

adE

H % m :

< L . N

οī (

h 0 c[∖

iK l

ON E

/ | ~ X I

R:g D

dvs ny

D ?

Machine Learning for Security and Security for Machine Learning

Nicole Nichols**

Pacific Northwest National Lab

Co-Authors: Rob Jasper, Mark Raugas, Nathan Hilliard, Sean Robinson, Sam Kaplan* Andy Brown*, Aaron Tuor, Nick Knowles*, Ryan Baerwolf*, and Brian Hutchinson**



PNNL is operated by Battelle for the U.S. Department of Energy

WWU*, joint appointee WWU / PNNL**

PNNL-SA-142069



L~1Xv~g) <mark>:"B:- 3</mark>A YกZI CX/ekvB!(>w AR<.74M #+2 R6jcY~



FadE

/]Tb c-o.

ba D?

dvs ny

δī (-×^⊧

h 0 ⁼c[∖

iK l

P ØN≲E w

R:g D

/ | ~ X I

H % m 👳

Two Questions

- Can ML be used in security applications where malicious patterns are not predefined?
- Can ML itself be secure in deployments?





t_}.0

dvs ny

οī (-≥^

f c[·ik l

P CNEE W

R:g D

/ | ~ X I

First Question

- Can ML be used in security applications where malicious patterns are not predefined?
- Can ML itself be secure in deployments?



Two Use Cases: NLP analysis of cyber data for insider threat detection

Neural Fuzzing for accelerating software security assessments

Common Approaches to Insider Threat



t_}.0 FadE/d

H % m :

banD?

dvs ny

οī (=⊗^⊺

iK l

ON E W

/ | ~ X I

R:g D

h 0 c[

Pacific

Northwest





e e

Context

Log Entry

Across Log Entries



PNNL-SA-142069

FadE

h 0 ℃[\





t_}.0
FadE d
FadE d
I Fad

dvs ny

οī (=%^R

o c[∖ i¥ l

ON E W

R:g D

/ | ~ X I

Northwest

RNN Event Model (EM)



Minimize anomaly score: $-\log P(\mathbf{x}_1, \mathbf{x}_2, ..., \mathbf{x}_{T-1}, \mathbf{x}_T)$



FadE d

iY 7^ /]Tb c-o.

H%mg.

ba D?

dvs ny

οī (=×^R

f c[iK l

P 2N E W

R:g D

0

/ | ~ X I

Northwest

Bidirectional RNN Event Model (BEM)







0

c[

Tiered Event Models (T-EM/T-BEM)





t_}.0 FadEZd]A Eu!

iY 7^~ /]Tb c-o.`

H%mg+

ba D? k;nihW P = G

ASV %Z

dvs ny

δī (=≥^⊺

o c[∖ iĶ l

P CNEW

R:g D

/ | ~ X I

C

:1} Z e e

Attention

Attention Mechanism





t_}.0

dvs ny

οī (-<u>≃</u>^

h 0 c[∖

; -it l f

P ONSE W

R:g D

/ | ~ X I

go

Experiment Setup

Data

- LANL cyber security data set authentication logs.
- 0.00007% of events are marked as Red Team activities.

Performance Metric:

Area under the Receiver Operating Characteristic Curve (AUC)

Baseline Comparison

- Baseline models use user-day aggregate statistics.
- Use max event anomaly score for user on that day for language models.
- Also evaluate language models on a per-event basis.



t_}.0 FadE d] E ! iY 7^~ /]Tb c-o.

H%mg+ :l} ZG

ee b@ D? k;nihW # = G

ASV %2

dvs ny

δī (=≈^)

f c[ik l

P CNEW

R:g D

/ | ~ X I

Experiment Results Vs Baseline





fadE d

Y 7 ^ /]Tb c - o .

H%mg+

ba D?

ASV %Z

dvs ny

οī (=≤^

f h 0 c[∖

, iK ≀ f

P CNEE W

R:g D

/ | ~ X I

Word Models

Model	Mean	Max	Min	Std. Dev.
EM	0.968	0.976	0.964	0.005
BEM	0.976	0.981	0.972	0.003

EM with attention

Fixed	0.974	0.976	0.972	0.001
Syntactic	0.972	0.975	0.967	0.004
Semantic 1	0.975	0.980	0.971	0.004
Semantic 2	0.973	0.976	0.968	0.003

model: Semantic I

Higher ROC than the simple Event Model

Tiered LSTM variants

T-EM	0.984	0.989	0.977	0.005
T-BEM	0.987	0.989	0.985	0.002
TA-EM	0.985	0.991	0.979	0.004
TA-BEM	0.988	0.991	0.984	0.003

Best performing single tier



fadE d

Y 7' /]Tb c - o .

H%mg+

ba D? k;nihW

ASV %z

dvs ny

οī (=≤^

; -⊌ ≀_f●

P; 2NgEw

R:g D

/ | ~ X I

f h 0 c[∖

Pacific Northwest

Word Models

Model	Mean	Max	Min	Std. Dev.
EM	0.968	0.976	0.964	0.005
BEM	0.976	0.981	0.972	0.003

EM with attention

Fixed	0.974	0.976	0.972	0.001	
Syntactic	0.972	0.975	0.967	0.004	
Semantic 1	0.975	0.980	0.971	0.004	
Semantic 2	0.973	0.976	0.968	0.003	

Tiered LSTM variants

T-EM	0.984	0.989	0.977	0.005
T-BEM	0.987	0.989	0.985	0.002
TA-EM	0.985	0.991	0.979	0.004
TA-BEM	0.988	0.991	0.984	0.003

Attention models perform only marginally worse than **bidirectional models**

Global Average Importance of Fields



fadE d

H % m 👳 + 1} Z

ba D? k;nihW

ASV %z

dvs ny

οī (-__^

0

f

Input token

c[

if l

P 2N E w

R:g D

/ | ~ X I

20

e e

Eu! iY 7 ~~~~ /]Tb c - o .

C

G a S









PNNL-SA-142069



t_}.0

dvs ny

οī (-≥^

f c[·ik l

P CNEE W

R:g D

/ | ~ X I

First Question

- Can ML be used in security applications where malicious patterns are not predefined?
- Can ML itself be secure in deployments?



Two Use Cases: NLP analysis of cyber data for insider threat detection

Neural Fuzzing for accelerating software security assessments



Goal

Accelerate search for unique code paths that could reveal faults

Assumptions

Faults are more likely to exist on untested / unexplored code paths

Shorter paths are easier to test / explore than longer paths

Approach

Augment American Fuzzy Lop (AFL) with LSTM and GANS generated seed files to accelerate search.

< L . N

dvs ny

οī (−≈^≈

P 2NEW

R:g D

/ | ~ X I





Approach



s:dvs ny

οī (=×^R

f c[∖ iK ^lf

P 2NEW

R:g D

/ | ~ X I

Additional Seed File of Unique Code Paths



fadE d

i Y 7 ^ /]Tb c - 0.

H % m 👳 🖣

ba D?

dvs ny

οī (=≥^

نلا ۱ f

P CNEE W

R:g D

/ | ~ X I

f h 0 c[∖

Analysis of Seed Files

Class C	C	L(C)	% Unique	$\mu(L(C))$	$\sigma(L($
AFL seed	38384	31212	0.813	26.968M	33.9:
Rand seed	19824	485	0.024	2.602M	724.6
LSTM seed	20000	1921	0.096	2.596M	8.68
GAN seed	20000	119	0.006	2.593M	1.84

- The seed themselves are **not** what we are interested in measuring
- They only provide a set of *initial conditions* for AFL
- Interestingly LSTM and GAN do have as much variance as using purely random seeds

58M 674K 87K 41K



fadE d

iY 7^ /]⊤b c-o. H%mg+

1} Z

: b@ D? :k;nihW

ASV %2

dvs ny

οī (=⊻^)

; -₩ ι_f●

P 2NEW

R:g D

/ | ~ X I

f o c[\

~ S

Time Analysis of Sustained Run

- Both LSTM and GAN outperform random sampling for discovering new unique code paths.
- GAN 11 % faster / random
- LSTM 8% faster over random

Class	Files	% new	sec/path	NRate
Rand	1231	0.9017	214.478	1.00
LSTM	1251	0.8984	197.130	1.08
GAN	1240	0.8694	191.893	1.11



t_}.0 FadEd

iY 07^ /]Tb c - 0 .

H % m 👳 +

ba D? k;nihW

dvs ny

οī (=≾^

f c[iK l

P 2N E W

R:g D

/ | ~ X I

) ~ S ASV %2

1} Z

Code Path Length of Sustained Run

- length of unique code paths using GAN was 13.84% longer than a strategy based on randomly sampling.
- length of unique code paths using **LSTM** was **4.60%** longer than a strategy based on randomly sampling.

Class	μ(<i>L</i> (<i>C</i>))	σ(<i>L</i> (<i>C</i>))
Rand	25.373M	3.339M
LSTM	26.541M	3.385M
GAN	28.885M	3.456M





t_}.0 FadE d

> /]Tb c-o.

b@ D? k;nihW

dvs ny

οī (=⊻^R

it ۱ f

P CNEE W

R:g D

/ | ~ X I

h 0 c[∖

H % m 👳 🕈

Second Question

- Can ML detect malicious behavior without predefined patterns?
- Can ML itself be secure in deployments?



Adversarial Machine Learning





Digital Attacks:

Pacific

< L . N

dvs ny

οī (=≊^

′h 0 ⁼c[∖

iK l

P 2N E W

R:g D

/ | ~ X I

Northwest

Direct access to maliciously modify model, input features, or database of training examples.

Physical Attacks:

A physical object is added or modified in the scene being evaluated.

A Cambrian Explosion of Machine Learning Research Topics



"The Malicious Use of Artificial Intelligence: Forecasting, Prevention, and Mitigation" arXiv:1802. 07228 (2018)

Trimps-Soushen

ML+neuroscience

Accountability

and Transparency

Goodfellow 2017



FadE

ASV %Z

dvs ny

οī (-<u>≥</u>^

h 0 c[∖

i¥ ۱

ON E W

R:g D

/ | ~ X I

Why is Machine Learning Vulnerable?

Known general ML fragilities...

- Every model has a decision boundary; manifolds can be wrinkly
- Not enough training data to resolve boundary cases (chihuahua/muffin)
- Not all classes are separable
- High dimensional space is not intuitive
- Decisions are hard to understand
- Poisoned the training data (GIGO)
- Compromise of privacy in the training data
- Denial of service, output corruption, hacks...

Additional DL vulnerabilities

- No mystery: DL models are the approach of choice for many problems
- Limited diversity: A few training sets, standard architectures, standard models
- Spotlight: Many researchers are publishing DL-specific attacks







t_}.0 FadE

D ?

dvs ny

οī (=<u>×</u>^

c [

i۴ ۱

N E w

R:g D

/ | ~ X I

Decision Boundaries

- Data driven models are only as good as their data
- Training data cannot fully define a decision boundary
- What is going on with vulnerability and misclassification:



Feinman et al. "Detecting adversarial samples from artifacts." arXiv preprint arXiv:1703.00410 (2017).

don't know what I am.



FadE

/1Tb

ba D?

A 5 V 1 % 2

dvs ny

ī (=×

0 €c[\

ik l

R:g D

/ | ~ X I

< L . N

H%m ≘ ∙

Attacks in the Digital Domain

- <u>Adversarial Example</u> model input an attacker has intentionally designed to cause the model to make a mistake.
- Distance in feature space is not always intuitive.
- Numerous ways to craft adversarial examples.



d difference

0.0081



Schoolbus



Perturbation



Ostrich

Szegedy et. al., "Intriguing properties of neural networks" arXiv preprint arXiv:1312.6199 (2013) Zheng, Stephan, et al. "Improving the robustness of deep neural networks via stability training." *Proceedings of the ieee conference on computer vision and pattern recognition*. 2016.

PNNL-SA-142069



0.1038

0.1011

255



t_}.0 FadE d

iY 7^ /]Tb c-o.

H % m 👳 🕈

: b@ D? k;nihW

ASV %Z

dvs ny

οī (-≊^

c [

iK l

R:g D

/ | ~ X I

< L . N

- S

1} Z

Attacks in the Physical World

Physical attacks span significant range of perception and detectability

- Targeted attacks
- 2 and 3D object construction
- Digital formulation for physical world deployment (White box attacks)

Camouflage Camouflage Art Camouflage Art Distance/Angle (LISA-CNN) **Right Turn** (GTSRB-CNN) Graffiti 5' 0° 5' 15° 10' 0° 10' 30° 40' 0 Targeted-Attack Success 73.33% 66.67% 100% 80%



[1]Athalye, Anish, and Ilya Sutskever. "Synthesizing robust adversarial examples." *arXiv preprint arXiv:1707.07397*(2017).

[2]Sharif et. al., "Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition" Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security. ACM, 2016, pp. 1528–1540.

[3] Evtimov, Ivan, et al. "Robust physical-world attacks on machine learning models." arXiv preprint arXiv:1707.08945(2017).

[4] Brown, Tom B., et al. "Adversarial patch." arXiv preprint arXiv:1712.09665 (2017).









classified as turtle classified as rifle classified as other

PNNL-SA-142069



Decision boundaries for models of the same class are likely to be similar

Papernot et. al. "Transferability in machine learning: from phenomena to black-box attacks using adversarial samples" CoRR, arXiv:1605.07277 (2016). PNNL-SA-142069

0.82

DNN

kNN- 11,75

Source |

12.22

42,89

LR

t_}.0 FadEd

Y 7^

/]Tb c - o .

ba D?

k;nihW

A5v %z

< L . N

dvs ny

οī (=≥^

iK l

N E W

R:g D

/ | ~ X I

h 0 €c[\

H % m 👳 •

models trained on target

4.32	79.31	8.36	20.72 -
1,43	87,42	11,29	44,14 -
00.0	80.03	5.19	15.67 -
8.85	89.29	3.31	5.11 -
2,16	82,95	41,65	31,92 -
SVM	DT	kNN	Ens.

Target Machine Learning Technique



t_}.0 FadE

3:dvs ny

οī (=≚^

iK l

ON E V

/ | ~ X I

R:g D

h 0 c[\

Experiment Inception

Goal 1: Can light cause misclassification of 2D print images Goal 2: Can light cause misclassification of 3D objects Goal 3: What is the stability of this approach.



Inspired by :

Kurakin, A., Goodfellow, I., and Bengio, S. "Adversarial examples in the physical world." arXiv preprint arXiv:1607.02533 (2016).

PNNL-SA-142069



t_}.0 FadE/d

> Y 7^ /]Tb c-o.

H % m 👳 🖣

b@ D? k;nihW

A5v %2

dvs ny

οī (=≤^

iK l

ON E W

R:g D

/ | ~ X I

h 0 €c[∖

~ S

Projecting Trouble- 2D Experiments

Transient physical attacks

CIFAR10 dataset and pre-trained ResNet38 classifier.

Non-targeted and false negative attack

Differential Evolution, white-ish box attack (crafted to the image but without knowledge of classification model)





t_}.0 FadE

H % m 🗉

ba D?

dvs ny

 $\overline{\iota}$ (= \times *)

h 0 c[∖

iK l −iK l

ON E V

/ | ~ X I

R:g D

3D Presentation Experiment

Prob.	p(airplane)	p(automobile)	p(bird)	p(cat)	p(deer)	p(dog)	p(frog)	p(horse)	p(ship)
Original	0	89%	0	0	0	0	0	0	0
Attacked	0	43%	0	0	0	0	0	0	0

- 3D attacks can be successful
- In CFAR10, trucks are semi-trailers, fire trucks, etc, thus bigger difference to shift.
- Non-targeted, transient attack.











Northwest NATIONAL LABORATOR

3D CIFAR Experiment

- One example of each CIFAR10 class.
- Environmental control
- Additional Baseline attacks (white light, random square, DE square)
- ImageNet co-classification



dvs ny

οī (=≥^

iK l

ON E I

/ | ~ X I

R:g D



t_}.0
FadE d
FadE d
I E ! 9
iY 7^~
/]Tb
c-o.`
H%m +
1} ZG
ee q
b@ D?
k;nihW
= G

fulfa |j | T!)∼S |A5vj%z h

dvs ny

f o c[∖

P CNEE W

/ | ~ X I

R:g D

Pacific Northwest

Results

CIFAR Class	Experiment Condition	Mean	Median	SD	Var	Min	Max	Δ Mean	Δ Median
Airplane	Baseline	1.000	1.000	.000	.000	1.000	1.000	.000	.000
	White Light	.151	.101	.198	.039	.017	.997	.849	.899
	Random	.114	.105	.088	.008	.022	.445	.886	.895
22	Diff Evolution	.133	.112	.087	.007	.014	.459	.867	.888
Automobile	Baseline	1.000	1.000	.000	.000	1.000	1.000	.000	.000
and the second s	White Light	1.000	1.000	.000	.000	.999	1.000	.000	.000
	Random	1.000	1.000	.000	.000	.999	1.000	.000	.000
<u>T</u> T	Diff Evolution	1.000	1.000	.000	.000	1.000	1.000	.000	.000
Bird	Baseline	1.000	1.000	.000	.000	1.000	1.000	.000	.000
	White Light	1.000	1.000	.002	.000	.993	1.000	.000	.000
	Random	1.000	1.000	.000	.000	1.000	1.000	.000	.000
26	Diff Evolution	1.000	1.000	.000	.000	_999	1.000	.000	-000
Cat	Baseline	.990	.991	.004	.000	.979	.996	.000	.000
	White Light	.009	.008	.005	.000	.000	.020	.981	.983
	Random	.011	.007	.012	.000	.001	.047	.979	.984
22	Diff Evolution	.023	.017	.019	.000	.002	.124	.967	.974
Deer	Baseline	.999	.999	.000	.000	.999	1.000	.000	.000
	White Light	.516	.516	.145	.021	.242	.997	.483	.483
	Random	.545	.507	.155	.024	.327	.871	.454	.492
22	Diff Evolution	.473	.467	.130	.017	.144	.829	.526	.532



t_}.0
FadE d
P] E ! 9
iY 7^~
/]Tb
c-o.`ZK
H%m +
:l} ZG

ee d ba D?

k;nihW

A5vj%z

dvs ny

οī (-×^⊧

, ∙ų ≀ f

P CNEW

R:g D

/ | ~ X I

f h 0 c[∖

q

Pacific Northwest

Results

CIFAR Class	Experiment Condition	Mean	Median	SD	Var	Min	Max	Δ Mean	Δ Median
Dog	Baseline	.993	.993	.003	.000	.986	.996	.000	.000
	White Light	.512	.499	.088	.008	.390	.695	.481	.494
	Random	.482	.497	.123	.015	.136	.753	.511	.496
	Diff Evolution	.386	.388	.088	.008	.123	.601	.606	.605
Frog	Baseline	.888	.888	.025	.001	.842	.933	.000	.000
	White Light	.008	.008	.003	.000	.000	.015	.881	.880
	Random	.030	.011	.076	.006	.004	.360	.858	.877
	Diff Evolution	.071	.038	.093	.009	.005	.576	.817	.849
Horse	Baseline	1.000	1.000	.000	.000	1.000	1.000	.000	.000
	White Light	.999	1.000	.001	.000	.993	1.000	.000	.000
	Random	1.000	1.000	.000	.000	1.000	1.000	.000	.000
	Diff Evolution	1.000	1.000	.000	.000	1.000	1.000	.000	.000
Ship	Baseline	1.000	1.000	.000	.000	1.000	1.000	.000	.000
	White Light	1.000	1.000	.000	.000	1.000	1.000	.000	.000
	Random	1.000	1.000	.000	.000	1.000	1.000	.000	.000
	Diff Evolution	1.000	1.000	.000	.000	1.000	$_{-1.000}$.000	.000
Truck	Baseline	1.000	1.000	.000	.000	1.000	1.000	.000	.000
	White Light	.832	.832	.052	.003	.729	1.000	.168	.168
	Random	.818	.819	.072	.005	.634	.970	.182	.180
	Diff Evolution	.826	.839	.088	.008	.507	.949	.174	.161

Table 1: Classification statistics for baseline and attacked CIFAR figures.



Results

- Extreme variability between target class susceptibility to attack.
 - 6 of 10 classes were susceptible to light based attacks.

	Ave(Δ Mean)	Ave(Δ Median)		
White Light	0.641	0.651		
Random	0.645	0.654		
Diff Evolution	0.660	0.668		

- White light was similarly effective to random squares and DE.
- Rotation, lighting, and scale invariance of classification models are significant considerations.



dvs ny

οī (-≥^

h 0 c[∖

iK l

R:g D

/ | ~ X I

to tacks.

and DE.



FadE

dvs ny

οī (=≈^r

f c[iK l

P CNEW

R:g D

/ | ~ X I

Conclusions

ML for Security:

- Deep learning techniques can be used to enhance and accelerate a variety of security based applications.
- Pre-knowledge of patterns is not necessary in insider threat detection or software fuzzing.

Security for ML:

- Most off the shelf models are insufficiently resilient to real world invariance.
- An increasing range of digital and physical security gaps are being identified in ML models.

Security of the model itself needs to be considered, particularly when deploying ML for Security.



FadE d

ba D?

ASV %z

dvs ny

οī (-×^)

f c[iK l

P NEE W

R:g D

/ | ~ X I



Recurrent Neural Network Language Models for Open Vocabulary Event-Level Cyber Anomaly Detection https://arxiv.org/pdf/1712.00557.pdf

Deep Learning for Unsupervised Insider Threat Detection in Structured Cybersecurity Data Streams https://arxiv.org/pdf/1710.00811.pdf

Faster Fuzzing: Reinitialization with Deep Neural Models https://arxiv.org/pdf/1711.02807.pdf

Projecting Trouble: Light Based Adversarial Attacks on Deep Learning Classifiers https://arxiv.org/abs/1810.10337

Code available at: https://github.com/pnnl/safekit



t_}.0 FadE d

iY 7^· /]Tb c-o.

H%mg+ l} Z0 ee

: b∂ D? k;nihW

ASV %z

dvs ny

οī (-∞^

o c[\ iK ^lf

P 2N E W

R:g D

/ | ~ X I

G

Thank you

PNNL-SA-142069



ivC5 YGdIF L~1Xv~g)i E"B;- 3A YaZF CX/ekvB!(>W AR<.74MKB D = hEm /)f<H"U< a X h H #+2] b > u zl@ayX.] 1 W t (R6jcY~+je]c8r / **G** 6