

#### **Compiled Models, Built-In Exploits: Uncovering Pervasive Bit-Flip Attack Surfaces in DNN Executables**

Yanzuo Chen<sup>†</sup>, Zhibo Liu<sup>†</sup>, Yuanyuan Yuan<sup>†</sup>, Sihang Hu<sup>‡</sup>, Tianxiang Li<sup>‡</sup>, Shuai Wang<sup>†</sup>

<sup>†</sup>The Hong Kong University of Science and Technology, <sup>‡</sup>Huawei Technologies



DEPARTMENT OF COMPUTER SCIENCE & ENGINEERING



#### **DNN Models**















Ninja in camouflage	<mark>95%</mark>
Spooky ghost	4%
Professional chef	1%





35

Ninja in camouflage	<mark>95%</mark>
Spooky ghost	4%
Professional chef	1%

- Bit flips via hardware fault injection
  - e.g., Rowhammer



Take car



Ninja in camouflage	<mark>95%</mark>
Spooky ghost	4%
Professional chef	1%

- Bit flips via hardware fault injection
  - e.g., Rowhammer
- Pervasive in DNN executables





35

Ninja in camouflage	<mark>95%</mark>
Spooky ghost	4%
Professional chef	1%



- e.g., Rowhammer
- Pervasive in DNN executables
- Gray-box, restricted attacker



**DNN** Executables



36

Ninja in camouflage	<mark>95%</mark>
Spooky ghost	4%
Professional chef	1%



- e.g., Rowhammer
- Pervasive in DNN executables
- Gray-box, restricted attacker
- 70% confidence;  $\sim$ 1.4 flips to succeed



## Preliminaries





- Compiled from DNN models
  - By "deep learning (DL) compilers"





- Compiled from DNN models
  - By "deep learning (DL) compilers"

- Wanted for their performance
  - Optimised at the computational graph level
  - Optimised for the target hardware platform





## Bit-Flip Attacks (BFAs)

• Flipping data bits in the memory (DRAM)

- Rowhammer: "Hardware fault injection" attack
  - Software-triggered hardware bug: Special access patterns
  - Current leakage between DRAM cells
  - DDR3 / DDR4 / ECC / DDR5 / ...

 $egin{array}{ccc} 0 
ightarrow 1 \ 1 
ightarrow 0 \end{array}$ 



#### BFAs on DNNs: Existing Work



## BFAs on DNNs: Existing Work

- Yes, researchers have done this
  - Flip bits in model weights
  - $\Rightarrow$  A type of white-box attack





## BFAs on DNNs: Existing Work

- Yes, researchers have done this
  - Flip bits in model weights
  - $\Rightarrow$  A type of white-box attack
- Few flips for full-precision models
- 12~24 flips for quantised models
  (Not anymore!)



Work	Avg. #Flips
BFA [69]	14.3
T-BFA (N-to-1) [71]	23.63
DeepHammer [88]	12.25
Ours	1.4



## Threat Model & Motivation



#### Attacker Objectives

- Model intelligence depletion
- Classification models (most frequently targeted)
  - Before attack: well-trained models, normal accuracy
  - After attack: random guessers (acc ightarrow 1/#classes)



#### Attack Flow













#### A More Restricted Attacker



#### A More Restricted Attacker

- Recall: White-box attackers in previous works
  - Model structures, weights ( $\rightarrow$  gradients), runtime setup...
  - Problem: Weights are often confidential



#### A More Restricted Attacker

- Recall: White-box attackers in previous works
  - Model structures, weights ( $\rightarrow$  gradients), runtime setup...
  - Problem: Weights are often confidential
- Our attacker: Weights  $\Rightarrow$  No gradient-based search

• How to identify which bits to flip?





- DNN executables are compiled code
  - We can flip bits in the code (.text section)
  - (Which are compiled DNN operators)





#### The Random Baseline



Knowledge: Offset of victim's .text  $\downarrow$  Bits to flip: Randomly choose in range



#### The Random Baseline (It was Bad)



#### The Random Baseline (It was Bad)

• Attack success rate (ASR): ~2%



#### The Random Baseline (It was Bad)

- Attack success rate (ASR): ~2%
- What happened in the remaining 98% of time?
  - Crash / no effect

segfault at 940c9 ip 00007f329a3df57b sp 00007f3299b54d10 error 6
segfault at 73249 ip 00007f329a3df57b sp 00007f3298b52d10 error 6
segfault at 20e09 ip 00007f329a3df57b sp 00007f329634dd10 error 6
segfault at 523c9 ip 00007f329a3df57b sp 00007f3297b50d10 error 6
segfault at ffffffffffff89 ip 00007f329a3df57b sp 00007f32290e9b9
segfault at 7f326a8ecc40 ip 00007f329a3df56f sp 00007f322a8ecb90 er
segfault at 48000028 ip 00007f329a3df577 sp 00007f32290e9b90 error
10909] trap invalid opcode ip:7f329a3df577 sp:7f32290e9b90 error:0 i
segfault at 7f329b34fdc0 ip 00007f329a3df56f sp 00007f329734fd10 er





#### Scan Survey Results

- Pervasive attack surface does exist
  - Different models / datasets / compilers

- How to identify vulnerable bits?
  - Do they change with training data? Yes
  - Do attackers have training data? No

Model	Dataset	Compiler	% Vuln.
ResNet50	CIFAR10	TVM	2.59
ResNet50	MNIST	TVM	3.14
ResNet50	Fashion	TVM	3.07
GoogLeNet	CIFAR10	TVM	2.56
GoogLeNet	MNIST	TVM	2.51
GoogLeNet	Fashion	TVM	2.59
DenseNet121	CIFAR10	TVM	2.66
DenseNet121	MNIST	TVM	2.10
DenseNet121	Fashion	TVM	2.29
<b>QResNet50</b>	CIFAR10	TVM	2.17
<pre>@GoogLeNet</pre>	CIFAR10	TVM	0.84
QDenseNet121	CIFAR10	TVM	0.52
LeNet	MNIST	TVM	2.52
DCGAN	MNIST	TVM	7.54
ResNet50	CIFAR10	Glow	2.48
ResNet50	MNIST	Glow	2.32



# Do vulnerable bits overlap?



#### Common Vulnerable Bits

- Trying to find recurring vulnerable bits
- Same model structure, trained on two datasets
  - ~45% vulnerable bits shared







## Transferable Vulnerable Bits

45% vulnerable bits transferable, despite different training sets



#### Building an Attack: In Seek of "Superbits"

• Using more local executables for profiling





## Building More Local Executables

- Train them on datasets of random noise!
  - Regulate weights
  - Unbiased choice





#### Attack Success Rate: 70%





#### Real World Examples on DDR4

Model	Dataset	#Flips	#Crashes	%Acc. Change
ResNet50	CIFAR10	1.4	0.0	$87.20 \rightarrow 10.00$
GoogLeNet	CIFAR10	1.4	0.0	$84.80 \rightarrow 10.00$
DenseNet121	CIFAR10	1.0	0.0	$80.00 \rightarrow 11.40$
DenseNet121	MNIST	1.2	0.0	$99.10 \rightarrow 11.20$
DenseNet121	Fashion	1.2	0.0	$92.50 \rightarrow 10.60$
<b>QResNet50</b>	CIFAR10	1.6	0.0	$86.90 \rightarrow 9.60$
<pre>@GoogLeNet</pre>	CIFAR10	1.4	0.0	$84.60 \rightarrow 11.20$
<b>QDenseNet121</b>	CIFAR10	1.6	0.0	$78.50 \rightarrow 10.20$
ResNet50	CIFAR10	1.4	0.0	78.80  ightarrow 10.00



#### Real World Examples on DDR4

Model	Dataset	#Flips	#Crashes	%Acc. Change
ResNet50	CIFAR10	1.4	0.0	$87.20 \rightarrow 10.00$
GoogLeNet	CIFAR10	1.4	0.0	$84.80 \rightarrow 10.00$
DenseNet121	CIFAR10	1.0	0.0	$80.00 \rightarrow 11.40$
DenseNet121	MNIST	1.2	0.0	$99.10 \rightarrow 11.20$
DenseNet121	Fashion	1.2	0.0	$92.50 \rightarrow 10.60$
<b>QResNet50</b>	CIFAR10	1.6	0.0	86.90  ightarrow 9.60
<b>QGoogLeNet</b>	CIFAR10	1.4	0.0	$84.60 \rightarrow 11.20$
<b>QDenseNet121</b>	CIFAR10	1.6	0.0	$78.50 \rightarrow 10.20$
ResNet50	CIFAR10	1.4	0.0	$78.80 \rightarrow 10.00$



#### What Does That Mean Exactly?

- Compare with DeepHammer, a SOTA method
  - Same attack objectives
  - Applicable to quantised models ("harder to attack")
  - ~12 flips were needed on average vs. ~1.4
  - White-box attacker was required vs. gray-box

Work	Avg. #Flips
BFA [69]	14.3
T-BFA (N-to-1) [71]	23.63
DeepHammer [88]	12.25
Ours	1.4



#### Conclusion

- DNN executables have large bit-flip attack surfaces
- $\bullet$  We achieve 70% confidence in vulnerable bits identification,  ${\sim}1.4$  flips to ruin model intelligence
- More security research on DNN executables please!



#### Thank You!

- PDF, source code, other materials
  - Visit <u>yanzuo.ch/debfa</u>

- Contact me
  - Yanzuo Chen: <a href="mailto:ychenjo@cse.ust.hk">ychenjo@cse.ust.hk</a>
  - Homepage: <u>yanzuo.ch</u>





Compiled Models, Built-In Exploits: Uncovering Pervasive Bit-Flip Attack Surfaces in DNN Executables. By Yanzuo Chen, Zhibo Liu, Yuanyuan Yuan, Sihang Hu, Tianxiang Li, and Shuai Wang.



#### Different Effects on GAN





#### Comparison with Existing Methods

#### TABLE XI

#### A COMPARISON OF ATTACK PERFORMANCE WITH PRIOR WORKS. FOR MITIGATIONS, ○, ●, AND ● DENOTE NO, PARTIAL, AND FULL MITIGATION, RESPECTIVELY.

Work	Attack	Avg.	Mitigable by				
VVUI K	Target	<b>#Flips</b>	<b>Q</b> [92]	A [85]	<b>D</b> [13]	W [49]	N [50]
BFA [69]	Weights	14.3	O				
T-BFA (N-to-1) [71]	Weights	23.63					
DeepHammer [88]	Weights	12.25					
Ours	Structure	1.4	0	0	0	0	0



#### Case Study Example: Control Flow Broken

Addr	Opcode bytes	x86 assembly instruction
0x70 0x73 0x76 0x78	83 F8 28 0F 4D C2 39 F0 0F 8D FA 00+ 00 00	<pre>cmp eax, 28h ;; max ID cmovge eax, edx ;; edx=28h cmp eax, esi ;; esi&lt;28h jge func_end</pre>
	(a) Assemb	oly code before BFA.
0x70 0x73 0x76 0x78	83 FC 28 0F 4D C2 39 F0 0F 8D FA 00+ 00 00	<pre>cmp esp, 28h ;; true cmovge eax, edx ;; true cmp eax, esi ;; true jge func_end ;; exit</pre>

(b) Assembly code after BFA.

