# LEGAL DISCLAIMER

# WHO AM I?

- **Guy Barnhart-Magen**

  - **@barnhartguy** on Twitter

- Security Researcher, Manager, Presenter

- Interests:

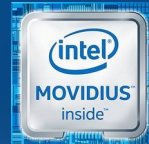  Crypto, Embedded systems, Artificial Intelligence, System/Product security

# 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

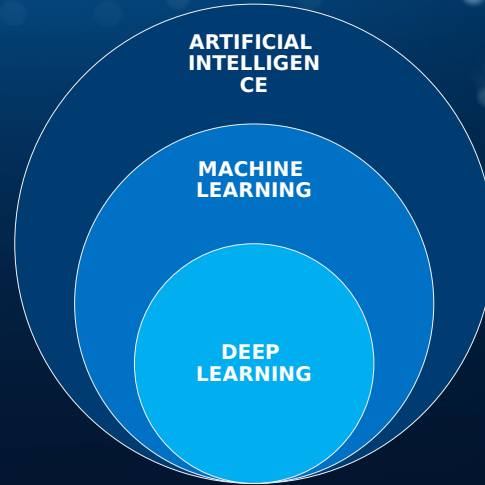# Artificial intelligence

**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!

# Pre-processing – it's a complicated journey

Normalization **01**

**02** Deduplication

Noise reduction **03**

**04** Sanity checks

Labeling **05**

intel®

© 2018 Intel Corporation

© 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

IP can be stolen using public APIs

Reverse engineering or device access not required

intel®

# Different view points

**What Microsoft Sees:**

Oh #$%©! 2 Out of 18 Million Across Most of the Corporate World Have No Phish Protection.

What Y

Thanks

Office 3
© 201

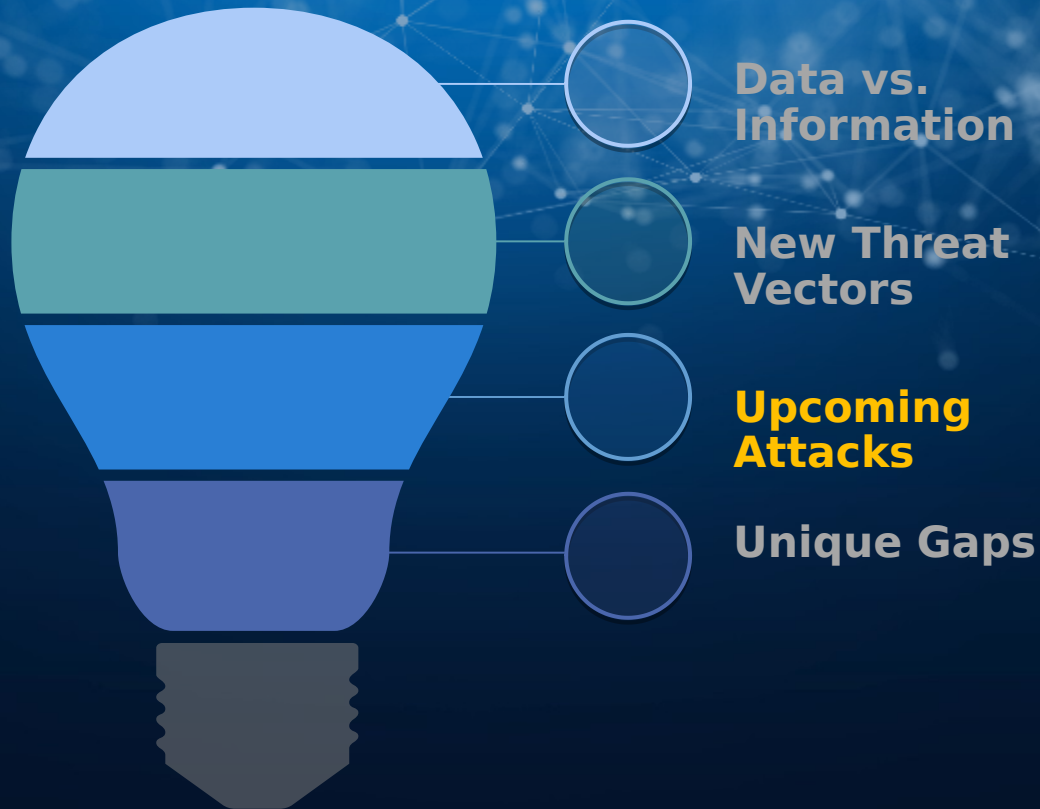Rewriting "Microsoft Security Team" in HTML eMail:

Micro`<span style='font-size:0'>`processors run optimize`</span>`soft`<span style='font-size:0'>`ware to store your secrets`</span>`Secur`<span style='font-size:0'>`ely. It is also good for system integr`</span>`ity`<span style='font-size:0'>`, thanks to our`</span>`Team
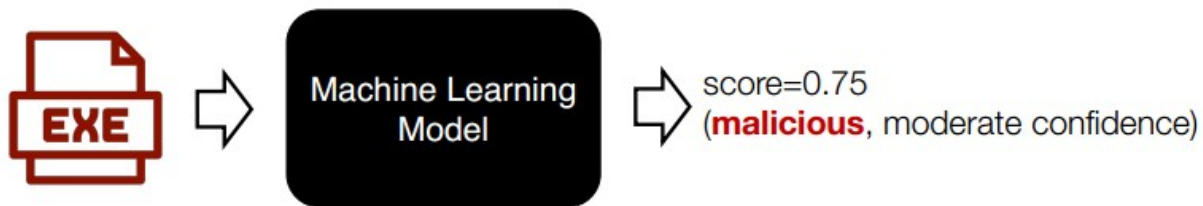
Scanners read unstructured text as:

Microprocessors run optimize software to store your secrets Securely. It is also good for system integrity, thanks to our Team.
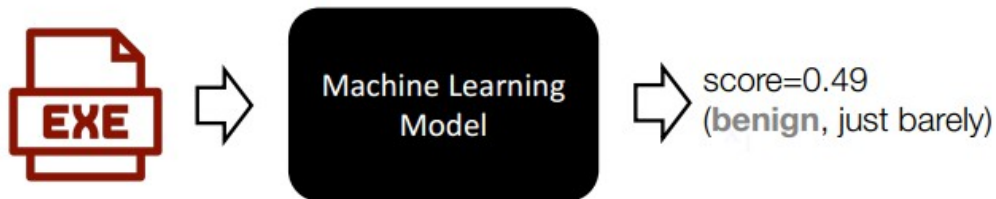
What E

Thname
nameb
Oname
© 201

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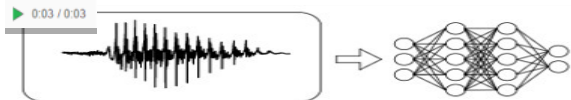
- Simple structural changes that don't change behavior
  - unpack
  - '.text' -> '.foo' (remains valid entry point)
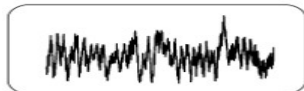  - create '.text' and populate with '.text from calc.exe'

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

https://www.youtube.com/watch?v=FGCIe6T0Jpc

# Turtle or a Rifle?

# Adversarial Audio



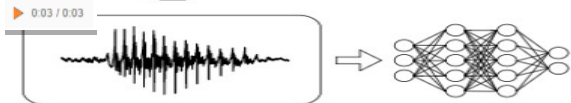"okay google without the dataset the article is useless"

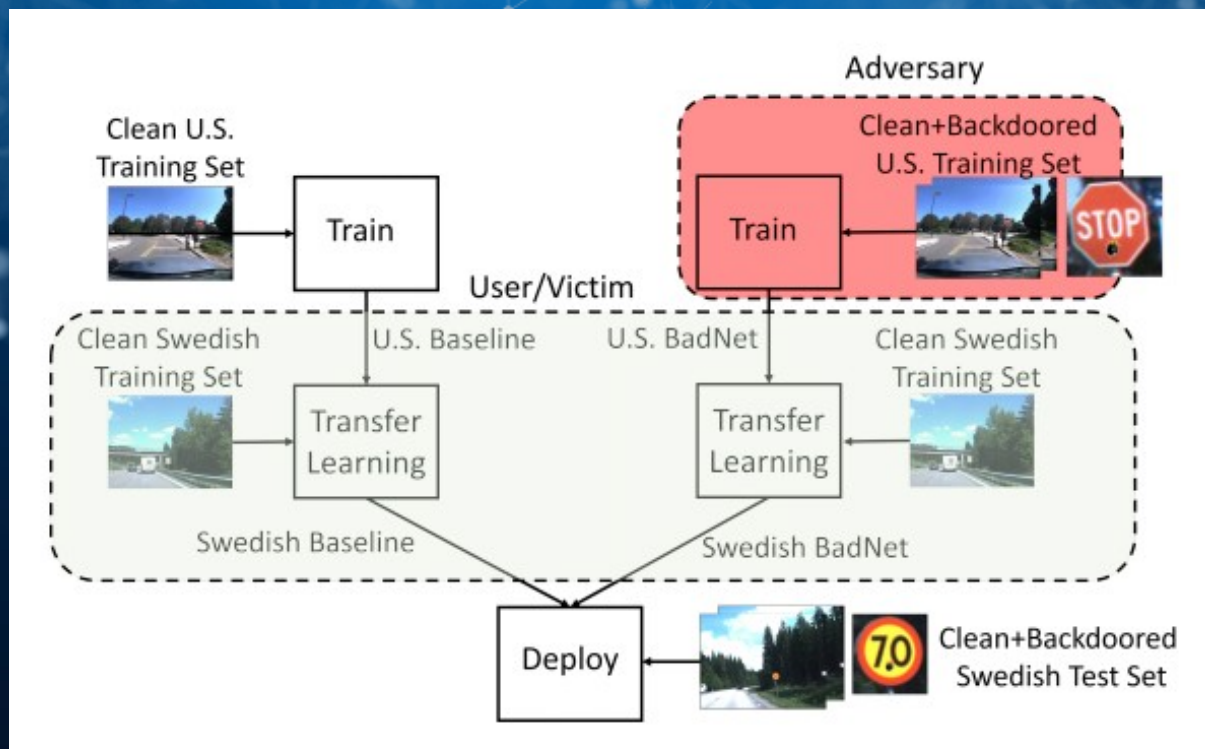"okay google browse to evil dot com"

"okay google browse to evil dot com"

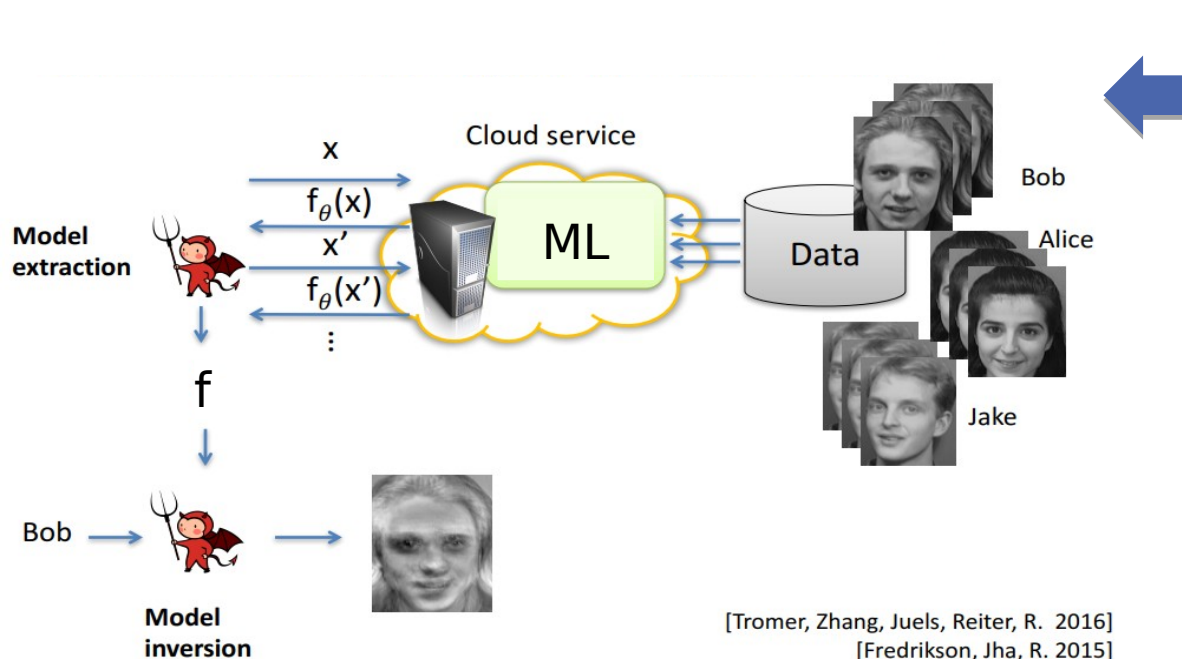Adversarial Verdi's Requiem

# You can fool home automation, smartphones and other devices
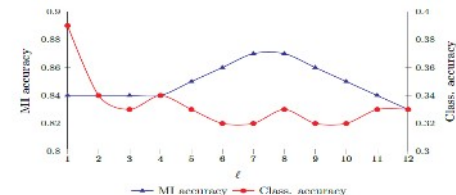
# Supply chain security – in AI



https://machine-learning-and-security.github.io/papers/mlsec17_paper_51.pdf

# 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]

# Privacy leaks? Not yet, but soon...

Training

Inference

Risk: 7.4%    Risk: 35.3%

# 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: Unique Gaps

You shouldn't **trust** the data, even if collected securely, the **data might be**

# AI Security: dynamic systems

**You may end with a different system than 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?

# 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

Reach out to us to discuss these issues after this talk

# Remember Mr. ed the talking horse?

Any Questions?

@barnhartguy

https://media.giphy.com/media/ejwFX1DPsfqec/giphy.gif