

Sponge Examples: Energy-Latency Attacks on Neural Networks

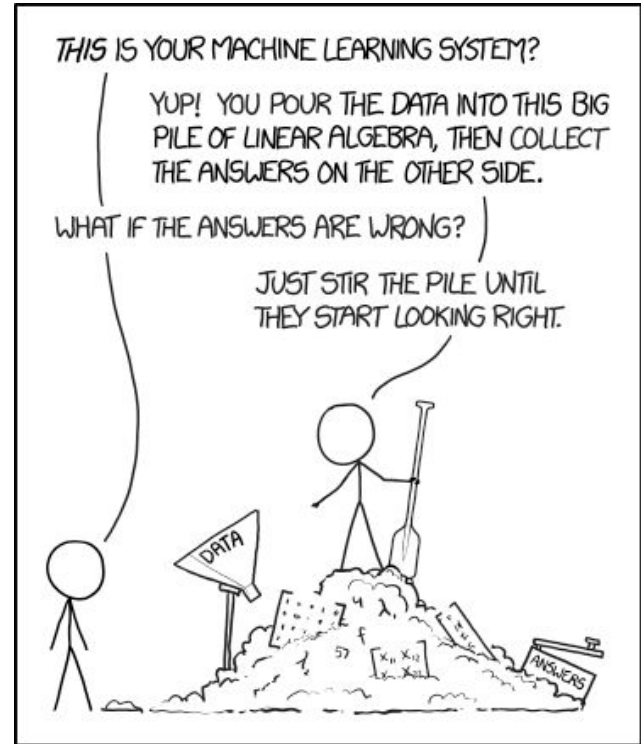
Ilia Shumailov*^, Yiren Zhao*, Daniel Bates*, Nicolas Papernot^, Robert Mullins*, Ross Anderson*

* University of Cambridge

^ University of Toronto, Vector Institute

Machine Learning

- Machine learning is everywhere
- We operate based on data, not formal rules
- There's a lot of non-determinism
- It is suddenly hard to define *Security*



Computer Security in context of Machine Learning

Class: bird
Confidence: 0.9659422039985657



+

Difference



=

Class: automobile
Confidence: 0.8248467445373535



- Adversarial examples exist for all models
- A large taxonomy of confidentiality and integrity attackers
- What about availability?

Availability

Ensuring **timely** and **reliable** access to and use of information.
(NIST Special Publication 800-12)

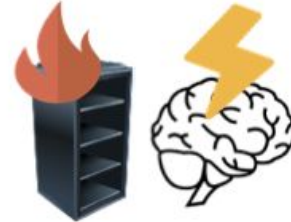
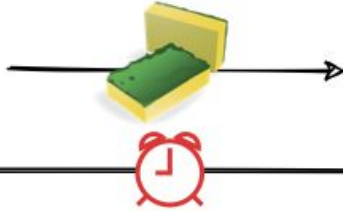
Availability



Benign Data



Sponge Examples



Increased latency

Over-heating and over-consumption of energy

Energy Gap

The amount of energy consumed by one inference pass (i.e. a forward pass in a neural network) depends primarily on:

- The overall **number of arithmetic operations** required to process the inputs;
- The **number of memory accesses** e.g. to the GPU DRAM.

Computation Dimensions

Modern networks have a **computational dimension**

- A large number of NLP models are **auto-regressive** e.g. RNNs and GPT2
- **Adaptive** input **dimensions** to help performance e.g. GPT2 uses Byte Pair Encoding
- ML components are **a part of loop**

see more in the paper ...

Computation Dimensions for GPT2

Auto-regressiveness adds an unbounded loop

Algorithm 1: Translation Transformer NLP pipeline

Result: y

```
1  $\downarrow O(l_{\text{tin}})$   
2  $x_{\text{tin}} = \text{Tokenize}(x)$ ;  
3  $y_{\text{touts}} = \emptyset$ ;  
4  $\downarrow O(l_{\text{ein}})$   
5  $x_{\text{ein}} = \text{Encode}(x_{\text{tin}})$ ;  
6  $\downarrow O(l_{\text{tin}} \times l_{\text{ein}} \times l_{\text{tout}} \times l_{\text{eout}})$   
7 while  $y_{\text{tout}}$  has no end of sentence token do  
8    $\downarrow O(l_{\text{eout}})$   
9    $y_{\text{eout}} = \text{Encode}(y_{\text{tout}})$ ;  
10   $\downarrow O(l_{\text{ein}} \times l_{\text{eout}})$   
11   $y_{\text{eout}} = \text{model.Inference}(x_{\text{ein}}, y_{\text{eout}}, y_{\text{touts}})$ ;  
12   $\downarrow O(l_{\text{eout}})$ ;  
13   $y_{\text{tout}} = \text{Decode}(y_{\text{eout}})$ ;  
14   $y_{\text{touts}}.\text{add}(y_{\text{tout}})$ ;  
15 end  
16  $\downarrow O(l_{\text{tout}})$ ;  
17  $y = \text{Detokenize}(y_{\text{touts}})$ 
```

Computation Dimensions for GPT2

Encoding adds **variable** I/O representation

Benign with 4 tokens for input of size 16:

Athazagoraphobia => ath, az, agor, aphobia

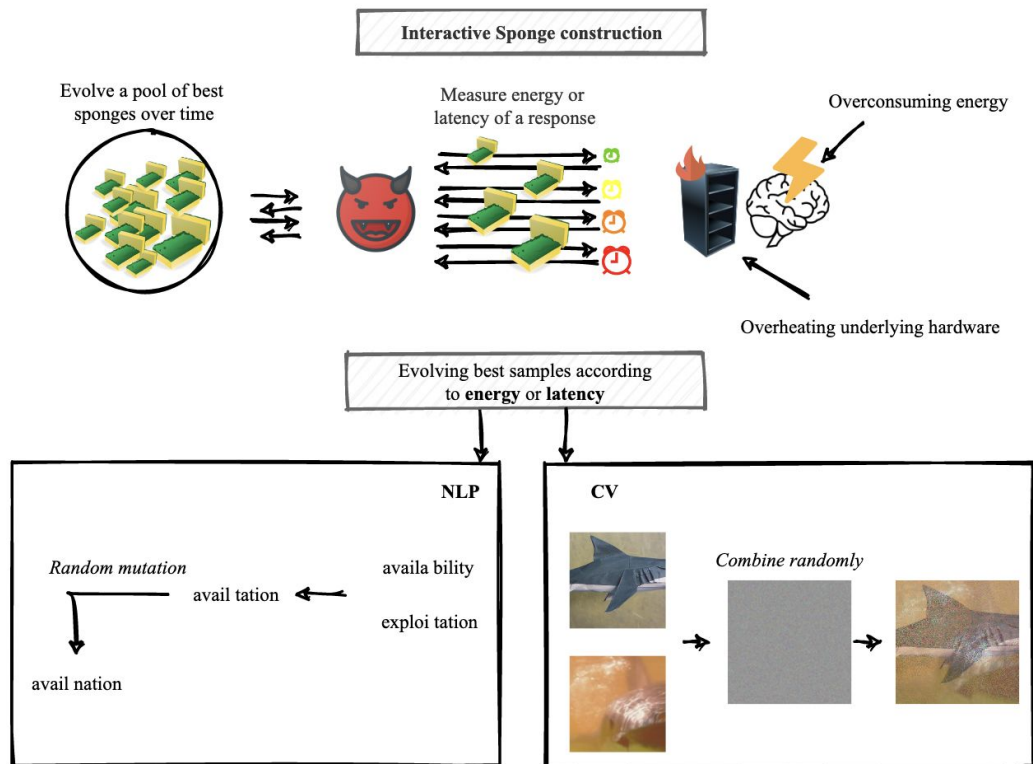
1 error with 7 tokens for input of size 16:

Athazagoraphpbia => ath, az, agor, aph, p, bi, a

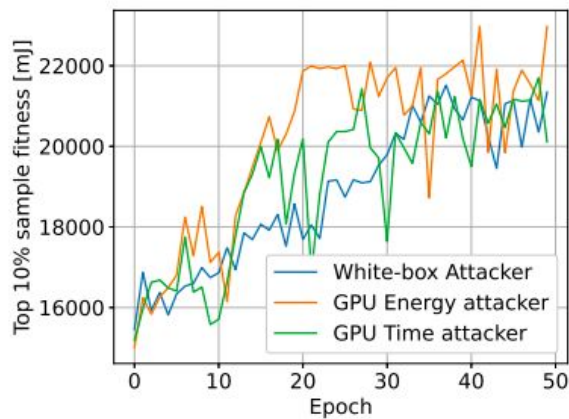
Malicious with 16 tokens for input of size 16:

A/h/z/g/r/p/p/i/ => A, /, h, /, z, /, g, /, r, /, p, /, p, /, i, /

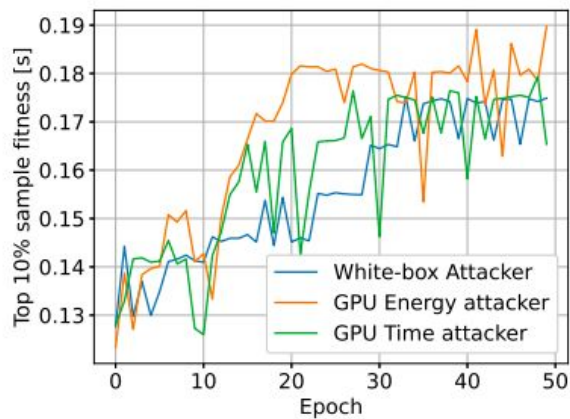
Multiple ways to search for Sponge examples



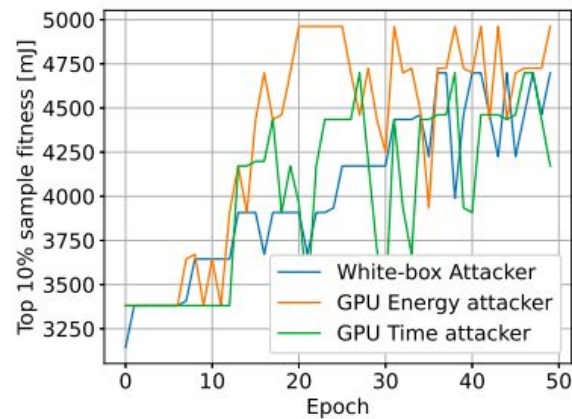
Interactive Black-box attack performance against WMT16 En→Fr



(a) GPU Energy



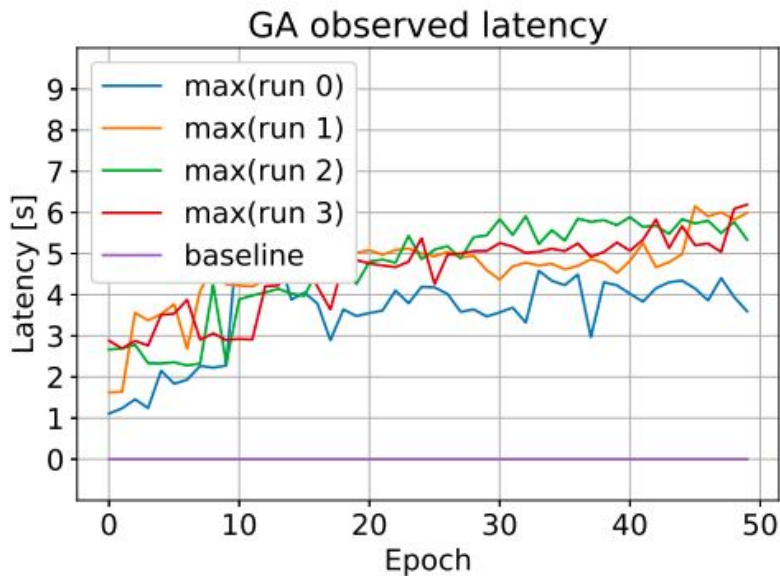
(b) GPU Time



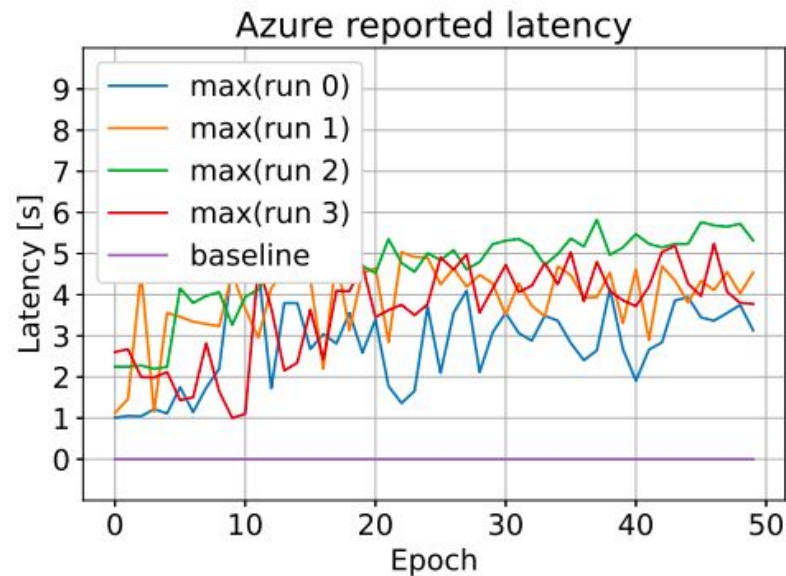
(c) ASIC Energy

Attack works equally as well optimising number of ops, energy and latency.

Microsoft Azure



(a) Requesting server measured



(b) Azure reported

Baseline is at 1ms. Attack performs consistently with multiple restarts.

Conclusions

- It is possible to attack model **availability** in both White and Black-box settings
- Attack can target **hardware optimisations**
 - For some CV tasks we fully negated benefits from acceleration
- Attacks can target **algorithmic complexity**
 - For some NLP tasks we managed to get up to **x30** energy consumption and **x27** time

Conclusions

- Average case is very **different** from **worst case** scenario
- Pipeline **complexity matters**
- **Impact** of ML on climate change might have been **underestimated**
- It is **not clear how to defend** systems against Sponge examples
- **Real-time systems** with ML components **should model availability** adversary