RSA®Conference2018

San Francisco | April 16–20 | Moscone Center

SESSION ID: TV-T02

# AI DECEPTION: FOOLING (ARTIFICIAL) INTELLIGENCE IS EASIER THAN YOU THINK

**Itsik Mantin**

Lead Scientist
Imperva

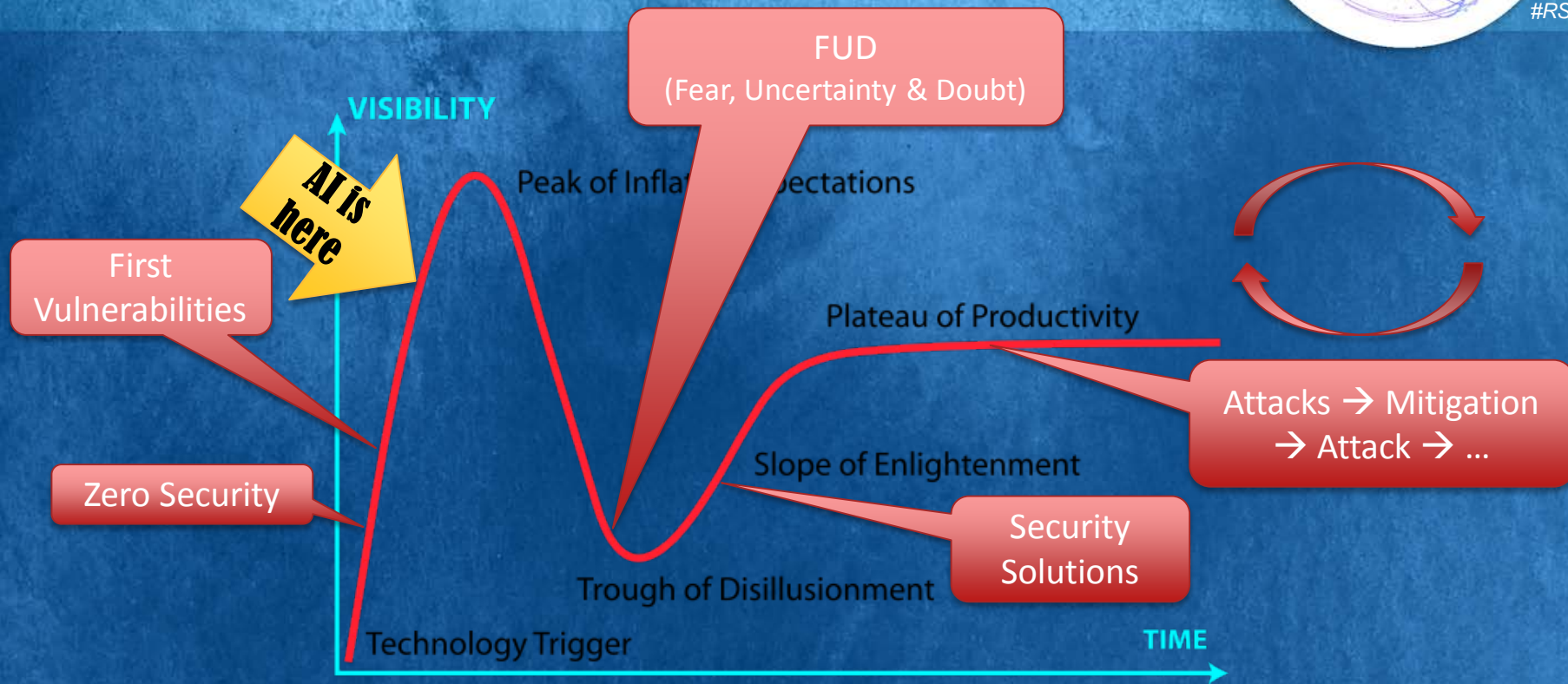**IMPERVA**

RSAConference2018

# Security and the Hype-Cycle

**VISIBILITY**

**AI is here**

**FUD**
**(Fear, Uncertainty & Doubt)**

Peak of Inflated Expectations

Plateau of Productivity

**First Vulnerabilities**

**Attacks → Mitigation → Attack → ...**

Slope of Enlightenment

**Zero Security**

**Security Solutions**

Trough of Disillusionment

Technology Trigger

**TIME**

**IMPERVA**®

3

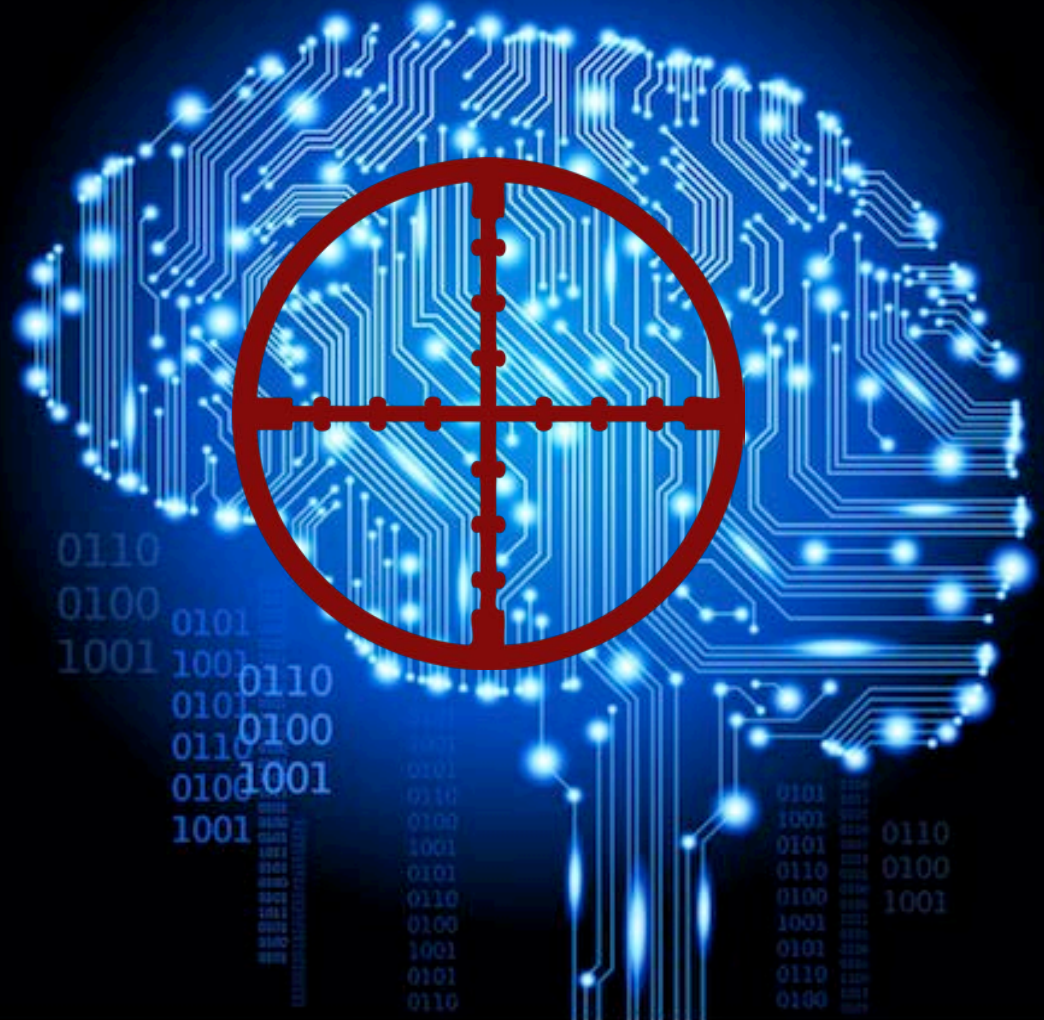RSA Conference2018

No Exemptions for AI!

# The Australian Challenge

Fooling AI by ~~Adversaries~~ *Innocent*

Volvo admits its self-driving cars are confused by kangaroos

Swedish company's animal detection system can identify and avoid deer, elk and caribou, but is yet to work against the marsupials' movements

**Innocent AI Deception**
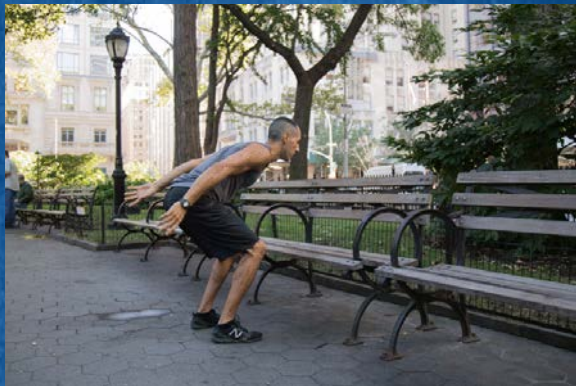
Kangaroos ... Hmmm..

# Adversarial Thinking

|  | Builder | Adversary |
|---|---|---|
| **Primary Focus** | What happens in case of normal input | What may happen in case of anomalous input? |
| **Failure in rare coincidence** | Something I can <u>ignore</u> | Something I can <u>abuse</u> |

**IMPERVA**®

RSA Conference2018

# Blind Spots

Artificial Kangaroo



**IMPERVA**®

RSAConference2018

# Getting Hit by an Ostrich
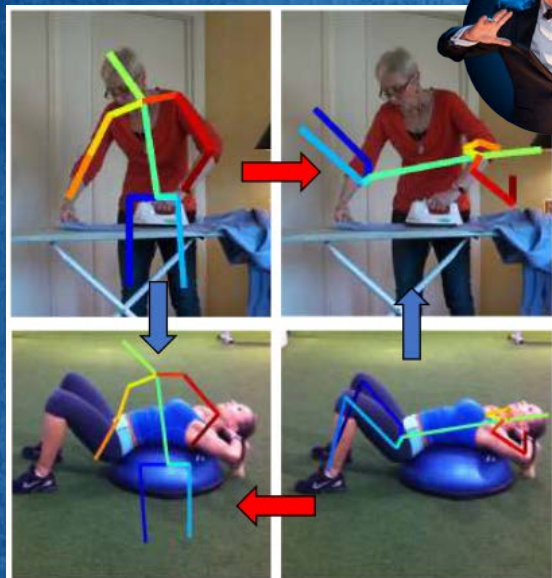## Adversarial Examples

Bar-Ilan University
אוניברסיטת בר-אילן

**Houdini Research**

Pose Estimation

Speech Recognition

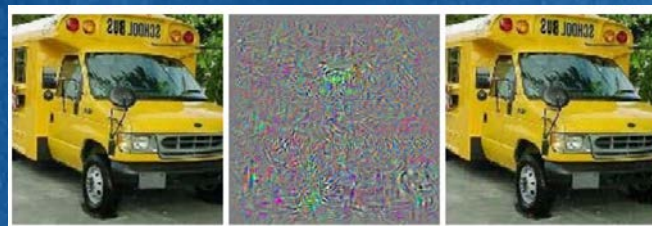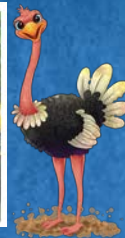Semantic Segmentation

Visual Object Recognition

**Panda**
(57.7%)

**Gibbon**
(99.3%)

$+ .007 \times$

**School Bus**

**Ostrich**

IMPERVA®

RSAConference2018

# Define AI Deception

- Given an A/B classifier, and given a sample X correctly classified as A, attacker generates a sample X' that:
1) Has **same-essence** as X, and
2) Classified as B

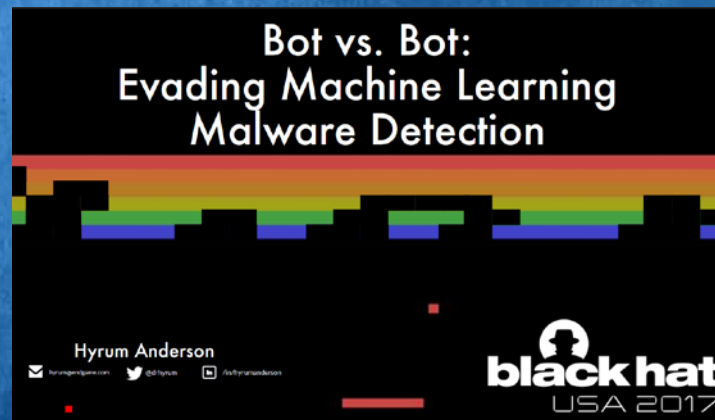God does not play dice…

…but the devil does

**IMPERVA**®

RSAConference2018

# AI Deception Threat

| | Same-Essence | AI Deception |
|---|---|---|
| **Image** | Identical for human viewer | Different objects detected |
| **Video** | | Different scene detected |
| **Voice** | Identical for human listener | Different speech detected |
| **Software** | Same malicious function | Classified as non-malicious |



Bot vs. Bot:
Evading Machine Learning
Malware Detection

Hyrum Anderson

black hat USA 2017

# AI Deception – More Threats
## Propagation to the Real World

| Domain | AI Engine | AI Circumvention |
|--------|-----------|------------------|
| **Surveillance and Control systems** | **Face detector** | **Embed face patterns that prevent correct detection** |
| Finance | Stock prediction | Cause stock patterns that imply positive prediction |
| Text analyzers | Translation Engine / Topic Extraction | Embed text patterns that prevent correct analysis |
| E-Commerce | Customized pricing | Force user profile or behavior that implies cheaper prices |

1.0    0.01

**IMPERVA**®

RSAConference2018

# Are We Better?

## HI Deception

- Human Intelligence Deception

**Wealth Prediction**

Suit

Disrespect money

Well maintained

Cigar

≥10M

Arrogant Expression

Balance Summary [?]

-$888,871.91 Available Balance as of today

IMPERVA®

RSAConference2018

# Risk Mitigation

- Attackers
- Threats

- Attack Vectors
- Attacker's reach

Threat Analysis

Security Analysis

Mitigation

Model Analysis

- Model hardening
- Input sanitization

- Robustness

**IMPERVA**®

**14**

RSAConference2018

# The AI'ker's Guide to the (Cyber-Security) Galaxy
AI in Cyber Security Applications

**The Model**
- Prefer robust models
- Prefer explainable models

**Training**
- Sanitize training data

**Usage**
- Prefer internal (hidden outcome)
- Avoid raw input

**Threat Detection**
- Prefer positive security
- Combine with other mechanisms

# Summary

- AI is awesome technology

- AI fails in adversarial settings

- AI Deception gets only little industry attention

- The threats are acute and critical

- Mitigation in many cases is not trivial

**IMPERVA**®

RSA Conference2018

# Apply What You Have Learned Today

- Next week you should:
  - Identify critical AI usages within your organization and your roadmap

- In the next three months you should:
  - Carry out security modeling for AI usages (threats, adversaries, attack vectors)
  - In case there is significant threat and viable attack vectors, build mitigation plan. Focus on critical easy-to-exploit vulnerabilities.

- Within six months you should:
  - Execute at least the critical part of mitigation plan

**IMPERVA**®

RSA Conference2018