

Fault Injection and Embedded Neural Networks: Models Integrity and Confidentiality

Pierre-Alain Moëllic

CEA LETI / Mines Saint-Etienne





Context: Security of Machine Learning

SotA: a decade of experience

→ mainly API-based attacks → Confidentiality & Privacy / Integrity / Availability

ALGORITHM / ABSTRACTION

→ Model / Data

Background – Deep Neural Network

Neurons, kernels, layers and... instructions to skip



Deep Neural Network

Feedforward models





Deep Neural Network

Feedforward models



Convolutional layer

Algorithm 1 Convolution layer (*K* kernels) **Input:** Tensor X of size $H \times H \times C$, parameters tensor Θ of size $Z \times Z \times C \times K$, bias tensor of size K **Output:** Tensor Y of size $H \times H \times K$ 1: **for** k in [1, K] **do** for x in [1, H] do 2: for y in [1, H] do 3: $Y_{x,y,c} = B_k$ 4: for m in [1, Z] do 5: for n in [1, Z] do 6: for c in [1, C] do 7: $Y_{i,j,c} + = \theta_{m,n,k,c} \cdot X_{x+m,y+n,k}$ 8: return Y

Fully-Connected layer

$$a_j^l(x) = \sigma\Big(\sum_{i \in (l-1)} \theta_{i,j} a_i^{l-1} + b_j\Big)$$

Instruction Skip

Impact of a single instruction skip

Very first experimental study¹ from C. Gaine (ANR PICTURE)

- Where to start?
- Critical attack paths? Adversarial goal?



- 2 injection means: Laser / EM pulse
- Cortex-M4 platform
- Inference of a standard CNN (trained on FashionMNIST)
- Impact of a single instruction skip
- → convolutional layer
- bias addition
- activation function



Instruction Skip

MIŅES

Saint-Étienn



question...



Parameters-based Adversarial Attacks

Some results about the Bit-Flip Attack (BFA)

Weight-based Adversarial Attacks

cks

Target internal parameters stored in memory

- Main reference: Bit-Flip Attack BFA¹
 - Theoretical "demonstration" on SotA CNN
 - First practical demonstration: RowHammer² attack (CPU, DRAM)
 - Former works³ on evaluating BFA
- ✤ Random bit-flips → many safety analysis
- ✤ Highly recommended works from D. Stutz⁴: Safety ← → Security





Weight-based Adversarial Attacks

Target internal parameters stored in memory

- ✤ BFA = Adversarial bit-flips → Faults on the most sensitive parameters
- Very similar to white-box white-box evasion attacks (adv exp)









State of the Art¹

- Most works rely on simulation only
- Practical exp: RowHammer attacks (DRAM)
- (very) Few and partial works on LFI on MCU against embedded DNN²

Practical Evaluation on MCU³?

- Security evaluator point of view
- Small MLP (compressed MNIST) embedded in 32-bit MCU, Cortex M3
 - Laser Fault Injection (LFI)
 - Bit-Set Fault Model
 - Explained and demonstrated for NOR-Flash memory of Cortex-M MCU by Colombier + Menu ^{3,4}



Qian, et al. A Survey of bit-flip attacks on deep neural network and corresponding defense.Electronics 2023
 Hou, et al. Security Evaluation of DNN Resistance against Laser Fault Injection, IPFA 2020
 Dumont et al. Evaluation of Parameter-based Attacks against Embedded Neural Networks with Laser Injection. SafeComp 2023
 Colombier, et al. Laser-induced Single-bit Faults in Flash Memory..., IEEE HOST 2019.
 Menu, et al. Single-bit laser fault model in NOR flash memories..., IEEE FDTC 2020









@ neuron-level → weighted sum

C code, weightedsum in a FC layer (NNoM)





 $;q7_t inB = *pB++$;Weight n+1 initialization r3, [r7, #80] ;Loading the 2 ldr address of the weight n ;Next weight r2, r3, #1 3 adds address r2, [r7, #80] ; Input value str 4 loading into r2 reg 5 ldrsb.w r3, [r3] ;Weight value loading. LASER SHOT ;Store of the r3, [r7,#23] 6 strb weight in SRAM reg

Assembler code, of line 6





Target a MLP model

- Brute-Force approach (4960 bits)
- Significant accuracy drops for MSB locations

- Guided-LFI with adapted BFA (BSCA)
- How practical LFI fit with simulation (BSCA)?



Exploitation?

Integrity threat... really?

- \clubsuit LFI \rightarrow powerful tool for practical robustness evaluation
- Use average adversarial accuracy as a task-quality metric
- #faults vs accuracy_drop is more interesting

Targeted BFA

(one targeted input / one targeted misprediction)

Backdoor attacks¹

Upcoming threats at training time for distributed paradigm

Federated Learning (model poisoning)





Exploitation → confidentiality threats

Model Extraction

Cez

MINES int-Étie

A growing confidentiality concern

Adversarial Goal¹

- FIDELITY (model cloning)
 Target: architecture & parameters
- TASK-PERFORMANCE (steal performance, save design & training time)
 Transfer knowledge
- **BLACK** \rightarrow WHITE BOX ATTACK







Theoretical framework (fidelity scenario)

- ♦ VICTIM MODEL M_W trained with $(X^{train}, Y^{train}), X^{train} \sim D_\chi$
- ***** SUBSTITUTE MODEL: M'_{θ}
- Threat Model (parameter extraction): Architecture is known / (very) limited access to X^{train}

SotA

- ♦ Active learning principle: feed M_W to build a substitute training dataset for $M'_{\theta} \rightarrow (x, M_W(x))$
- Hard work... Need of a huge amount of query/output pairs
- **DEEPSTEAL**¹: use BFA+Rowhammer (RamBleed) to guess some parameter values and improve the training of M'_{θ}

Exploiting Safe Error Attack (SEA)







Exploiting Safe Error Attack (SEA)





 \checkmark Specific inputs lead to very successful SEA \rightarrow predictions with some uncertainty

✓ Train M'_{θ} with constrains from the partial extracted values → very efficient



CONCLUSION



- ✓ Very active context regarding safety / security concerns for AI systems
- ✓ Growing needs of
 - ✤ Threat models → risk analysis & exploitation
 - Practical demonstrations and evaluations (or platform-based simulation)
 Models, platforms...
 - Robust implementations
 - Defenses (practical) evaluation
- Vumerous attack vectors and paths \rightarrow works for everyone \odot



Thank you for your attention

Support & Funding

EU project InSecTT

InSecTT

PICTURE

NANOELEC.

French ANR, IRT Nanoelec

French ANR PICTURE program

This work benefited from the French **Jean Zay** supercomputer with the AI dynamic access program.