

FaultyGarble:

Fault Attack on Secure Multiparty Neural Network Inference

Mohammad Hashemi, Dev Mehta, Kyle Mitard, Shahin Tajik, Fatemeh Ganji

Worcester Polytechnic Institute

Let's do some computations





Multi-party computation

3

- MPC flourishing after the introduction of garbled circuits (GCs) by Yao [1].
 - Efficient computation, moderate communication complexity
- Real-world applications: banking, law, defense, medical applications, etc. [2].
 - Applications in cryptography: secure function evaluation, functional encryption, key-dependent message security, and recently, quantum circuits [2].



Attacks against Edge NNs



GC-based NN inference at the edge





Adversary model

- NN inference engines on an FPGA with a general-purpose processor
 - Usual implementation as in, e.g., [1]-[5]
 - Layer-by-layer inference
 - Alternating linear (fully connected, convolutional, etc.) and non-linear ReLU layers
- A client-server setting with the malicious client attempting to extract the model's weights held by the server
 - The neural network configuration is known to both client and server
 - The processor is not the IP to be protected: the malicious adversary knows the processor layout or can profile it to target points of interest
 - Adversary is capable of mounting physical fault attacks





Much simpler than one thinks!

- A simple example •
 - At the last layer
 - No bias

7

- The encryption is not shown
- What if we turn the AND to XOR? •
 - $-0 \operatorname{XOR} w_3 = w_3$



LW1

0

 w_2

x=0

W3

 $v = w_3$

At an intermediate layer

- Recall the layer-by-layer evaluation of the • NN
- When targeting the *k*-th layer, forcing the ReLU functions in the layers *k*-1,...,1 to behave linearly or like a buffer.
- After decoding, the output of the NN model: •

8



W1

 W_2

0

W3

 $v = w_2 w_3$

x=0

Experimental setup

- Target: A Genesys 2 development kit
 - AMD/Xilinx Kintex 7 FPGA in a flip-chip package
 - Clock frequency: 200 MHz
 - User-programmable outputs used as flags to show a successful fault
 - MIPS I instruction set: a family of RISC instruction sets
- Laser setup: ALPhANOV S-LMS
 - Wavelength: 1064 nm
 - 20x and 50x magnification lenses



Some results

Malleating and cryptanalysis against garbled NN Not protected against active attacks			Pure cryptanalysis against NN		
Network Dimensions	# Parameters	[1]	#Queries	Ours	#Fault
784-128-1	100,480	100,480	$2^{21.5}$	100,480	200,832
784-32-1	25,120	25,120	$2^{19.2}$	25,120	50,208
10-10-10-1	210	210	2^{16}	210	610
10-20-20-1	620	620	$2^{17.1}$	620	1820

#faults: $O(p\ell^2)$

Total number of the NN model parameters

Number of the layers



10 [1] Lehmkuhl (USENIX'21)

[2] Carlini (Crypto'20)

Take home messages

- Conventional theoretical countermeasures against active attacks does not work
- Is the attack limited to MIPS instruction set?
 - fetch-decode-execute cycle as a core operational process
 - Potential risks for garbled NNs implemented using other instruction sets
- Countermeasures
 - At the protocol level: protecting the instructions
 - At the hardware level: increasing the Hamming distance between binary codes in the instruction set
 - Reducing the predictability of the process by introducing random delays
 - Impact: no precise time-base to determine the best time to inject the fault



Conclusion

- What are other possible attacks against garbled circuits?
 - Timing side-channel analysis [1]
 - Photon emission analysis [2]
- Implementation attacks have been overlooked!
 - Even advanced protocols with protection against active attacks
- Hardware-based countermeasures
- Is it possible to further reducing the number of faults?





Thank You!









