# BOSTON UNIVERSITY

# Affordable and Practical Acceleration of CKKS-based Fully Homomorphic Encryption

Rashmi Agrawal, Leo de Castro, Rabia Yazicigil,

Anantha Chandrakasan, Vinod Vaikuntanathan, Chiraag Juvekar, Ajay Joshi

rashmi23@bu.edu

Fully Homomorphic Encryption (FHE)



### **FHE Applications**

- Medical, Financial, Supply chain, Marketing, etc.
- Machine learning: Logistic regression training



Problem: Bootstrapping is the major bottleneck.

#### SimFHE: Custom Simulator

- SimFHE:
- A python-based architectural simulator for modeling CKKS operations at subroutine level.
- Set scheme parameters.
- Set architecture parameters.
- Keeps track of compute operations as well as DRAM transfers for a given cache size.
- Study compute & memory trade-off.
- Explore various optimizations and select parameters to optimize throughput.

#### Analysis using SimFHE

- Arithmetic intensity analysis for CKKS operations and an end-to-end application.





- Memory access pattern optimization:
- Regorganize low-level memory access pattern for faster data access.

Limb index Slot index 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 11 10 9 8 7 6 5 4 3 2 1 0 (a) Baseline address mapping

Lower-order bits Higher-order bits Higher-order bits Lower-order bits for limb index for limb index for slot index for slot index 31 30 29 28 27 26 25 24 23 22 21 20 19 18 17 16 15 14 13 12 11 10 9 8 7 6 5 4 3 2 1 0 (b) Optimized address mapping

Byte index Column Bank Bank Rank Row within a row Index Index Group Index

|          | Limb-wise<br>access | Slot-wise<br>access | Full<br>Switch |
|----------|---------------------|---------------------|----------------|
| Baseline | 2.3 ms              | 9.2 ms              | 11.5 ms        |

#### • Hierarchical caching optimizations:

2.5 ms

Optimized

 Compute as much as possible to maximize data reuse.

2.2 ms

4.7 ms

 Re-order operations to maximize cache utilization.



Algorithmic optimizations:

500

• Double-hoisting to reduce orientation switch.

 Improve arithmetic intensity of low level operations.



### **Key Observations**

8.2x improvement in bootstrapping throughput.

#### **Bootstrapping Throughput Comparison**

#### **FAB: FPGA-based Accelerator**

- FAB an FPGA-based accelerator for bootstrappable FHE workloads.
- First ever bootstrapping implementation on FPGA.
  - Secure and practical parameter set
- Balanced FPGA design.
  - Not memory-bound.
- Only 256 functional units operating at 300MHz
- High data rates from/to main memory at 450MHz.
- Effective utilization of limited on-chip memory. Modified datapath for basic FHE operations like KeySwitch.
  - Smart Operation scheduling.

#### **FAB: Overall Architecture**

- Four components:
- Host CPU offloads RTL design to the FPGA.
- RTL design packaged as kernel code.
- 8GB HBM2 memory stacks.
- 100G Ethernet CMAC subsystem.



#### **FAB: Evaluation**

- Designed in Verilog and implemented on Xilinx
- Alveo U280 deployed in cloud environment.
- Basic FHE operations' performance:

|         | FAB     | GPU     | Speedup<br>(Time) |
|---------|---------|---------|-------------------|
| Add     | 0.04 ms | 0.16 ms | 3.85x             |
| Mult    | 1.71 ms | 2.96 ms | 1.73x             |
| Rescale | 0.19 ms | 0.49 ms | 2.62x             |
| Rotate  | 1.57 ms | 2.55 ms | 1.62x             |

Bootstrapping performance:

|     | Time<br>(micros) | Speedup<br>(Time) | Speedup<br>(Cycles) |
|-----|------------------|-------------------|---------------------|
| CPU | 101.78           | 213x              | 2485x               |
| GPU | 0.740            | 1.55x             | 6.35x               |
|     |                  |                   |                     |



- Bootstrapping has low arithmetic intensity.
- Limited memory bandwidth is the bottleneck.



• 3.2x improvement in arithmetic intensity.



Memory bandwidth is still a bottleneck.

• CPU/GPU solutions are limited by main memory bandwidth.

Existing ASIC proposals are too expensive:

• Need large on-chip memory and register file. Need 12nm/7nm technology nodes.

| FAB 0.4// |
|-----------|
|-----------|

Logistic regression model training with 8 FPGAs:

|       | Time<br>(sec) | Speedup<br>(Time) | Speedup<br>(Cycles) |
|-------|---------------|-------------------|---------------------|
| CPU   | 37.05         | 456x              | 5318x               |
| GPU   | 0.775         | 9.5x              | 39x                 |
| FAB-1 | 0.103         | 1.3x              | 1.3x                |
| FAB-2 | 0.081         | -                 | -                   |

## Conclusion

 Affordable and practical FHE acceleration solution.

• MAD techniques provide Optimizations feasible with small cache sizes, agnostic of platforms.



• FAB provides Practical performance using FPGA at a fraction of ASIC cost.

