Secure Digital In-Memory Compute for Machine Learning Applications

Maitreyi Ashok

PI: Professor Anantha Chandrakasan

Work in collaboration with Saurav Maji (formerly MIT), Xin Zhang (IBM), John Cohn (IBM)





In Memory Compute



mismatch limit precision

Does not limit precision and benefits from technology scaling

Need for Security



Proposed System



Boolean Sharing Compute

Split each data bit and computation into separately computed *shares Power consumption decorrelated from data*

 $Original Function \\ t = F(b,c) \longrightarrow (t^1, t^2, t^3) = F^{sh}((b^1, b^2, b^3), (c^1, c^2, c^3))$

• Properties (for Practical Levels of Security)

- Correctness:
 - If $b = b^1 \oplus b^2 \oplus b^3$, $c = c^1 \oplus c^2 \oplus c^3$, Then $t = t^1 \oplus t^2 \oplus t^3$
- Non-Completeness:
 - If $F^{sh} = \{F^1, F^2, F^3\}$, Then each of F^j does not include all shares of each input
- Approximate Uniformity:
 - For each sub-circuit, each shared output has the same distribution bias as the unshared output
 - Outputs that are not jointly uniform are not combined directly

Multiplier Optimization

Conventional: Shared AND gate



- X Random bit refreshing to maintain uniformity
- Registers to maintain non-completeness
- ✗ 1 multiply = 48 gate-equivalents
- Standard bit-serial multiply for digital IMC

Proposed: Shared XNOR gate



- \checkmark
 - Shares maintain uniformity without random bits
 - Can cascade to next gate without registers
- 1 multiply = 6 gate-equivalents
- Need data format conversion at macro interface Negligible effect on NN accuracy

Addition, Accumulation use similar methods to reduce latency/eliminate random bits

School of Engineering



Model Decryption for BPA Security



Secure Key Generated On-Chip



Secure Write Reset
 Write fixed value to remove data dependence
 2 PUF Evaluation
 Cut and reconnect feedback transistor
 Settles to 0 or 1 since + or – side stronger due to local mismatch
 3 Standard Read
 Use differential temporal majority voting for security

Fabricated IC + Evaluation Setup





IMC Performance



Throughput (GOPS) 0.55 V, 80 MHz	Unprotected	41.0 (4b weight, 1b act) 9.10 (4b weight, 8b act)
	Protected	81.9 (4b weight, 1b act) 10.2 (4b weight, 8b act)
Energy Efficiency (TOPS/W) 0.55 V, 80 MHz	Unprotected	90.2 (4b weight, 1b act) 14.4 (4b weight, 8b act)
	Protected	6.94 (4b weight, 1b act) 0.89 (4b weight, 8b act)
Area Efficiency (TOPS/mm2) 0.55 V, 80 MHz	Unprotected	3.01 (4b weight, 1b act) 0.67 (4b weight, 8b act)
	Protected	0.49 (4b weight, 1b act) 0.061 (4b weight, 8b act)

Security Evaluation



PUF Security Analysis



Temporal Majority Voting SCA Security (CNN attack)

Key	Not Shared	Shared	Shared & Diff.
RMSE from FSM	250K train	1M train	
Unshared HW	1.38	4.90	7.20
Shared HW	N/A	2.35	32.13

Conclusion + Future Work

- Generalized IMC solution for ML with *Privacy & Integrity*
 - Side Channel and Bus Probing Attack Security for In Memory Compute
 - No random bits from PRNGs required
 - No limitations on neural network accuracy
- Future Improvements
 - More exploration of tradeoffs between security and area/energy overheads
 - Usage of approximate compute for further exploitation of natively secure compute gates
- Acknowledgements
 - MIT-IBM Watson AI lab for funding, helpful discussion, and support
 - NSF GRFP (Grant No. 1745302) and MathWorks Engineering Fellowship

