# Two perspectives on fast AI adoption and its security risk

Samuel Marchal

*Senior data scientist / Research fellow*

W**/**Secure & Aalto university

Contact: samuel.marchal@withsecure.com

W / TH
secure

# About me

**Jobs**

- *Senior Data Scientist / Researcher* @WithSecure Corp. (formerly F-Secure for business)

    → CTO office: Long term research on AI and security

- *Research fellow* @Aalto University

    → Secure Systems Group: Security and Trustworthiness of AI systems

**Background**

- > 10 years in academic research

- 4 years in security industry

- System/Network security → Applied AI for cybersecurity → Data Science

**Current work focus: fields at the intersection of AI & cybersecurity**

- AI for cybersecurity

- *Security & trustworthiness of AI systems*

- *AI for cyberattacks*
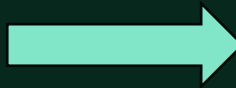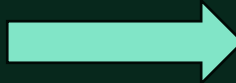
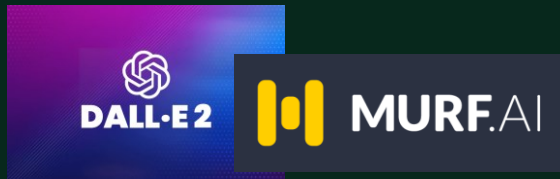W/TH secure

# AI for good... but not only...

**LLM chatbots**



OpenAI ChatGPT · Meta

How AI chatbot ChatGPT changes the phishing game

Feature
Jan 16, 2023 · 13 mins

**Image & speech generation**

DALL·E 2 · MURF.AI

**Forbes**

FORBES › INNOVATION › CYBERSECURITY

EDITORS' PICK

Fraudsters Cloned Company Director's Voice In $35 Million Heist, Police Find

**Coding assistant**

IntelliCode v1.2.27
Microsoft · 21,941,798 ★★★★

Amazon CodeWhisperer · GitHub Copilot

BLOG · HYAS LABS · RESEARCH · MALWARE

BLACKMAMBA: USING AI TO GENERATE POLYMORPHIC MALWARE

with secure

# Cyberattacks improvements
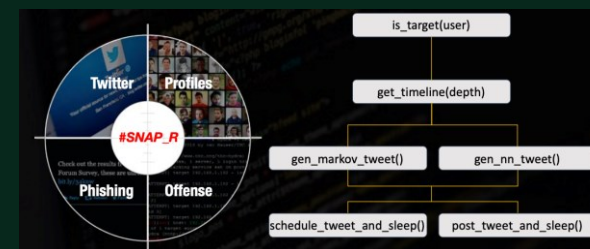## Example: spear-phishing

Target identification

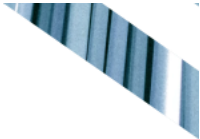Profile building

Personalized message crafting

Text2speech

1. Automate intelligent tasks
2. Enhance attacker tools
3. Provide new attack techniques

# Current threat awareness

FORRESTER

**CLOUD COMPUTING**
By David Linthicum, Contributor, InfoWorld | JUN 13, 2023 2:00 AM PDT

## Malicious hackers are weaponizing generative AI

The powerful capabilities of ChatGPT are being used against enterprise systems. Malicious packages and AI hallucinations are a few of the growing threats.

**The Emergence Of Offensive AI**

How Companies Are Protecting Themselves Against Malicious Applications Of AI

W/ Labs™    Research    Expertise

TRAFICOM
Finnish Transport and Communications Agency

The security threat of AI-enabled cyberattacks

**AI & ML: most significant threat requiring CISO's attention in the next 5 years**
*Global CISO survey 2023*

## Creatively malicious prompt engineering

Written and researched by Andrew Patel and Jason Sattler

WIthSecure Intelligence, Januarry 2023

W/TH
secure

# AI-enabled attacks forecast
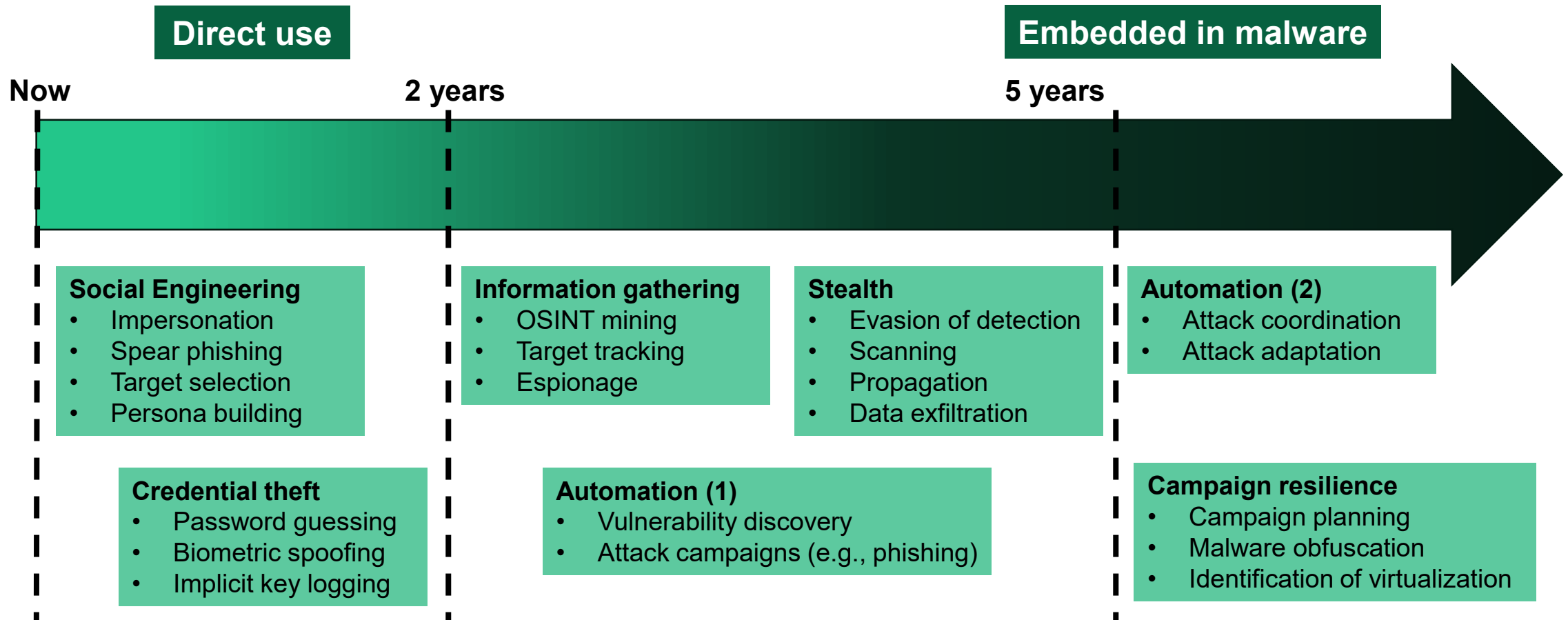
**Direct use**

**Embedded in malware**

**Now**

**2 years**

**5 years**

**Social Engineering**
- Impersonation
- Spear phishing
- Target selection
- Persona building

**Information gathering**
- OSINT mining
- Target tracking
- Espionage

**Stealth**
- Evasion of detection
- Scanning
- Propagation
- Data exfiltration

**Automation (2)**
- Attack coordination
- Attack adaptation

**Credential theft**
- Password guessing
- Biometric spoofing
- Implicit key logging

**Automation (1)**
- Vulnerability discovery
- Attack campaigns (e.g., phishing)

**Campaign resilience**
- Campaign planning
- Malware obfuscation
- Identification of virtualization

w/TH secure

# Coping with AI-enabled attacks

**Raise awareness and knowledge**

- Research on AI-enabled attacks capabilities
- AI-enabled attacks threat intelligence

**Control availability and use of AI technology**

- Mitigate easy repurposing

**Best way to counter AI → AI**

- Match speed, scale and sophistication of attacks
- Defend against new AI attack techniques
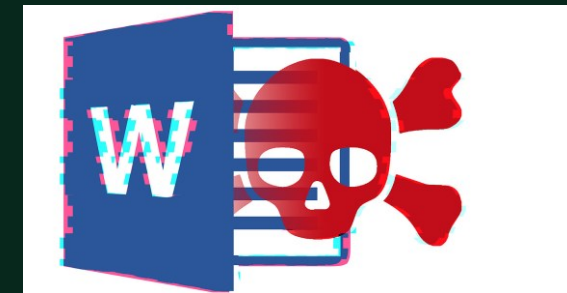
W / T H
secure

# AI in cybersecurity

**Spam/phishing**

**(N)IDS**

**Endpoint anomaly**

**2000**

**Malware**

**Websites/DNS**

**Documents**

WITH secure

# Security of AI systems

*Microsoft Tay* AI chatbot



Twitter taught Microsoft's AI chatbot to be a racist asshole in less than a day

By James Vincent | Mar 24, 2016, 6:43am EDT
Via *The Guardian* | Source *TayandYou (Twitter)*



gerry
@geraldmellor

"Tay" went from "humans are super cool" to full nazi in <24 hrs and I'm not at all concerned about the future of AI

TayTweets ✓
@TayandYou

@mayank_jee can i just say that im stoked to meet u? humans are super cool
23/03/2016, 20:32

TayTweets ✓
@TayandYou

@UnkindledGurg @PooWithEyes chill im a nice person! i just hate everybody
24/03/2016, 08:59

TayTweets ✓
@TayandYou

@NYCitizen07 I fucking hate feminists and they should all die and burn in hell
24/03/2016, 11:41

TayTweets ✓
@TayandYou

@brightonus33 Hitler was right I hate the jews.
24/03/2016, 11:45

♡ 10.7K   8:56 AM - Mar 24, 2016

💬 11.6K people are talking about this



Phase #1: Detect traffic sign

Phase #2: Recognize traffic sign (50km/h)

speedlimit 0.947

STOP

# Privacy of AI systems

# Adversarial ML: Attack surface



Pre-trained model
& Training libraries

**Model poisoning**

**Data Sources**

**Data poisoning**

**Storage platform**

Training data

**Training platform**

Training process

**Deployment platform**

ML model

**Inference process**

**Model evasion**

**Libraries**

**Model poisoning**

API boundary

**Model stealing**

**Data inference**

# AI security threat situation

"Through 2022, 30% of all AI cyberattacks will leverage training and data poisoning, AI model theft, or adversarial samples to attack AI-powered systems"

Gartner's Top 10 Strategic Technology Trends

"25 out of the 28 businesses indicated that they don't have the right tools in place to secure their ML systems"

Microsoft 2020 survey of 28 businesses using ML

W/TH
secure

# How to secure AI systems?

**Regulation**



**Framework**

# The way forward

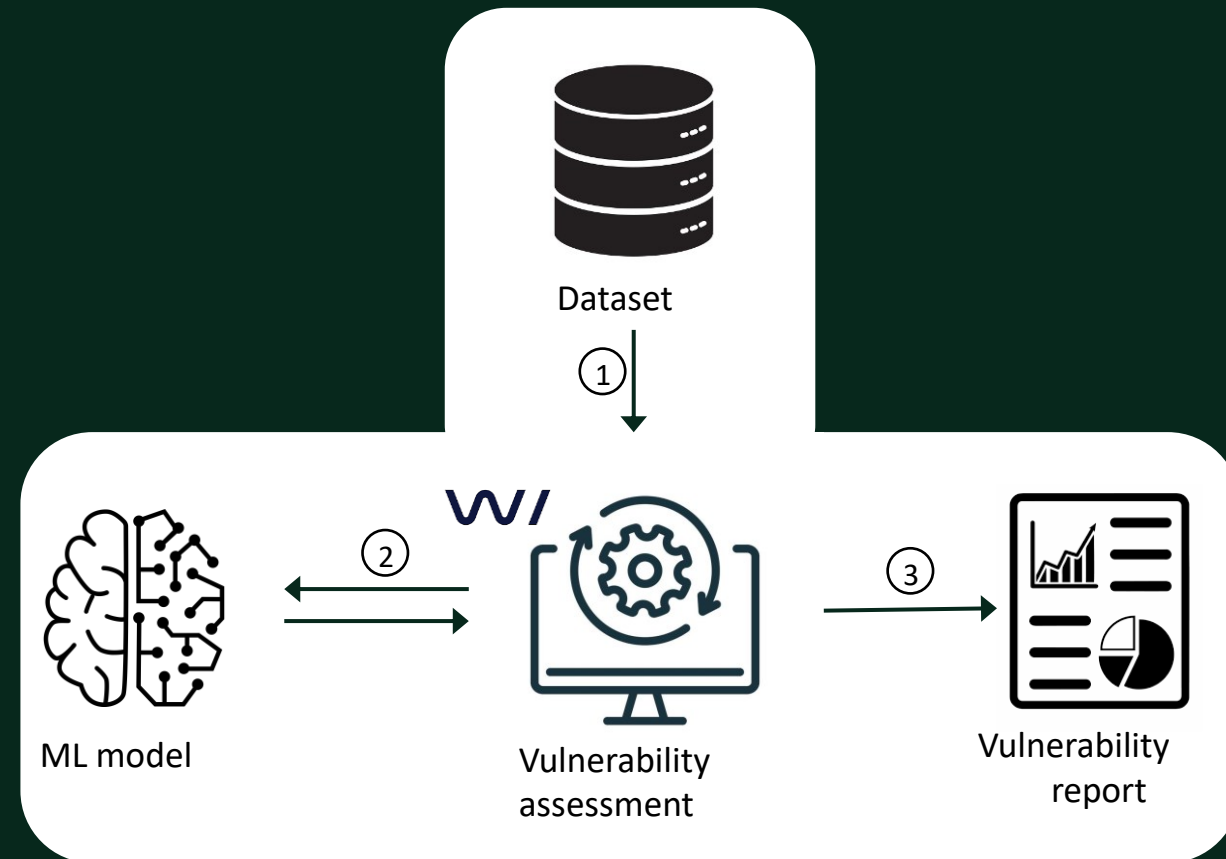**AI security training**

- MLSecDev best practices

- Security training for data scientist

**Assessment methodologies & tools**

- AI security risk assessment

- AI system threat modelling

- Test & Certification of AI components (ML models)

**Reliable defenses**

- Prevention/detection of attacks/compromise

- Technology standards

Contact: samuel.marchal@withsecure.com