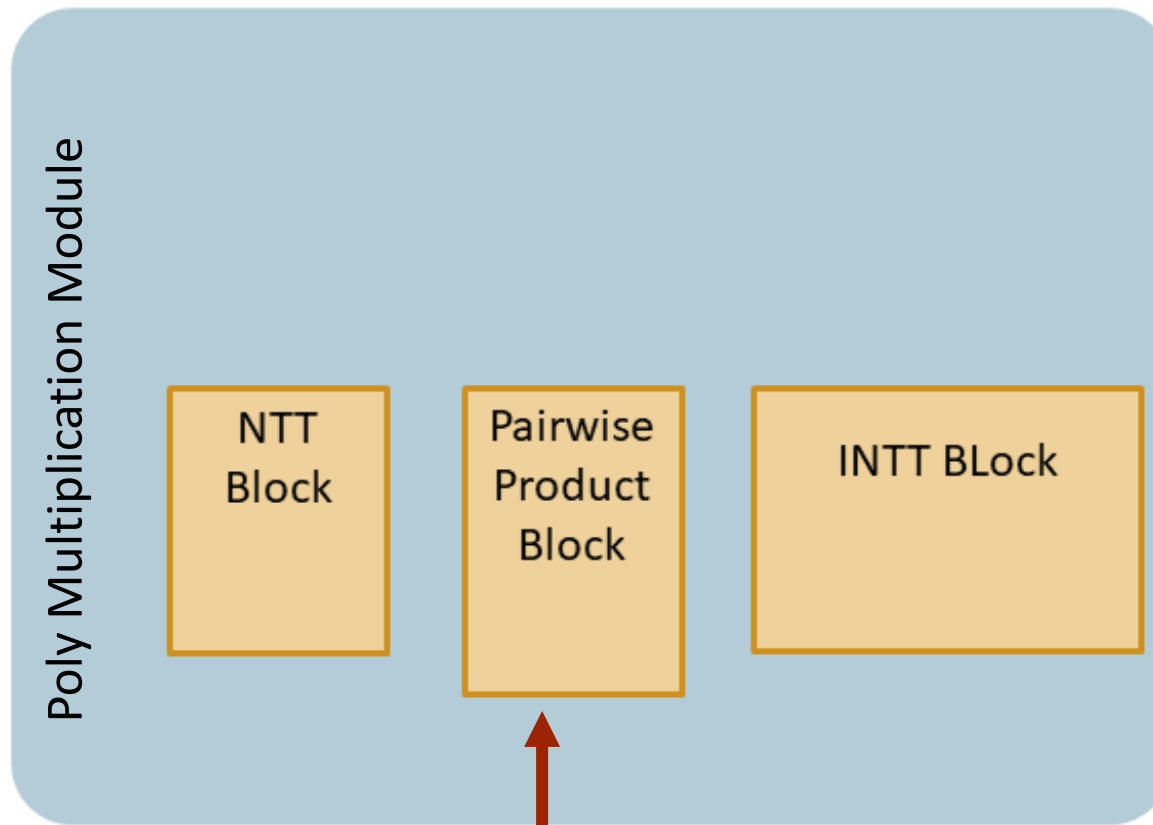# CHOCO-TACO Encryption & Decryption Hardware

# CHOCO-TACO Hardware Optimization
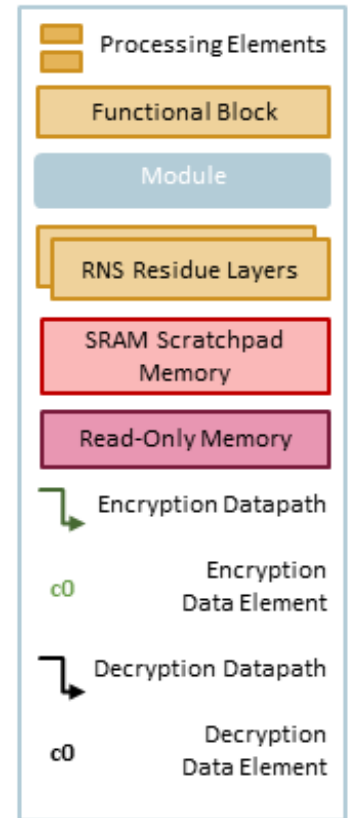
# CHOCO-TACO Hardware Optimization



**Poly Multiplication Module**

**Modules & Functional Blocks**
Specialized & Comprehensive

Legend:
- Processing Elements
- Functional Block
- Module
- RNS Residue Layers
- SRAM Scratchpad Memory
- Read-Only Memory
- Encryption Datapath
- c0 — Encryption Data Element
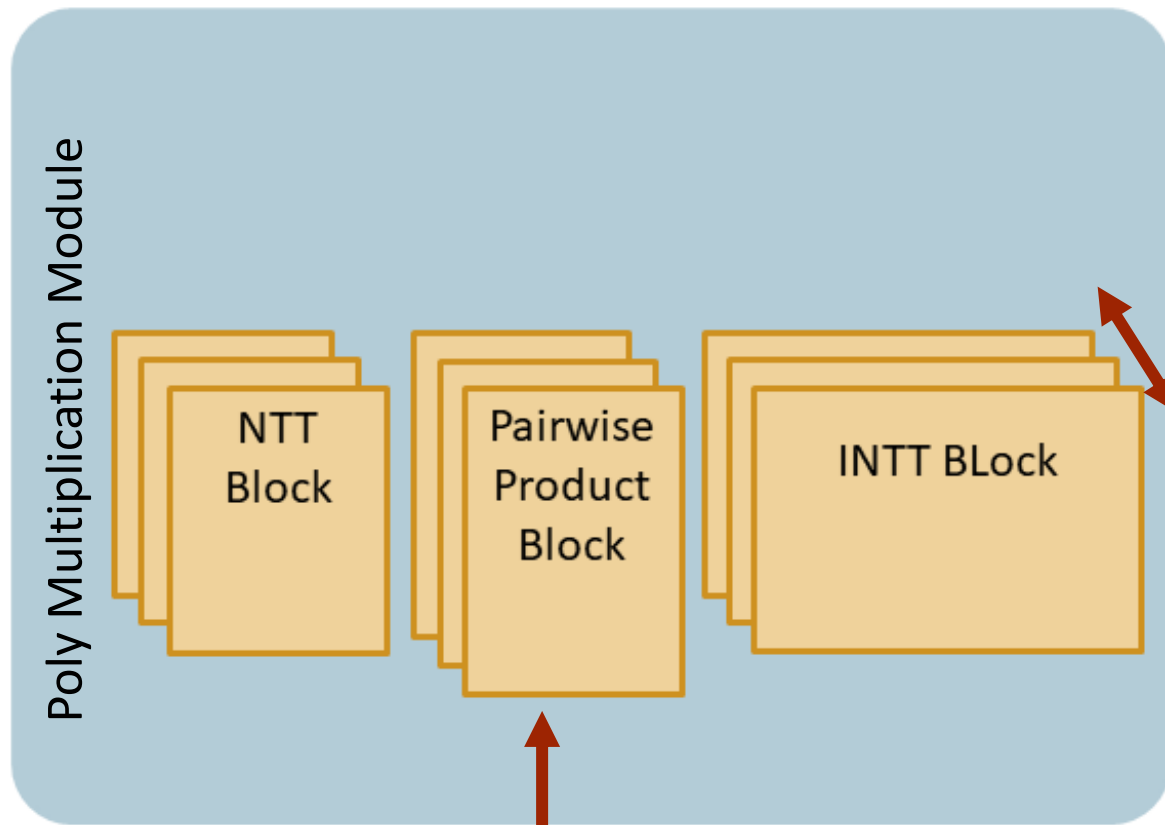- Decryption Datapath
- c0 — Decryption Data Element

# CHOCO-TACO Hardware Optimization



**Modules & Functional Blocks**
Pipelined & Comprehensive

# CHOCO-TACO Hardware Optimization



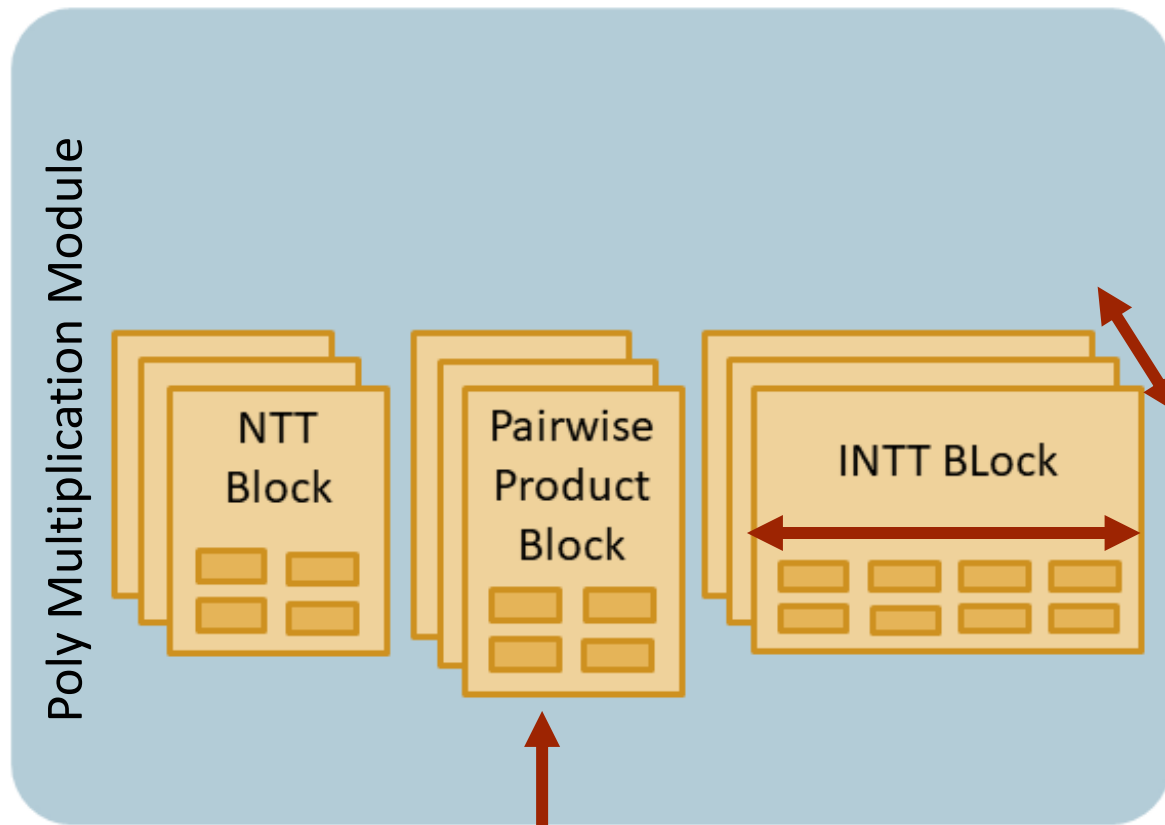**Conceptual Layers**
Polynomial Parallelism

**Modules & Functional Blocks**
Specialized & Comprehensive

# CHOCO-TACO Hardware Optimization



**Conceptual Layers**
Polynomial Parallelism

**Processing Elements**
Coefficient Parallelism

**Modules & Functional Blocks**
Specialized & Comprehensive

44

# CHOCO-TACO Hardware Optimization



**SRAM Buffers**
Parallelism, Pipelining,
Minimal Data Movement

**Conceptual Layers**
Polynomial Parallelism

**Processing Elements**
Coefficient Parallelism

**Modules & Functional Blocks**
Specialized & Comprehensive

Electrical & Computer ENGINEERING

# CHOCO-TACO Hardware Optimization



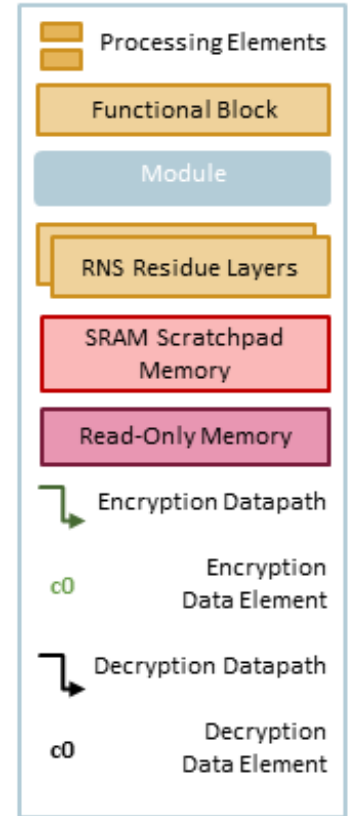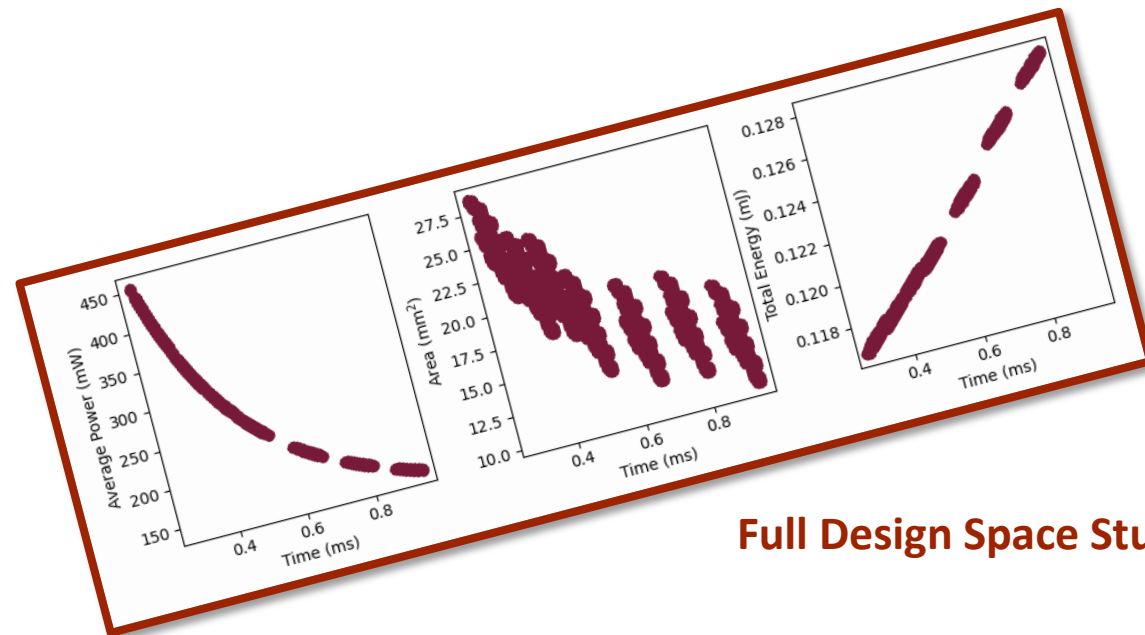**Full Design Space Study in Paper!**

Electrical & Computer ENGINEERING

46

# CHOCO-TACO Encryption & Decryption Hardware



19.3 mm$^2$ area. Consumes 200 mW power, .1228 mJ to perform a single encryption in .66 ms.

# Outline

- Introduction

- Encrypted Computing Background

- **CHOCO System Implementation**
  - **Algorithm Optimizations & Hardware Acceleration**

- Results & Conclusions

Electrical & Computer
ENGINEERING

# Outline

- Introduction

- Encrypted Computing Background

- CHOCO System Implementation

- **Results & Conclusions**
  - **Active Client Computation, Communication & Applications**

# CHOCO-TACO Accelerates Client Compute



- **CHOCO HW** = CHOCO SW + CHOCO-TACO Encryption/Decryption Simulated ASIC
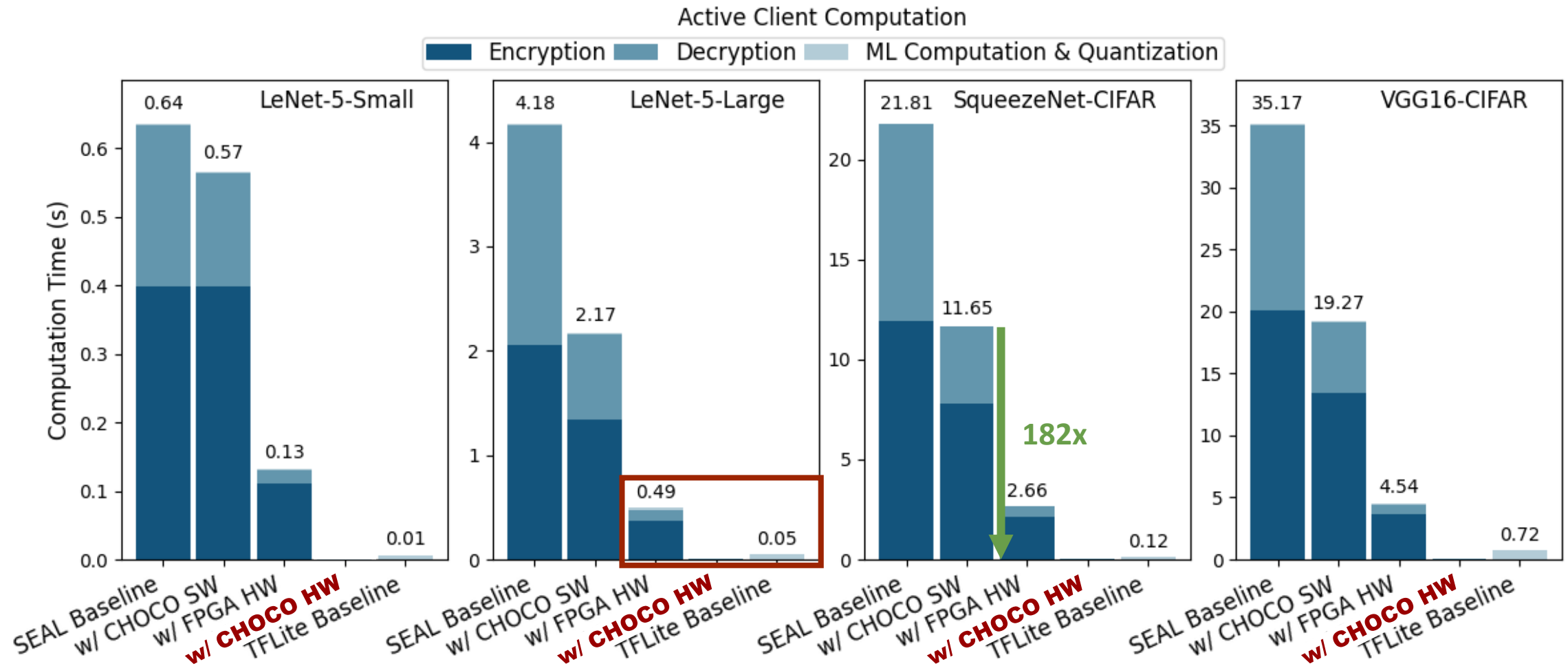
- Average **123.3x** Improvement over CHOCO software alone

# CHOCO-TACO Accelerates Client Compute



Active Client Computation

- Average **29x** better than FPGA accelerators

- Average **2.2x** better than local compute via

# CHOCO Makes Client End-to-End Costs Feasible



Local Compute w/ TFLite vs Offloaded Compute w/ CHOCO

Electrical & Computer
ENGINEERING

# CHOCO Makes End-to-End Client Costs Feasible
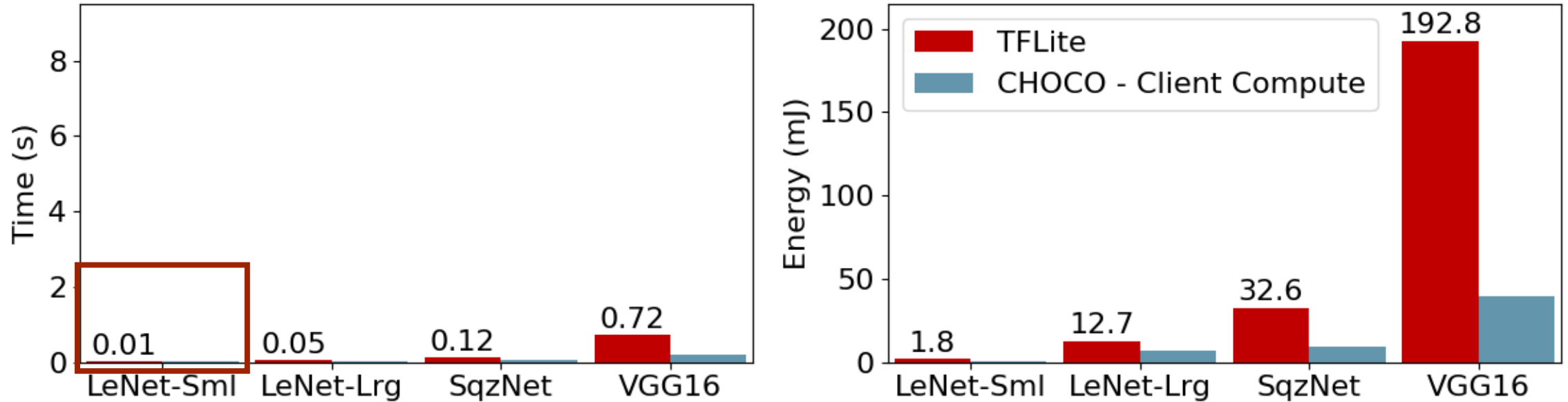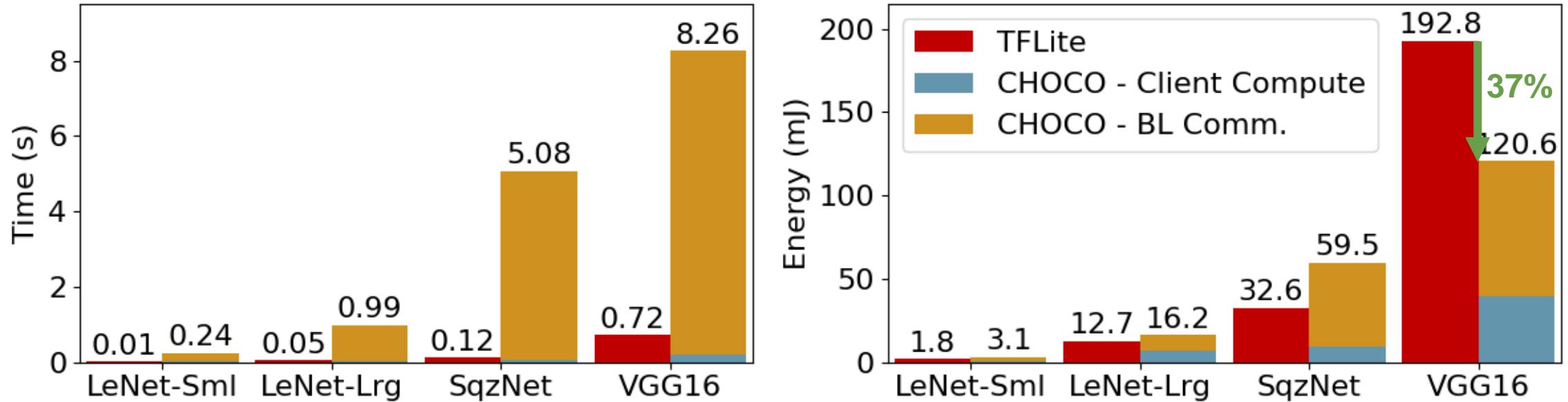


Local Compute w/ TFLite vs Offloaded Compute w/ CHOCO

Legend: TFLite, CHOCO - Client Compute, CHOCO - BL Comm.

Time (s): LeNet-Sml 0.01 / 0.24, LeNet-Lrg 0.05 / 0.99, SqzNet 0.12 / 5.08, VGG16 0.72 / 8.26

Energy (mJ): LeNet-Sml 1.8 / 3.1, LeNet-Lrg 12.7 / 16.2, SqzNet 32.6 / 59.5, VGG16 192.8 / 120.6 (37%)

- Privacy-Preserving Offload can be **Competitive with Local Compute**

- **37%** decrease in energy consumption for VGG16

- Up to **66% communication reduction** from SEAL baseline

# CHOCO Algorithms Reduce Communication



Total Commnication Costs

- Up to three orders of magnitude improvement in communication

- Nearly **90x** improvement over Gazelle [Juvekar `18]
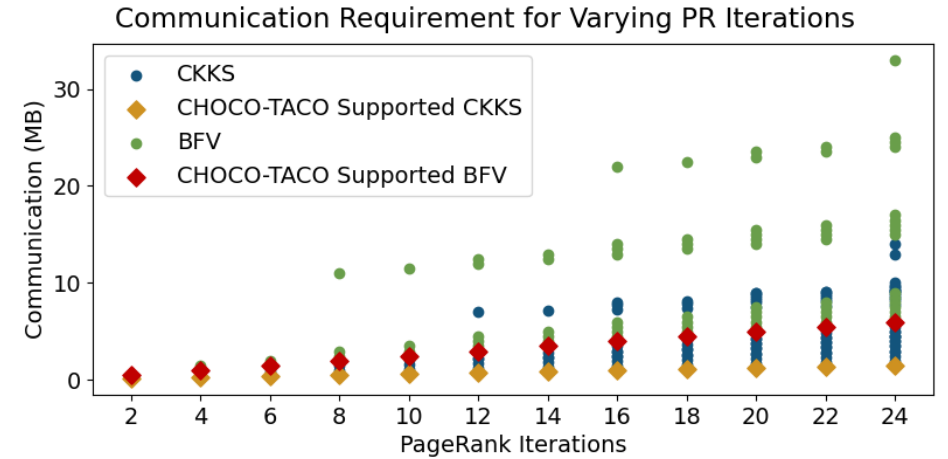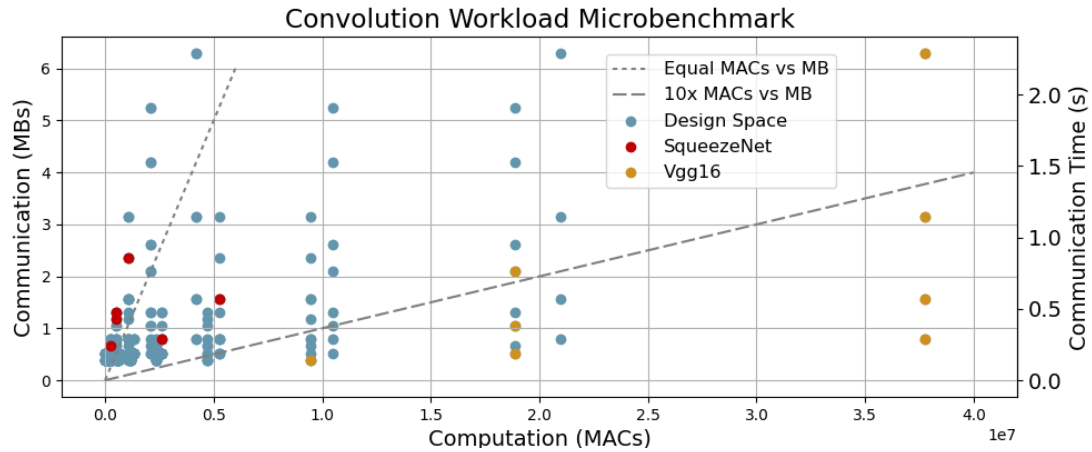
- **28x** better than LoLa (*not* client-aided) [Brutzkus `19]

- C. Juvekar, V. Vaikuntanathan, and A. Chandrakasan. 2018. *GAZELLE: A Low Latency Framework for Secure Neural Network Inference.* In Proceedings of the 27th USENIX Conference on Security Symposium (Baltimore, MD, USA) (SEC'18). USENIX Association, USA, 1651–1668.
- Alon Brutzkus, Ran Gilad-Bachrach, and Oren Elisha. 2019. *Low Latency Privacy Preserving Inference*. In Proceedings of the 36th International Conference on Machine Learning (Proceedings of Machine Learning Research, Vol. 97), Kamalika Chaudhuri and Ruslan Salakhutdinov (Eds.). PMLR, Long Beach, California, USA, 812–821.

Electrical & Computer ENGINEERING

# See Paper for More Applications & Results



Convolution Workload Microbenchmark

**Unmodified DNN Networks**



Communication Requirement for Varying PR Iterations

**PageRank**



Performance of Encrypted Distance Calculation Variations

**Encrypted Distance Calculations (K-Means & KNN)**

Electrical & Computer ENGINEERING

# Conclusions

- CHOCO motivates and prioritizes **client-aware optimizations**

- CHOCO alogorithm optimizations **reduce communication by orders of magnitude** over prior work

- CHOCO-TACO hardware **comprehensively accelerates client-side cryptographic primitives**

- CHOCO **enables participation from resource-constrained devices** in client-aided encrypted computation

- CHOCO makes client responsibility **competitive with local compute**

- CHOCO benefits **generalize to diverse applications**

Electrical & Computer
ENGINEERING

# Client-Optimized Algorithms & Acceleration for Encrypted Compute Offloading

## Thank You! Questions?

McKenzie van der Hagen – mckenziv@andrew.cmu.edu

Brandon Lucia – blucia@andrew.cmu.edu

Electrical & Computer
ENGINEERING