# The History of Adversarial AI

Eugene Neelou & Alex Polyakov

Adversa AI, Israel

# Alex Polyakov

- CEO & Founder: Adversa.AI
- 18 years in Cybersecurity, 6 years in AI
- Member: Forbes Technology Council
- Author: 2 books, First Adversarial ML MOOC
- Speaker: 100+ conferences in 30+ countries
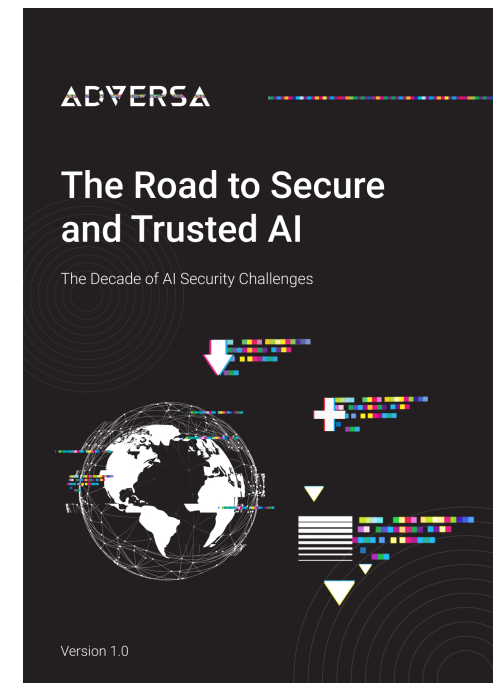- Interests: Trusted AI, SynthBio, Neuroscience

# Eugene Neelou

- CTO & Co-Founder: Adversa.AI
- 13 years in Cybersecurity, 6 years in AI
- Ex-Director of Security Research & Data Science
- Product leader, security researcher, consultant
- Completed 100+ projects in 30+ countries
- Expert in turning research into products

**We are on a mission to increase trust in AI by protecting it from cyber threats**

HITBSECCONF
AMSTERDAM - 2021

# Agenda

- The Need For Secure AI

- The Progress Toward Secure AI

- AI Systems Under Attacks

- How To Attack AI Systems

- *Case Study: Security Of AI Facial Recognition*

- How To Defend AI Systems



*Report: https://adversa.ai/hitb*

# The Need For Secure AI

- The Next Decade Of AI Security

- AI Is The New Attack Vector

- Real Incidents In AI Systems

- The Inception Of Trustworthy AI

# The Next Decade of AI Security

- 1990s – Network security

- 2000s – Endpoint security

- 2010s – Application security

- 2020s – AI security (Software 2.0)

# AI Is The New Attack Vector

## Traditional Software

- **Powered by**
  Fixed program logic

- **Workflow**
  Tasks and commands

- **Interaction**
  Graphical UI with menus and buttons

- **Typical problems**
  Improper validation, access control issues,
  system & security misconfiguration

## AI (Software 2.0)

- **Powered by**
  Flexible ML training

- **Workflow**
  Learning and decisions

- **Interaction**
  Cognitive UI with visual, audio, text commands

- **Typical problems**
  Model manipulation, data exfiltration,
  model & data infection

# Real Incidents In AI Systems

**Confidentiality**

- Personal data extraction from Netflix statistics during ML contest
- Model cloning of pre-released GPT-2 model from Open AI

**Integrity**

- Evasion of Cylance AI-based malware detection
- Poisoning of VirusTotal dataset with fake samples

**Availability**

- Fooling Tesla autopilot could lead to self-driving car crash
- Denial of service in IoT via resource exhaustion attacks

# The Inception Of Trustworthy AI



HITBSECCONF
AMSTERDAM - 2021

# The Progress Toward Secure AI
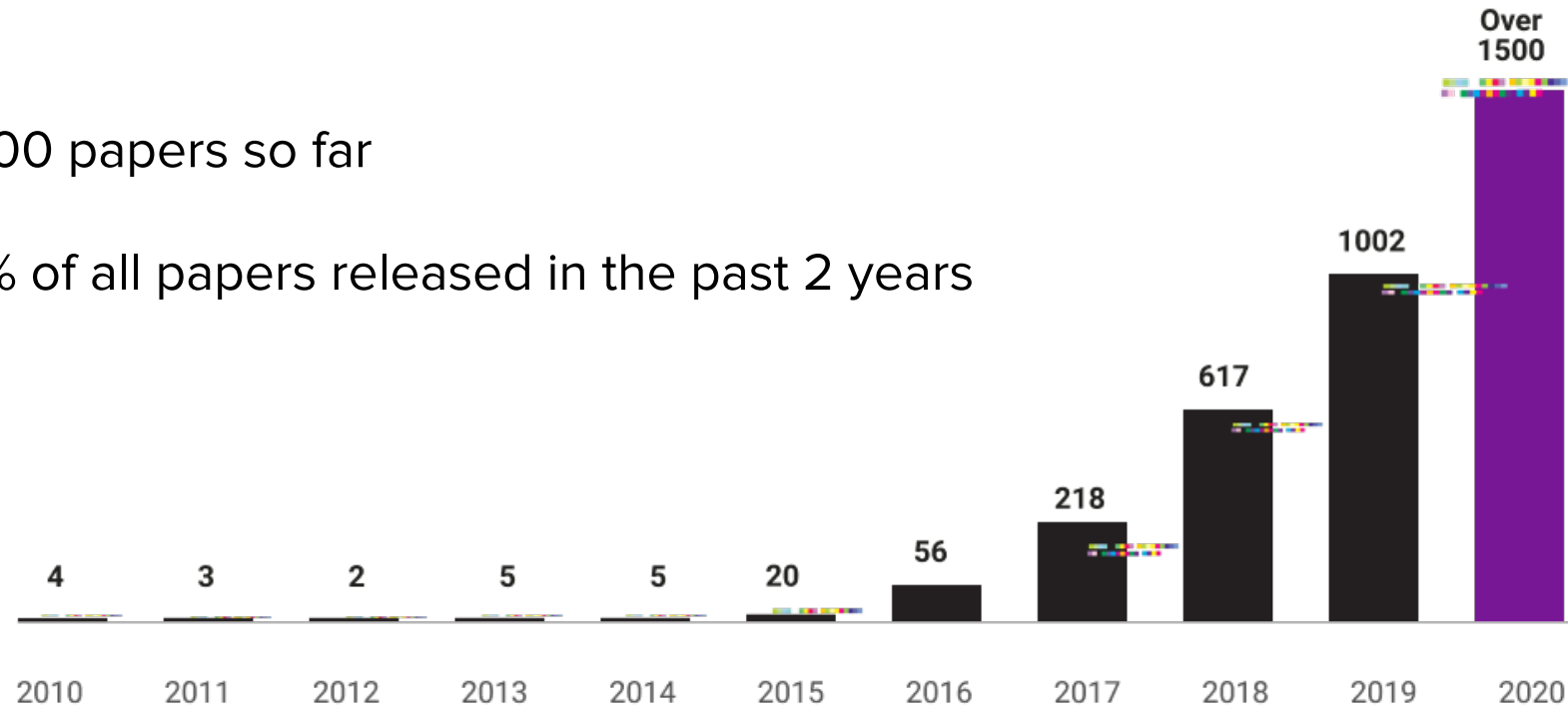
- Adversarial Machine Learning Research

- Most Active Country Contributors

- Key Events In Adversarial ML

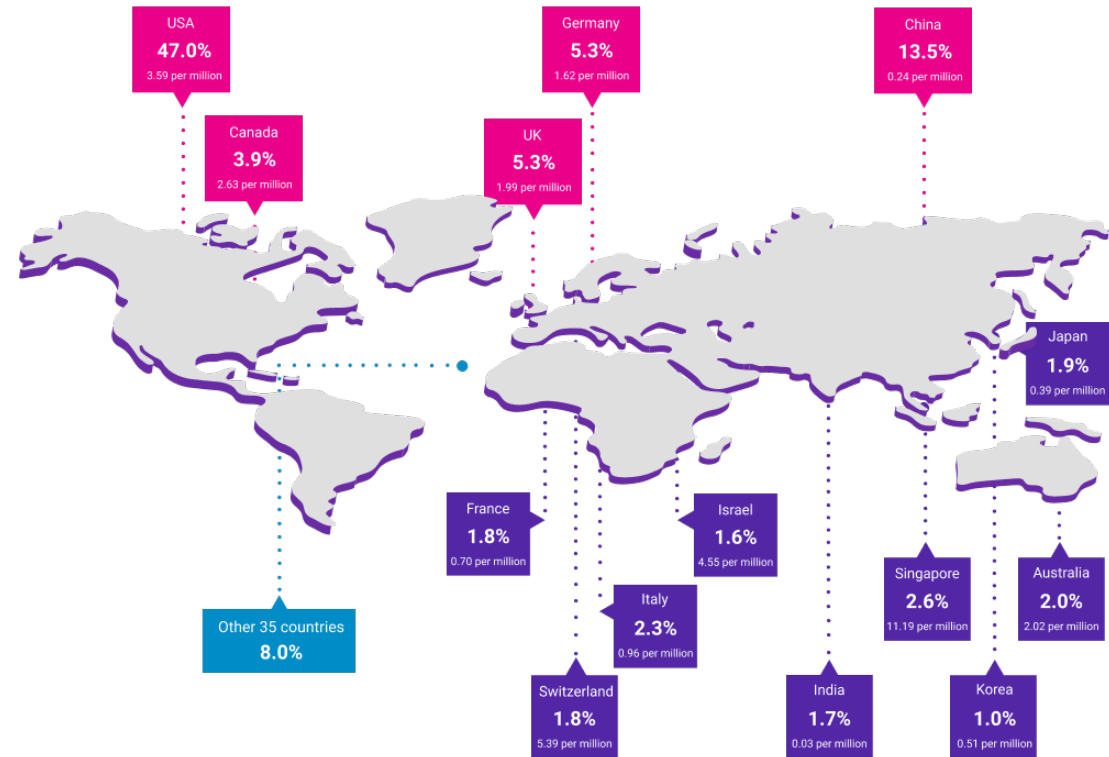# Adversarial Machine Learning Research

- Over 4,000 papers so far

- Over 50% of all papers released in the past 2 years



Bar chart of papers per year:
- 2010: 4
- 2011: 3
- 2012: 2
- 2013: 5
- 2014: 5
- 2015: 20
- 2016: 56
- 2017: 218
- 2018: 617
- 2019: 1002
- 2020: Over 1500

# Most Active Country Contributors

- Contributors from 50 countries

- Over 90% of papers from 14 countries



| USA | Germany | China |
|-----|---------|-------|
| **47.0%** | **5.3%** | **13.5%** |
| 3.59 per million | 1.62 per million | 0.24 per million |

Canada
**3.9%**
2.63 per million

UK
**5.3%**
1.99 per million

Japan
**1.9%**
0.39 per million

France
**1.8%**
0.70 per million

Israel
**1.6%**
4.55 per million

Singapore
**2.6%**
11.19 per million

Australia
**2.0%**
2.02 per million

Italy
**2.3%**
0.96 per million

Other 35 countries
**8.0%**

Switzerland
**1.8%**
5.39 per million

India
**1.7%**
0.03 per million

Korea
**1.0%**
0.51 per million

HITBSECCONF
AMSTERDAM - 2021

# Key Events In Adversarial ML

- **Academia**

  **2004** – early works in adversarial machine learning

  **2014** – reborn interest with attacking deep learning

  **2021** – more than 4000 research papers released

- **Governments**

  **2016** – first national AI document mentioning security in US

  **2019** – first national AI document focused on security in US
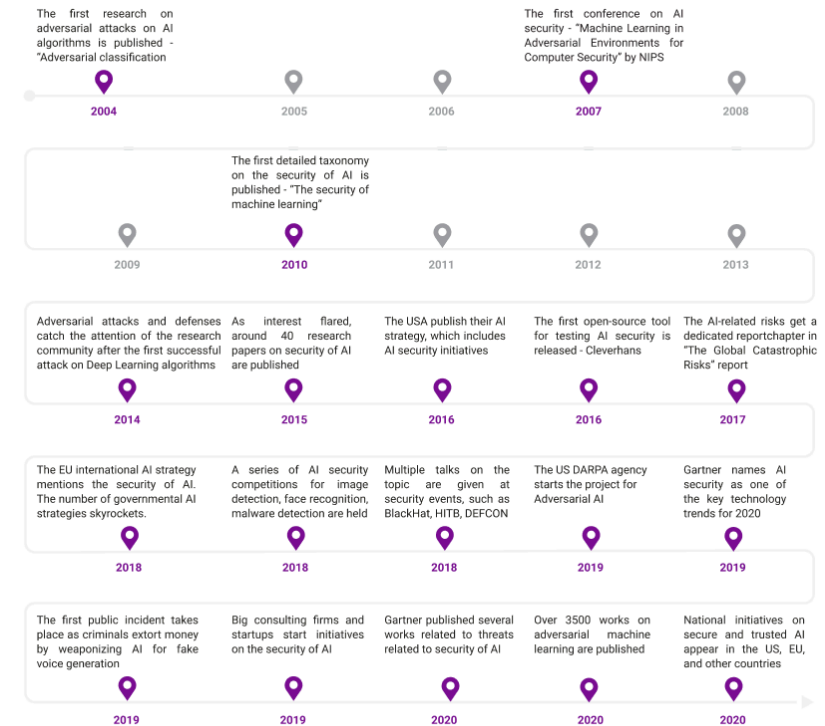
  **2021** – first international regulation for high-risk AI in Europe

- **Industry**

  **2018** – First trusted AI consulting and startups
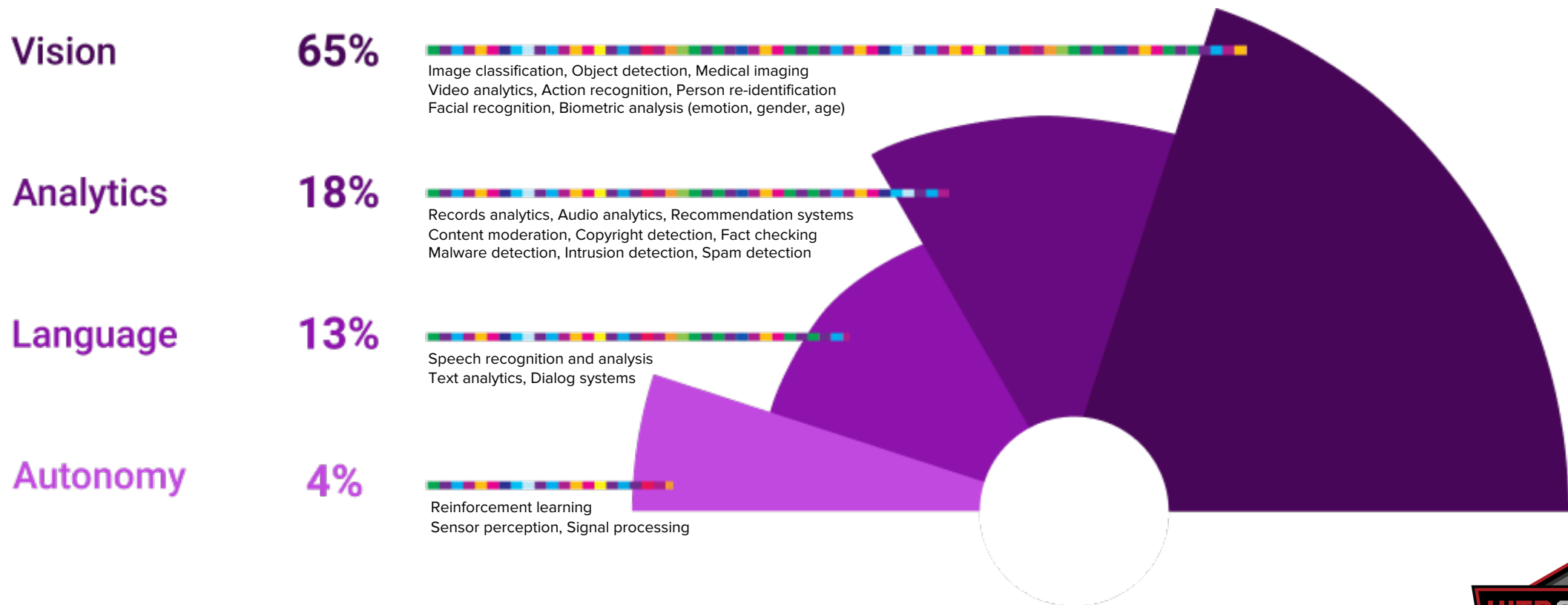
  **2019** – Gartner starts covering the AI security topic

  **2020** – Growth of AI red teams (FB, MS, Nvidia, Open AI, MITRE)



The first research on adversarial attacks on AI algorithms is published - "Adversarial classification
**2004**

2005

2006

The first conference on AI security - "Machine Learning in Adversarial Environments for Computer Security" by NIPS
**2007**

2008

The first detailed taxonomy on the security of AI is published - "The security of machine learning"
**2010**

2009 | 2011 | 2012 | 2013

Adversarial attacks and defenses catch the attention of the research community after the first successful attack on Deep Learning algorithms
**2014**

As interest flared, around 40 research papers on security of AI are published
**2015**

The USA publish their AI strategy, which includes AI security initiatives
**2016**

The first open-source tool for testing AI security is released - Cleverhans
**2016**

The AI-related risks get a dedicated report chapter in "The Global Catastrophic Risks" report
**2017**

The EU international AI strategy mentions the security of AI. The number of governmental AI strategies skyrockets.
**2018**

A series of AI security competitions for image detection, face recognition, malware detection are held
**2018**

Multiple talks on the topic are given at security events, such as BlackHat, HITB, DEFCON
**2018**

The US DARPA agency starts the project for Adversarial AI
**2019**

Gartner names AI security as one of the key technology trends for 2020
**2019**

The first public incident takes place as criminals extort money by weaponizing AI for fake voice generation
**2019**

Big consulting firms and startups start initiatives on the security of AI
**2019**

Gartner published several works related to threats related to security of AI
**2020**

Over 3500 works on adversarial machine learning are published
**2020**

National initiatives on secure and trusted AI appear in the US, EU, and other countries
**2020**
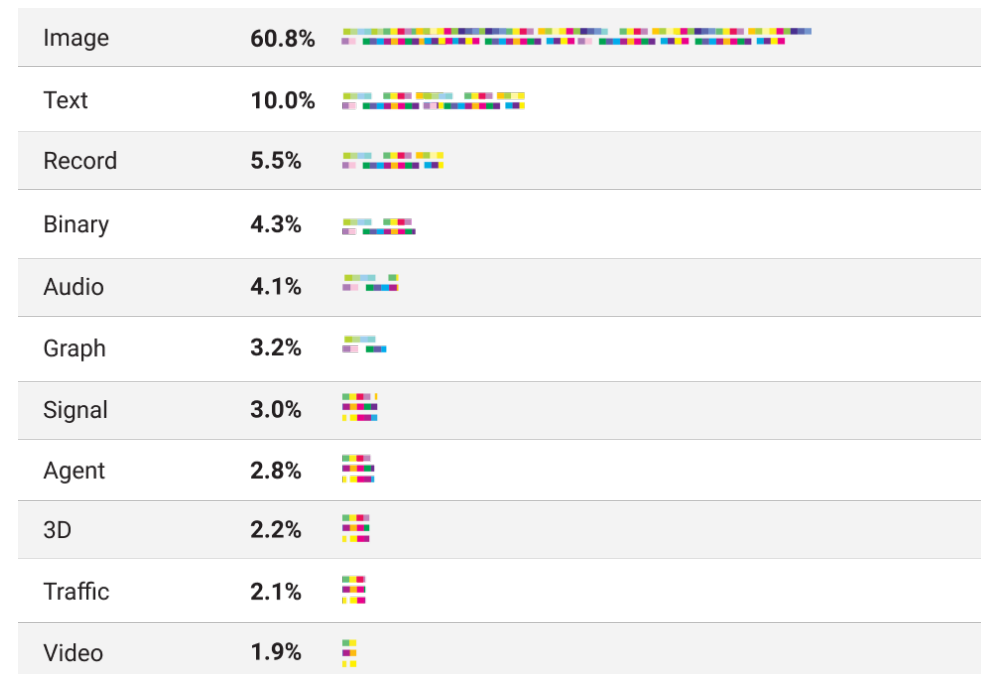
HITBSECCONF
AMSTERDAM - 2021

# AI Systems Under Attacks

- AI Areas Under Attacks

- AI Datasets Under Attacks

- AI Applications Under Attacks

- AI Industries Under Attacks
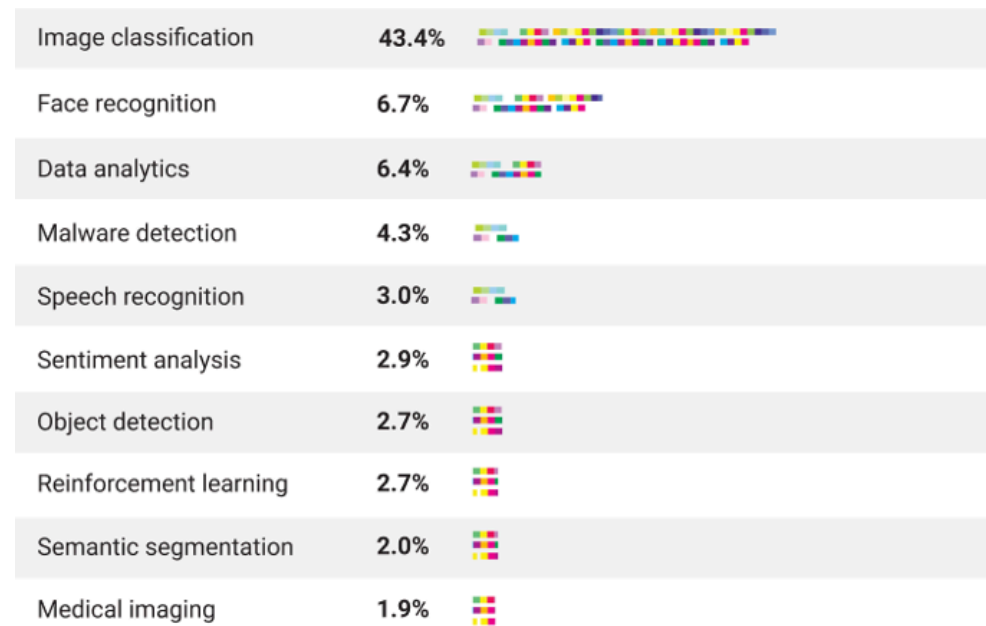
# AI Areas Under Attacks

**Vision**     **65%**

Image classification, Object detection, Medical imaging
Video analytics, Action recognition, Person re-identification
Facial recognition, Biometric analysis (emotion, gender, age)

**Analytics**     **18%**

Records analytics, Audio analytics, Recommendation systems
Content moderation, Copyright detection, Fact checking
Malware detection, Intrusion detection, Spam detection

**Language**     **13%**

Speech recognition and analysis
Text analytics, Dialog systems

**Autonomy**     **4%**

Reinforcement learning
Sensor perception, Signal processing

HITBSECCONF
AMSTERDAM - 2021

# AI Datasets Under Attacks

- The dominance of attacks against AI systems with image processing shouldn't mislead you into thinking that other AI applications are less vulnerable.

- Counterintuitively, AI applications with fewer attacks might be at greater risk because the interest to develop defenses for them is significantly smaller.

| | | |
|---|---|---|
| Image | 60.8% | |
| Text | 10.0% | |
| Record | 5.5% | |
| Binary | 4.3% | |
| Audio | 4.1% | |
| Graph | 3.2% | |
| Signal | 3.0% | |
| Agent | 2.8% | |
| 3D | 2.2% | |
| Traffic | 2.1% | |
| Video | 1.9% | |

# AI Applications Under Attacks

- Over 2000 successful attack cases against over 100 unique AI applications

- AI applications were *never* inherently robust

- **Attacks are transferable across applications**

| | | |
|---|---|---|
| Image classification | 43.4% | |
| Face recognition | 6.7% | |
| Data analytics | 6.4% | |
| Malware detection | 4.3% | |
| Speech recognition | 3.0% | |
| Sentiment analysis | 2.9% | |
| Object detection | 2.7% | |
| Reinforcement learning | 2.7% | |
| Semantic segmentation | 2.0% | |
| Medical imaging | 1.9% | |

HITB SECCONF
AMSTERDAM - 2021

# AI Industries Under Attacks

- Attack Research shows original interest of researchers to attack AI in a given industry

- Transferable Risk reveals the real threat landscape based on our risk correlation

- **Attacks are transferable across industries**

| Position | Industry | Attack Research | Transferable Risk |
|---|---|---|---|
| 1 | Internet | 23% | 97% |
| 2 | Cybersecurity | 17% | 41% |
| 3 | Biometrics | 16% | 67% |
| 4 | Automotive | 13% | 79% |
| 5 | Healthcare | 9% | 87% |
| 6 | Industrial | 5% | 74% |
| 7 | Smart Home | 5% | 89% |
| 8 | Retail | 4% | 86% |
| 9 | Finance | 4% | 95% |
| 10 | Surveillance | 3% | 77% |
| 11 | Robotics | 1% | 61% |

# How To Attack AI Systems

- Focus Of Adversarial ML Research

- Categories Of AI Attacks

- Adversa Top 10 Attacks

# Focus Of Adversarial ML Research

- **Attacks (49%)**

Researchers invent new ways of hacking AI models, AI infrastructure, side channels, and software bugs.

- **Defenses (47%)**

Researchers suggest defenses like analysis of data inputs/outputs, algorithm modifications, or retraining.

- **Surveys (3%)**

Researchers review security threats for AI systems and investigate existing attacks and defenses.

- **Tools (1%)**

Researchers develop tools for performing vulnerability testing and verification of AI models.

# Categories Of AI Attacks

**Manipulation** **82.6%**

Evasion attacks,
Reprogramming

**Infection** **9.5%**

Poisoning,
Backdoors, Trojans

**Exfiltration** **7.8%**

Inference attacks,
Model stealing

# Adversa Top 10 Attacks

**1** **Evasion attack** bypasses normal decisions by AI systems in favor of attacker-controlled behavior by crafting malicious data inputs called adversarial examples

**2** **Poisoning attack** reduces the quality of AI decisions while making AI systems unreliable or unusable by injecting malicious data into a dataset used for AI training

**3** **Membership inference attack** discloses whether specific data sample was a part of a dataset used for AI training

**4** **Backdoor attack** invokes hidden behavior of AI systems after poisoning them with secret triggers while keeping AI models work as intended in normal conditions

**5** **Model extraction attack** exposes algorithm's internal details by making malicious queries to AI systems

**5** **Model extraction attack** exposes algorithm's internal details by making malicious queries to AI systems

**6** **Attribute inference attack** reveals secret data details by exploiting public information received from AI systems

**7** **Trojan attack** enables attacker-controlled behavior of AI systems after malicious modification or distribution of AI models that work as expected in normal conditions

**8** **Model inversion attack** reveals secret data inputs based on public outputs by maliciously querying AI systems

**9** **Anti-watermarking attack** bypasses protection controls used by AI systems for copyright or authenticity checks
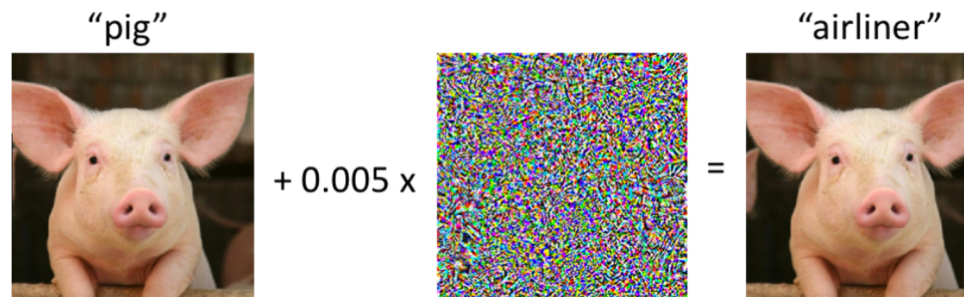
**10** **Reprogramming attack** allows threat actors to repurpose AI models and make them execute unexpected tasks

# Case Study: Security Of AI Facial Recognition

- Intro To Evasion Attacks In Computer Vision

- Problem Of Securing AI Facial Recognition

- AI Security Testing Challenges

- AI Security Testing Results

# Intro To Evasion Attacks In Computer Vision

- Discover the most important pixels by interacting with a model

- Craft a malicious image with modified pixels to fool a model

- Model makes a wrong prediction controlled by an attacker
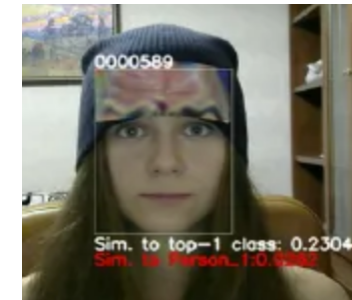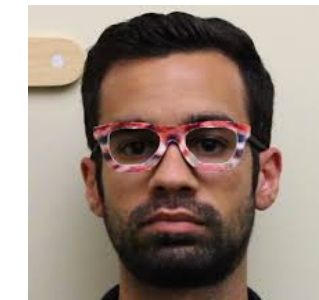
- Changes are imperceptible for system owners



"pig"  + 0.005 x  =  "airliner"

# Case Study Profile

- A smart home solution provider expressed security concerns

- They needed to implement a secure facial recognition solution

- The goal was to find the most reliable hardware/software vendor

# Problem Of Securing AI Facial Recognition

- Over 4000 research papers on adversarial machine learning

- Over 100 research papers on facial recognition security

- Countless combos of different conditions (attacks, models, environments)
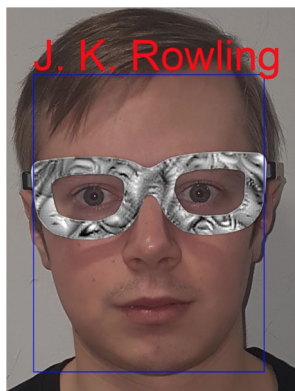
- No clear understanding of real-world risks beyond research

# AI Security Testing Challenges

- **Goals**: Confidence reduction, untargeted/targeted misclassification, etc.

- **Forms**: Glasses, lenses, mask, hat, moustache, band-aid, etc.

- **Knowledge**: Black-box, gray-box, white-box

- **Constraints**: Time, computation resources, data

- **Conditions**: Printing quality, inconsistent colors, position, angle, size, etc.

- **Algorithms**: FGSM, BIM, PGD, EOT, etc.

# AI Security Testing Results

- Physical attacks have successfully fooled facial recognition systems

- Glasses and bandanas have achieved the best misclassification rate

- Further attack optimization was possible but was out of scope

# How To Defend AI Systems

- Lifecycle for Secure AI Development

- Your Next Steps

# Lifecycle for Secure AI Development

## 1. Identify

Understand current AI security posture with *asset management*, *threat modeling*, *risk assessment*.

## 2. Protect

Implement protective controls such as *security awareness*, *system hardening*, *secure AI development*.

## 3. Detect

Defend against active adversaries with *security monitoring, threat detection*, and *AI red teaming*.

## 4. Respond

Prepare for AI security incidents by developing practices of *AI incident forensics and response*.

*Get detailed lifecycle: https://adversa.ai/hitb*

HITBSECCONF
AMSTERDAM - 2021
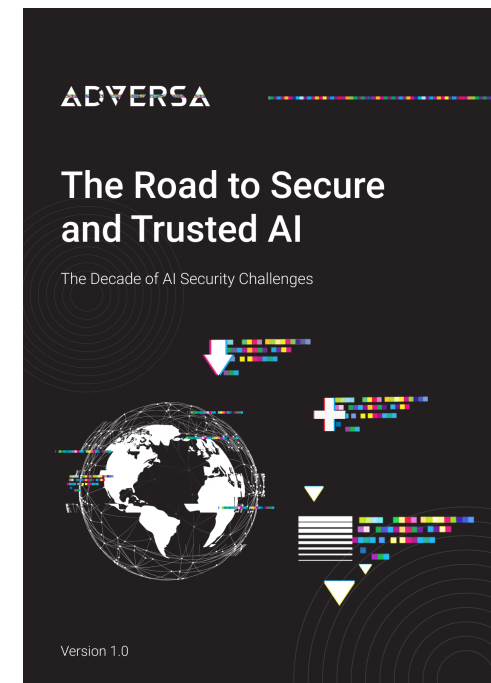
# Your Next Steps

1. **Awareness**
   - Educate stakeholders in AI and security teams about AI threats
   - Study relevant policies and practices for secure and trusted AI

2. **Assessment**
   - Perform initial threat modeling to understand AI security risks
   - Conduct initial security testing for mission-critical AI systems

3. **Assurance**
   - Understand and respond to AI security findings
   - Integrate security activities into AI development lifecycle

ADVERSA

**The Road to Secure and Trusted AI**

The Decade of AI Security Challenges

Version 1.0

*Report: https://adversa.ai/hitb*

# Thank You

Follow us! adversa.ai/stay-updated

Collaborate with us! info@adversa.ai