# LEGAL DISCLAIMER

# ARTIFICIAL INTELLIGENCE

**SOFTWARE**

**HARDWARE**

**COMMUNITY**

# ARTIFICIAL INTELLIGENCE & SECURITY

## SECURITY FOR AI

## AI FOR SECURITY



INTEL

XEON inside

NERVANA inside

STRATIX 10 inside

CORE i7 inside

MOVIDIUS inside

* Representation image only

# Security for AI considerations

# Clever Hans

"We have reached the point where machine learning works, but *may easily be broken*"

Nicolas Papernot, Google PhD Fellow in Security
Ian Goodfellow, Research scientist at Google Brain

# Artificial intelligence

**Machine Learning**

- Study many images labeled as flamingo

- Identify the flamingo in the image

**Deep Learning**

- Study many images

- Identify the flamingo, hedgehog, etc.

**Artificial Intelligence**

- Is she hugging the flamingo, or playing cricket?

- Is she happy, sad?

ARTIFICIAL INTELLIGENCE

MACHINE LEARNING

DEEP LEARNING

# Training Compute is not the bottleneck, data is

**40%**
Pre-processing

**20%**
Compute

**40%**
Optimization and Deployment

Preparing the data for analysis, Finding the right model for the problem

Training the model

Optimizing and deploying the model

# Inference is a different story!

© 2018 Intel Corporation

# Pre-processing – it's a complicated journey

**01** Normalization

**02** Deduplication

**03** Noise reduction

**04** Sanity checks

**05** Labeling

© 2018 Intel Corporation

**Machine Learning**

Features

Machine Learning Model

Pre-processing

ML Optimization

Backdoor

Retraining

Accuracy

Training Data

Data Poisoning

Malicious Retraining

Inference

Malicious Data

Real World Data

Machine Learning Model

Deployment

Information Leaks

Rouge Expert

Action

Expert          Rule Engine

Cross model attack vectors

IP Extraction

Model Modification

Output

© 2018 Intel Corporation

# Backdoors

**Validation** of ML is an open problem

We don't have a method for
**detecting backdoors**

Reverse engineering, code review are
not applicable to ML

# IP Extraction

© 2018 Intel Corporation.

IP can be stolen using public APIs

Reverse engineering or device access not required

# Evading next generation AV using AI



- Static machine learning model trained on millions of samples

EXE → Machine Learning Model → score=0.75 (**malicious**, moderate confidence)

- Simple structural changes that don't change behavior
  - unpack
  - '.text' -> '.foo' (remains valid entry point)
  - create '.text' and populate with '.text from calc.exe'

EXE → Machine Learning Model → score=0.49 (benign, just barely)

# Turtle or a Rifle?

# Adversarial Audio



"okay google without the dataset the article is useless"

"okay google browse to evil dot com"

"okay google browse to evil dot com"

Adversarial Verdi's Requiem

## You can fool home automation, smartphones and other devices

# Supply chain security – in AI

# Information Disclosure by



Figure 10: Reconstruction of the individual on the left by Softmax, MLP, and DAE.

[Tromer, Zhang, Juels, Reiter, R. 2016]
[Fredrikson, Jha, R. 2015]

What about

privacy?

# Privacy leaks? Not yet, but soon…

Training

Inference

Risk: 7.4%    Risk: 35.3%

© 2018 Intel Corporation

# Privacy leaks? Not yet, but soon…

Training

Inference

Risk: 96.2%

© 2018 Intel Corporation

# What are the issues?

- Data vs. Information
- New Threat Vectors
- Upcoming Attacks
- Unique Gaps

© 2018 Intel Corporation

# AI Security: Unique Gaps

**IP protections are early stage (at best)**

# AI Security: Unique Gaps

**AI Validation** **is a major issue**

**Pretty clear if the AI does what it claims, does it do more?**

**Will it fail unexpectedly?**

# AI Security: dynamic systems

**You may end with a different system** then what you started with

# AI Security: Unique Gaps

**Humans in the loop pose a security risk, we don't have sufficient controls during Machine Learning development**

So, what can we do?

© 2018 Intel Corporation.

intel®

# Our Recommendations

1. Start having conversations about Security and AI
2. Machine learning needs to be protected against attackers
3. Checks and balances, don't trust blindly
4. Protecting AI is a journey, join us in our Keynote on Wednesday

Set an agenda in your org to discuss your AI and Data security plans

Schedule meetings with your teams to discuss your posture on AI security

Reach out to us to discuss these issues after this talk

# Remember Mr. ed the talking horse?

# Please reach out for more information

Guy Barnhart-Magen

guy.barnhart-magen@intel.com

@barnhartguy (twitter)