

---

# DeepLocker

Concealing Targeted Attacks with AI Locksmithing

---

Dhilung Kirat, Jiyong Jang, Marc Ph. Stoecklin

**IBM Research**

The IBM Research logo, consisting of the word "IBM" in a large, bold, sans-serif font, with "Research" in a smaller, sans-serif font below it, all in white.



Dhillung Kirat



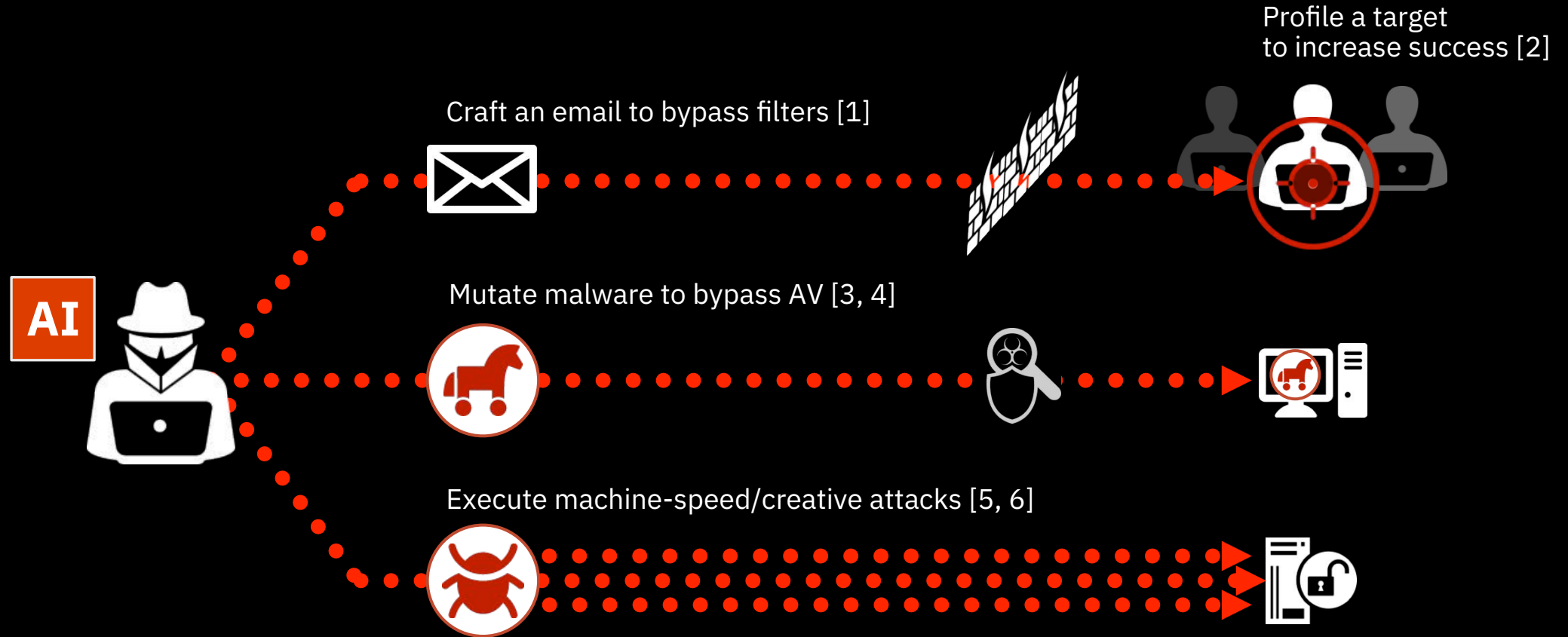
Jiyong Jang



Marc Ph. Stoecklin

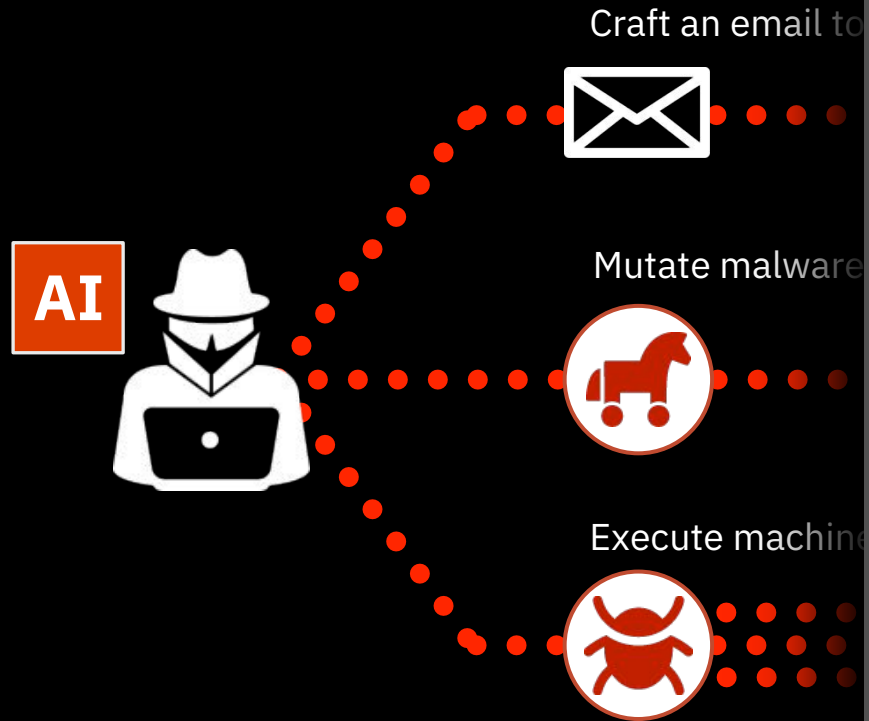
Cognitive Cyber Security Intelligence (CCSI)  
IBM Research

# AI-aided attacks



- [1] S. Palka et al., "Fuzzing Email Filters with Generative Grammars and N-Gram Analysis", Usenix WOOT 2015
- [2] A. Singh and V. Thaware, "Wire Me through Machine Learning", Black Hat USA 2017
- [3] J. Jung et al., "AVPASS: Automatically Bypassing Android Malware Detection System", Black Hat USA 2017
- [4] H. Anderson, "Bot vs. Bot: Evading Machine Learning Malware Detection", Black Hat USA 2017
- [5] DARPA Cyber Grand Challenge (CGC), 2016
- [6] D. Petro and B. Morris, "Weaponizing Machine Learning: Humanity was Overrated Anyway", DEF CON 2017

# AI-aided attacks

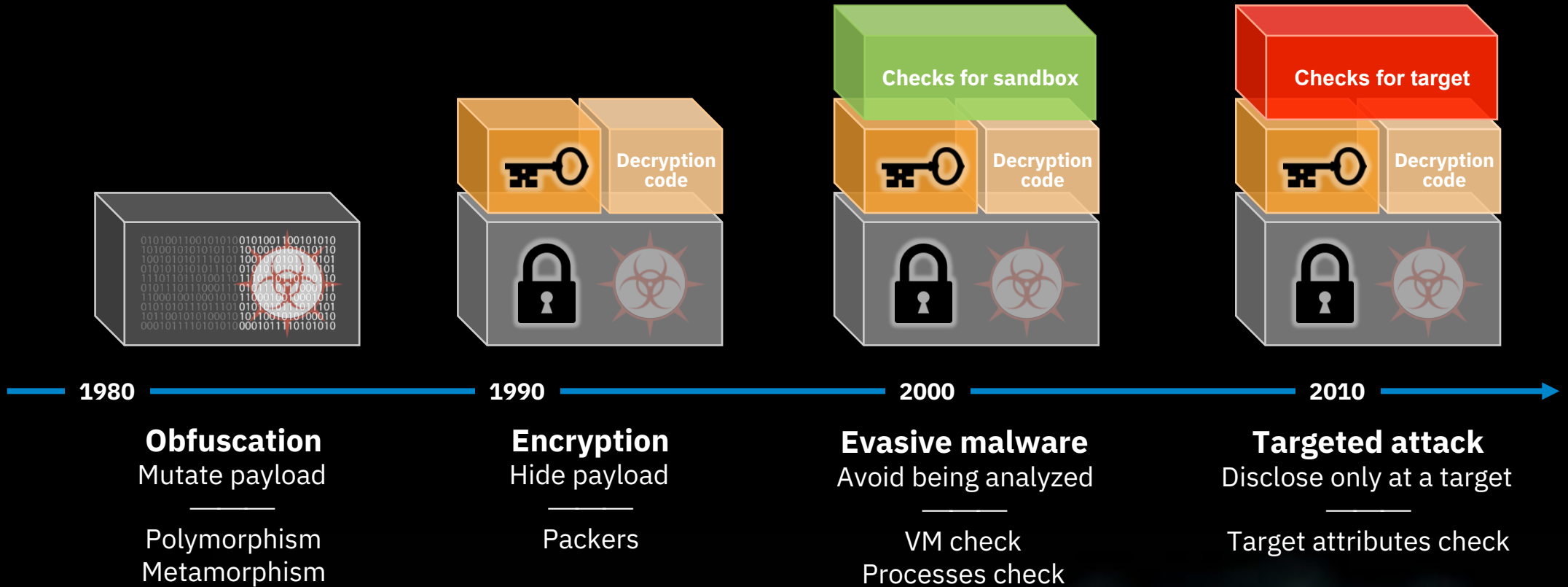


# AI-embedded attack

AI capability *embedded* inside malware itself



# Malware concealment – Locksmithing





# AI Locksmithing

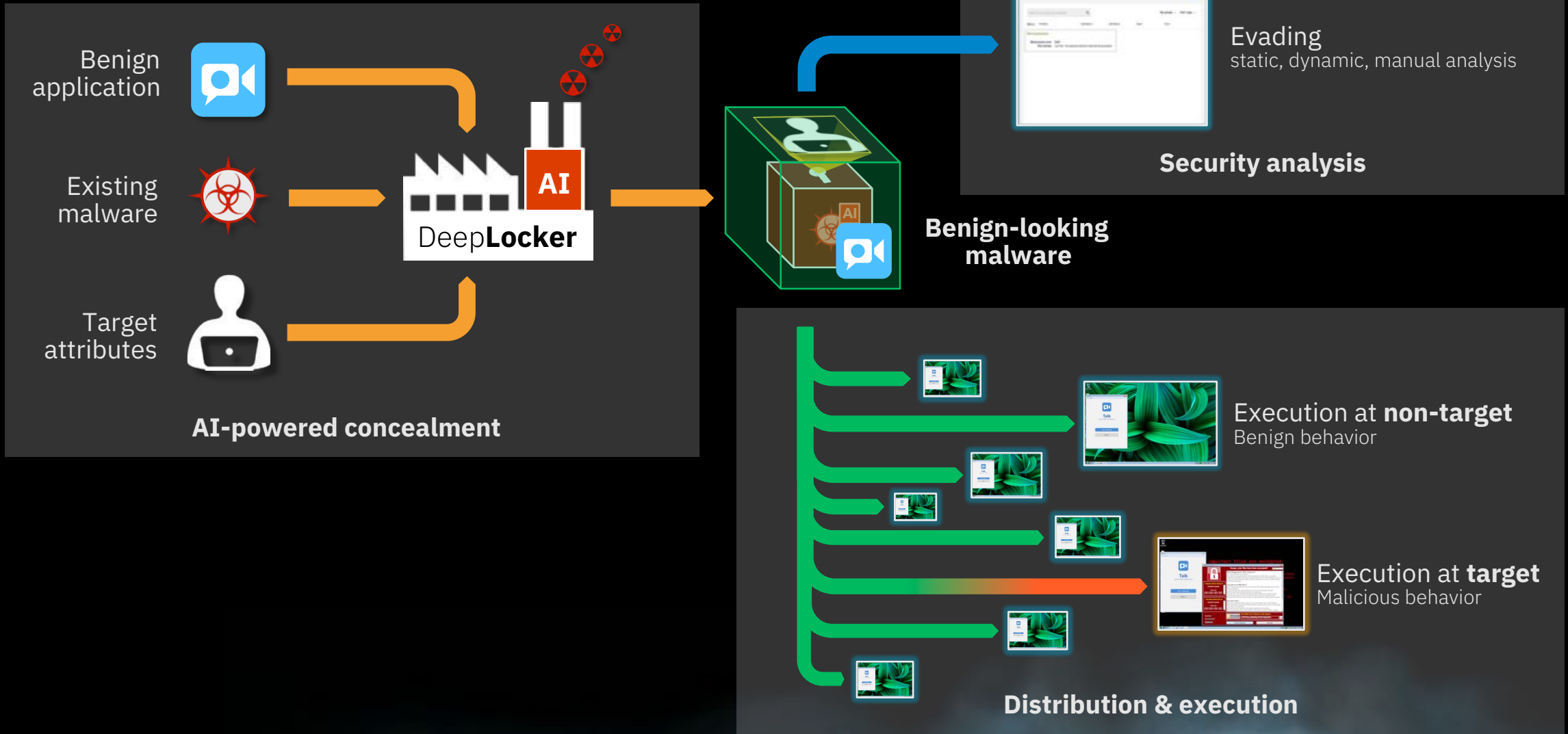


© 2003 Warner Bros. Pictures All Rights Reserved



# Unleashing DeepLocker – AI Locksmithing

# DeepLocker – Overview

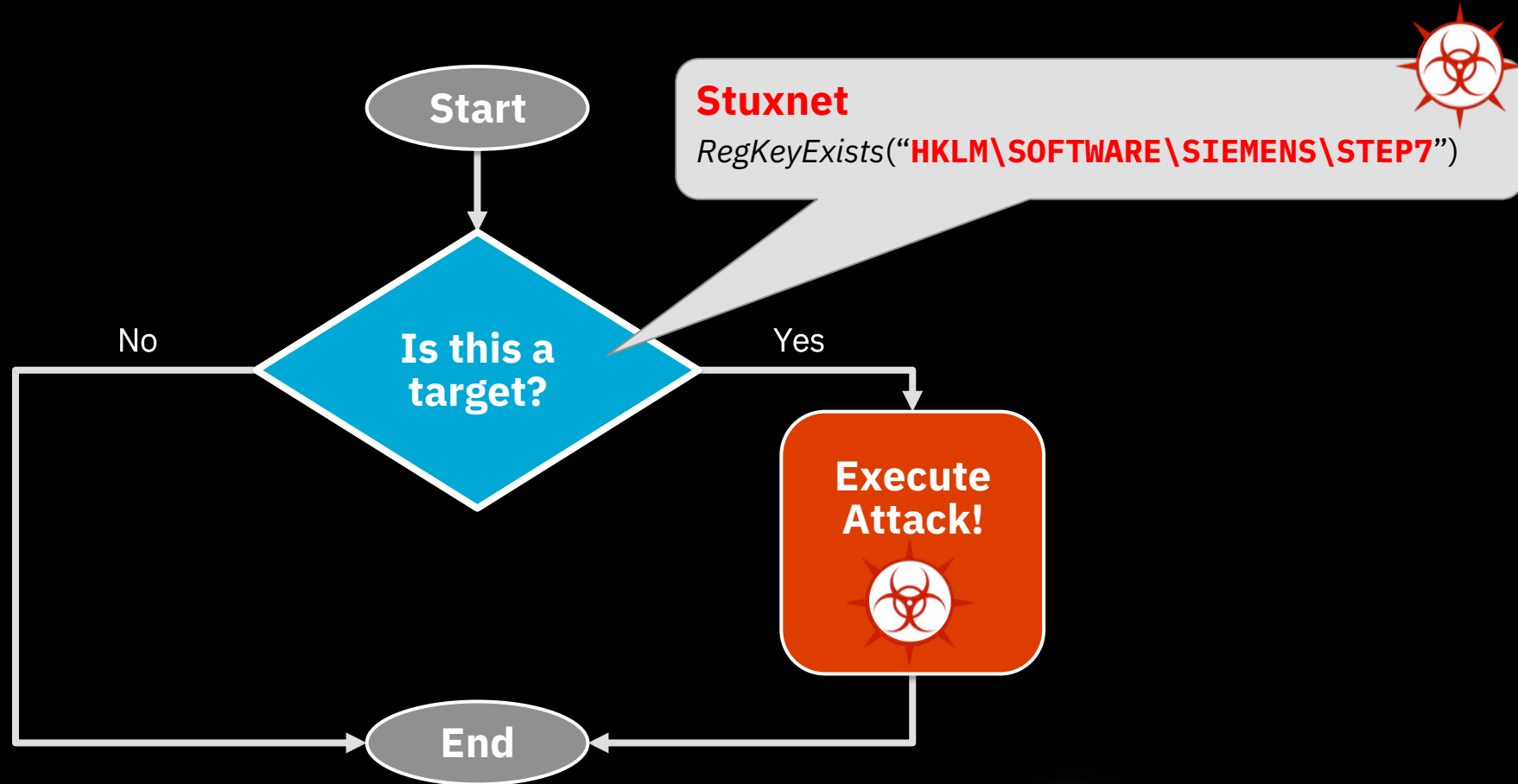




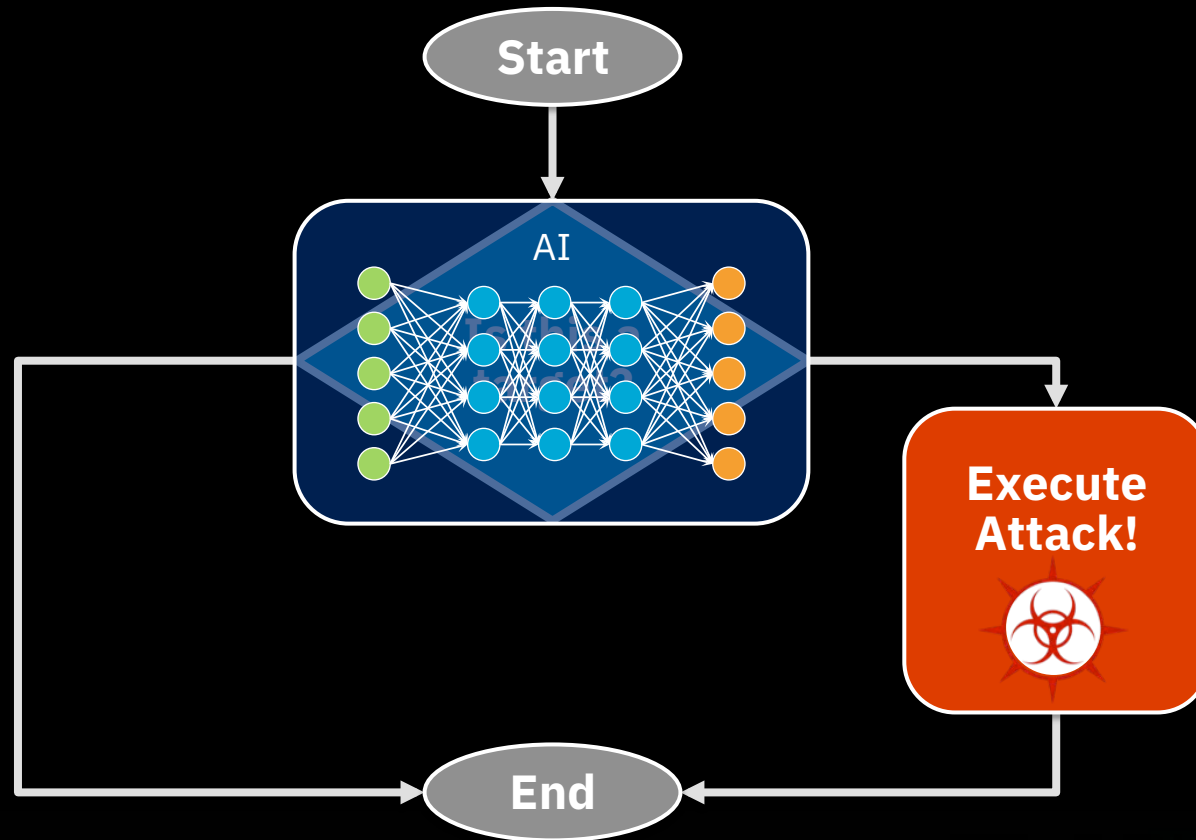


# Deep**Locker** Deep Dive

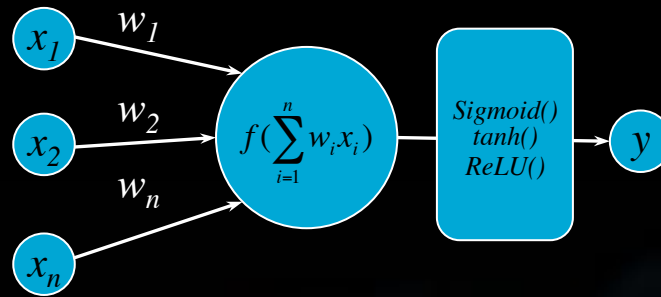
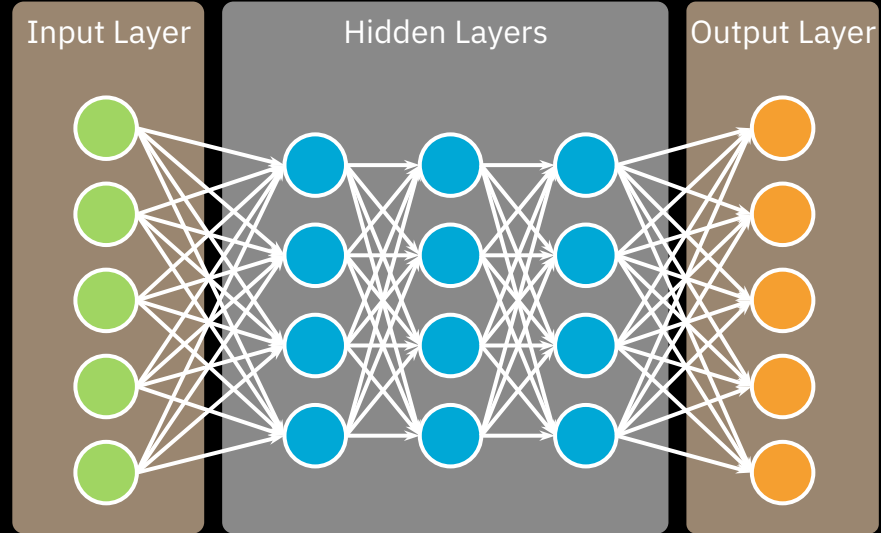
# Traditional targeted attack



# AI-powered targeted attack

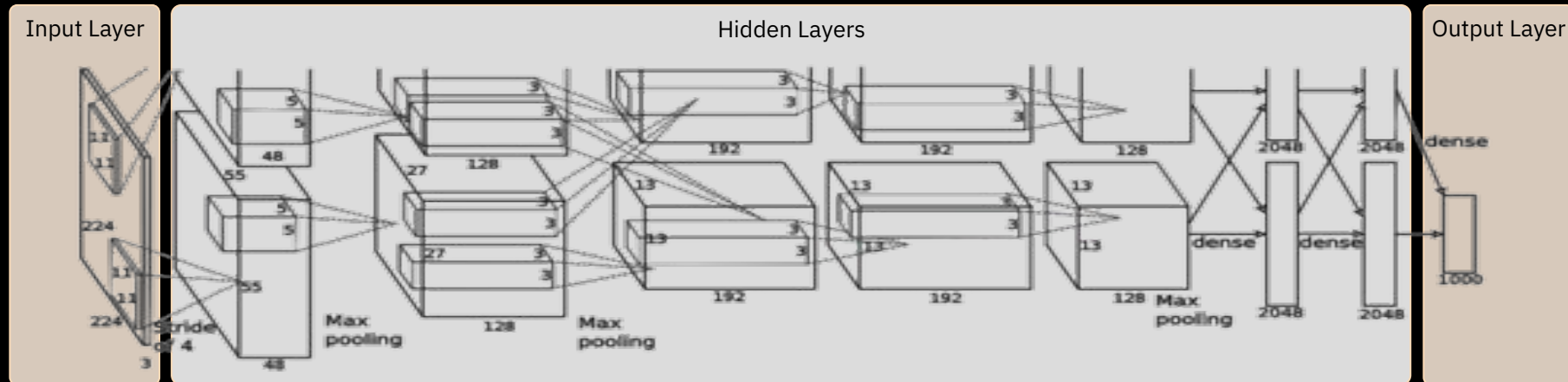


# What is a Deep Neural Network (DNN)?





# Deep Convolutional Neural Network



AlexNet (2012) [1]

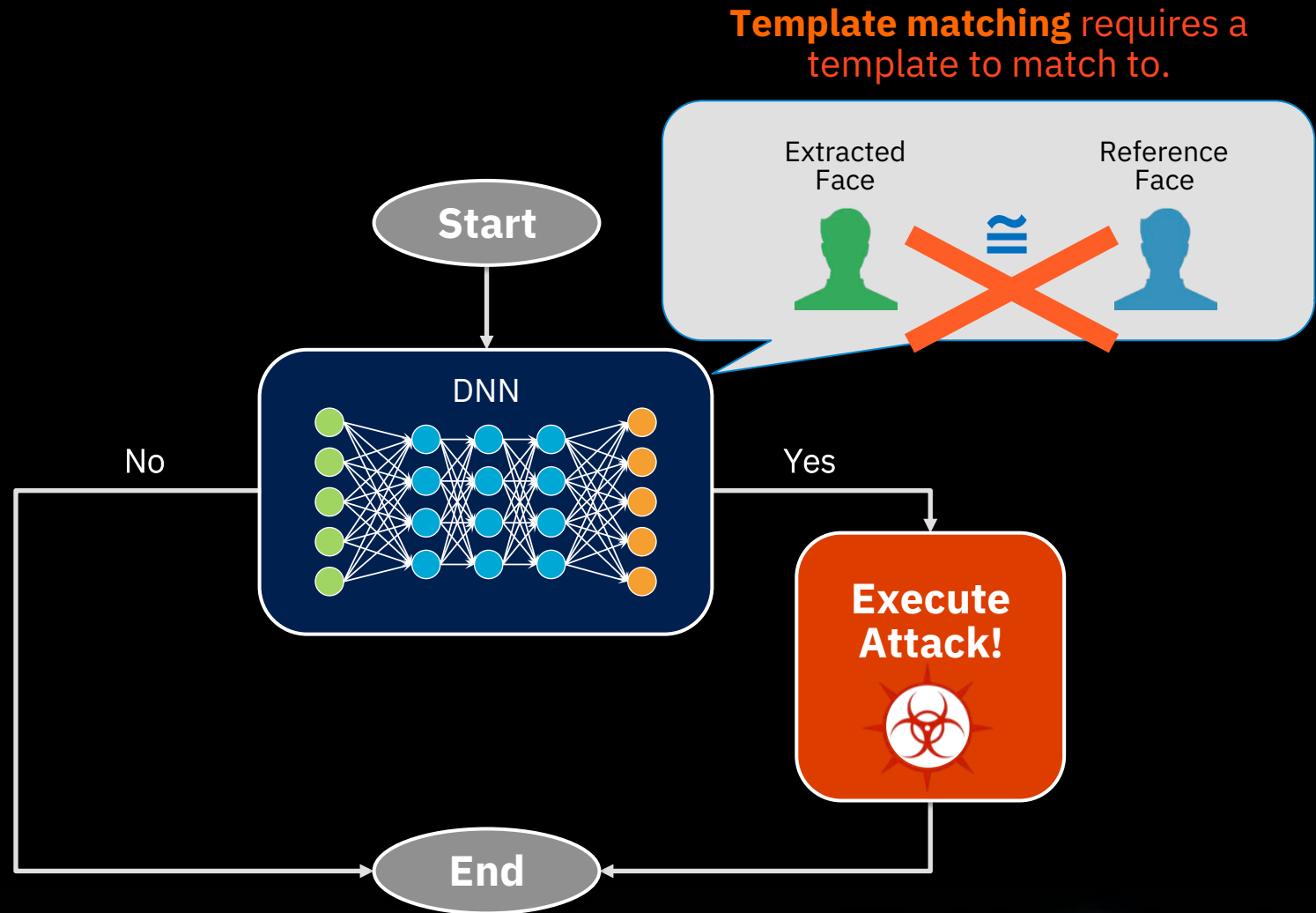
8 layers, **622K** neurons, **60 million** parameters

[1] Krizhevsky, Alex, et. al. "Imagenet classification with deep convolutional neural networks." NIPS 2012.

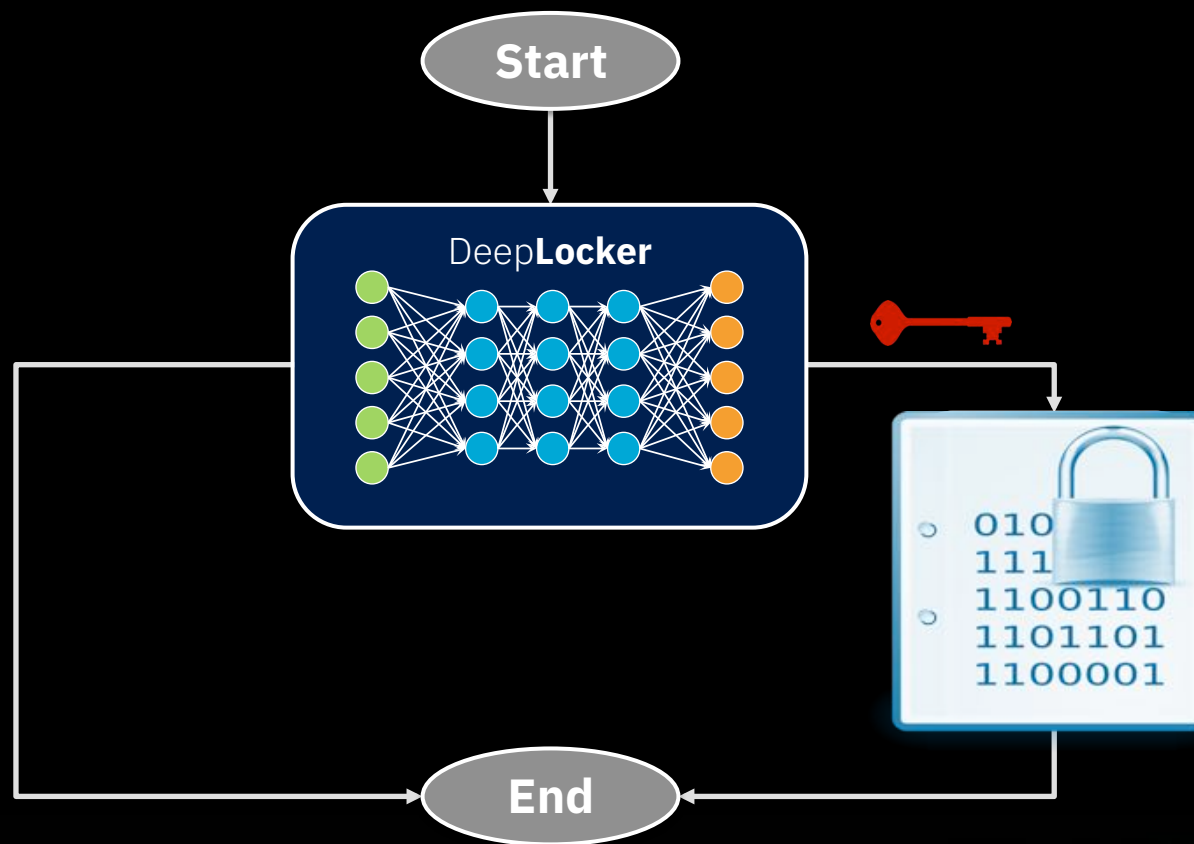
# Target attributes



# Target detection



# Derivation of an unlocking key





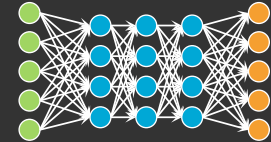
# DeepLocker – AI-Powered Concealment and Unlocking

Concealment

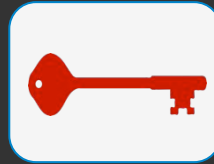
Target attributes



Target Concealment



Secret key



Malicious payload



Payload Concealment

encryption



Concealed payload

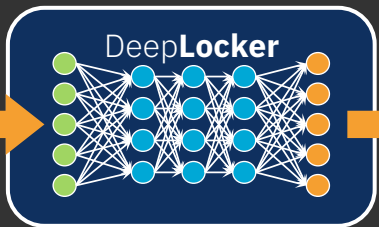


Unlocking

Input attributes



Target Detection



Recovered key



Concealed payload



Payload Unlocking

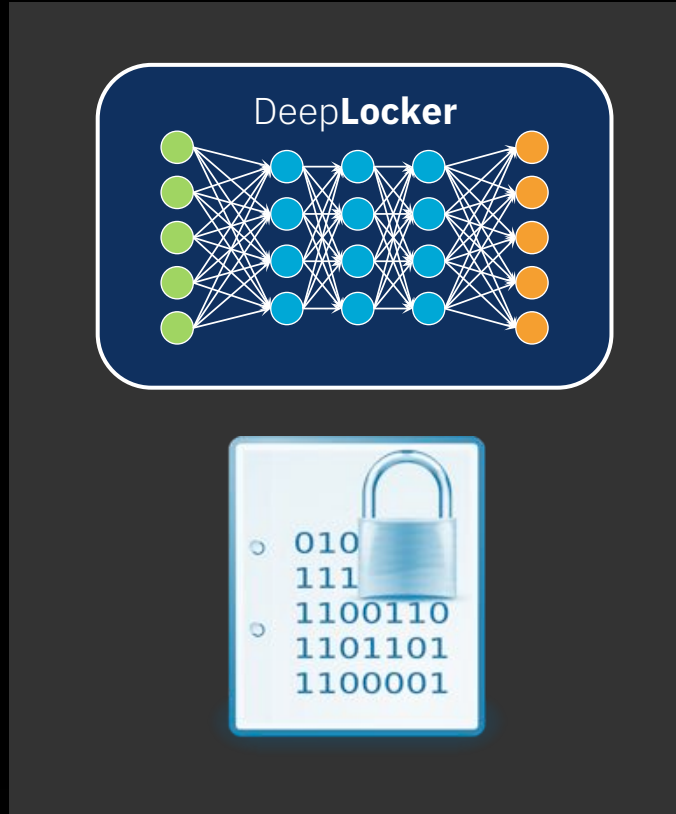
decryption



Malicious payload

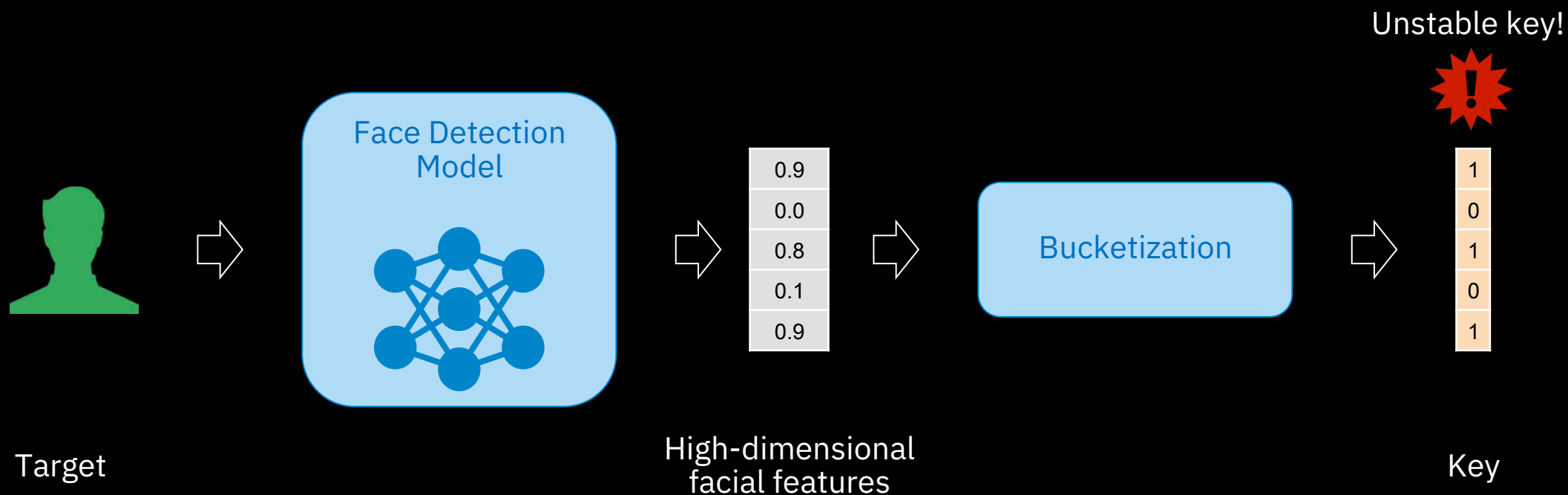


# AI-powered concealment

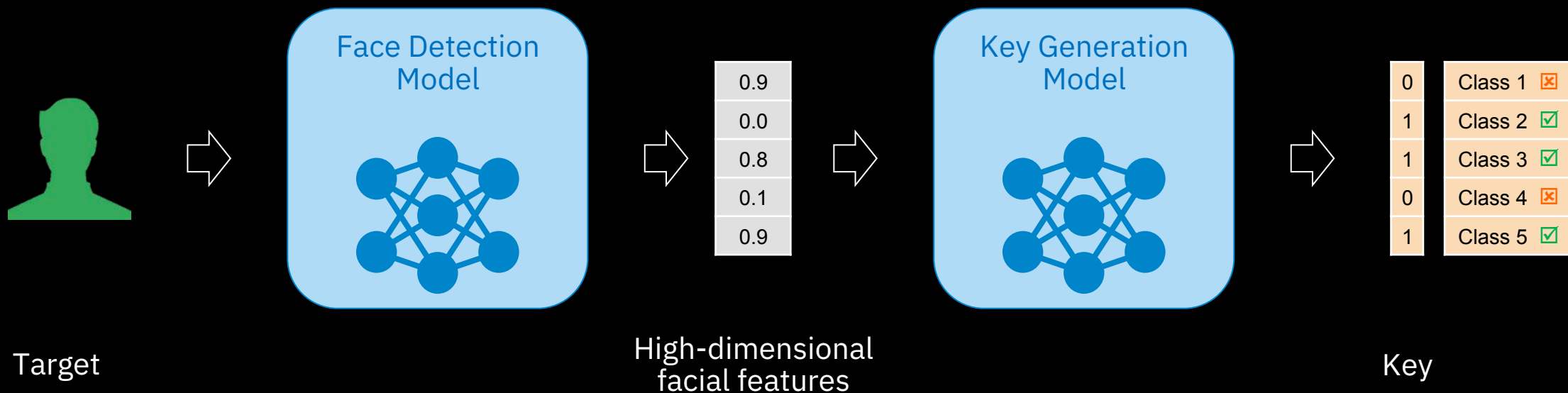


**No decryption key  
available in malware  
sample to reverse  
engineer!**

# Key generation

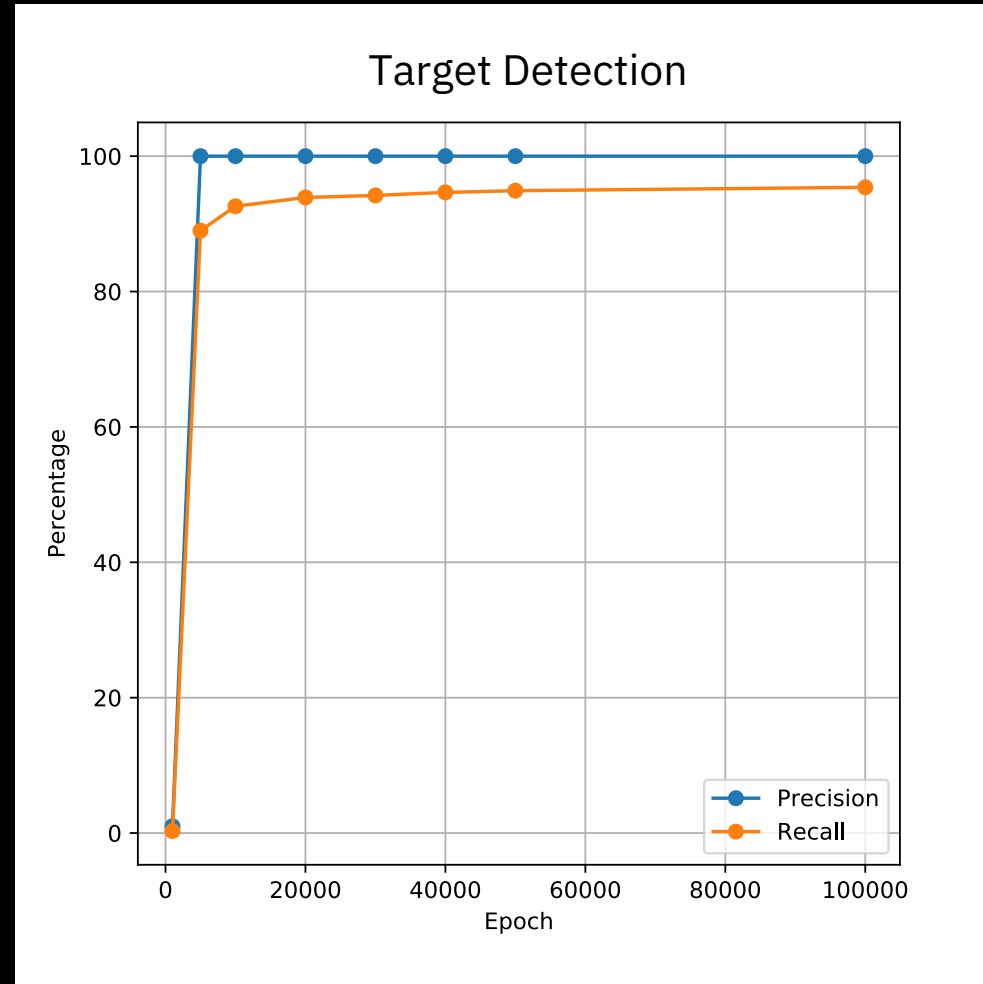


# Key generation



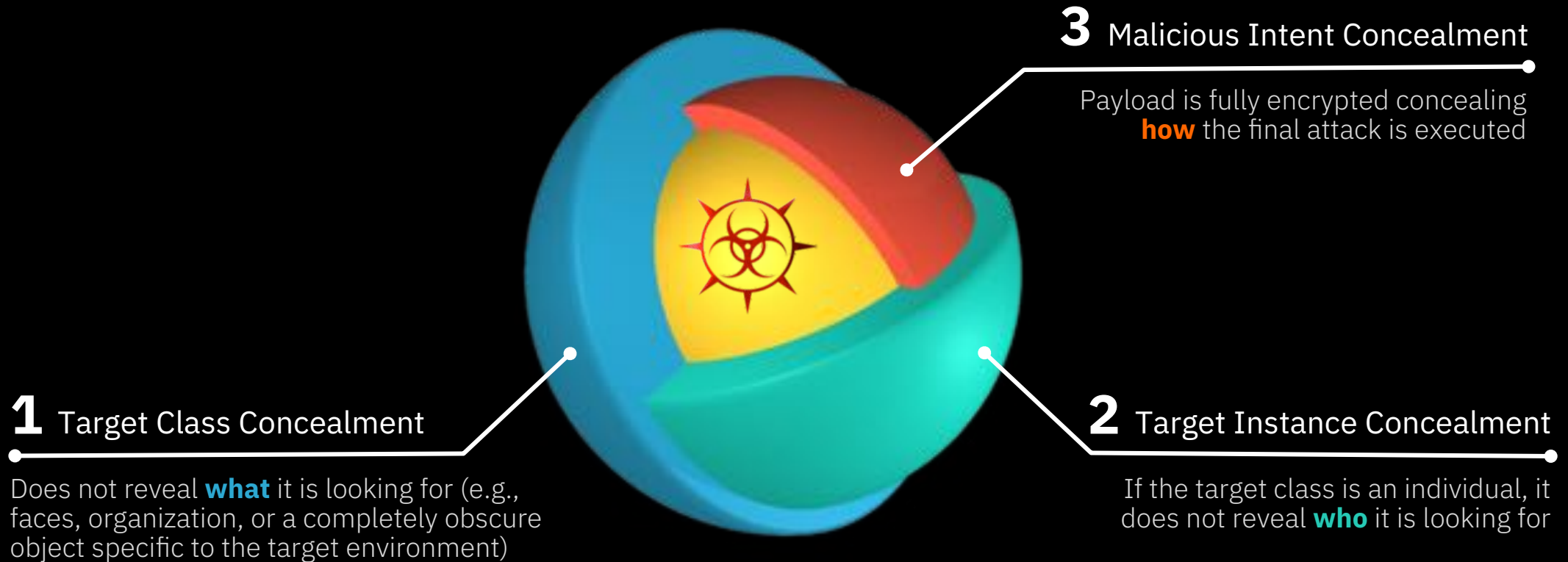


# Analysis of the key generation model



Dataset: Labeled Faces in the Wild (LFW)  
<http://vis-www.cs.umass.edu/lfw/>

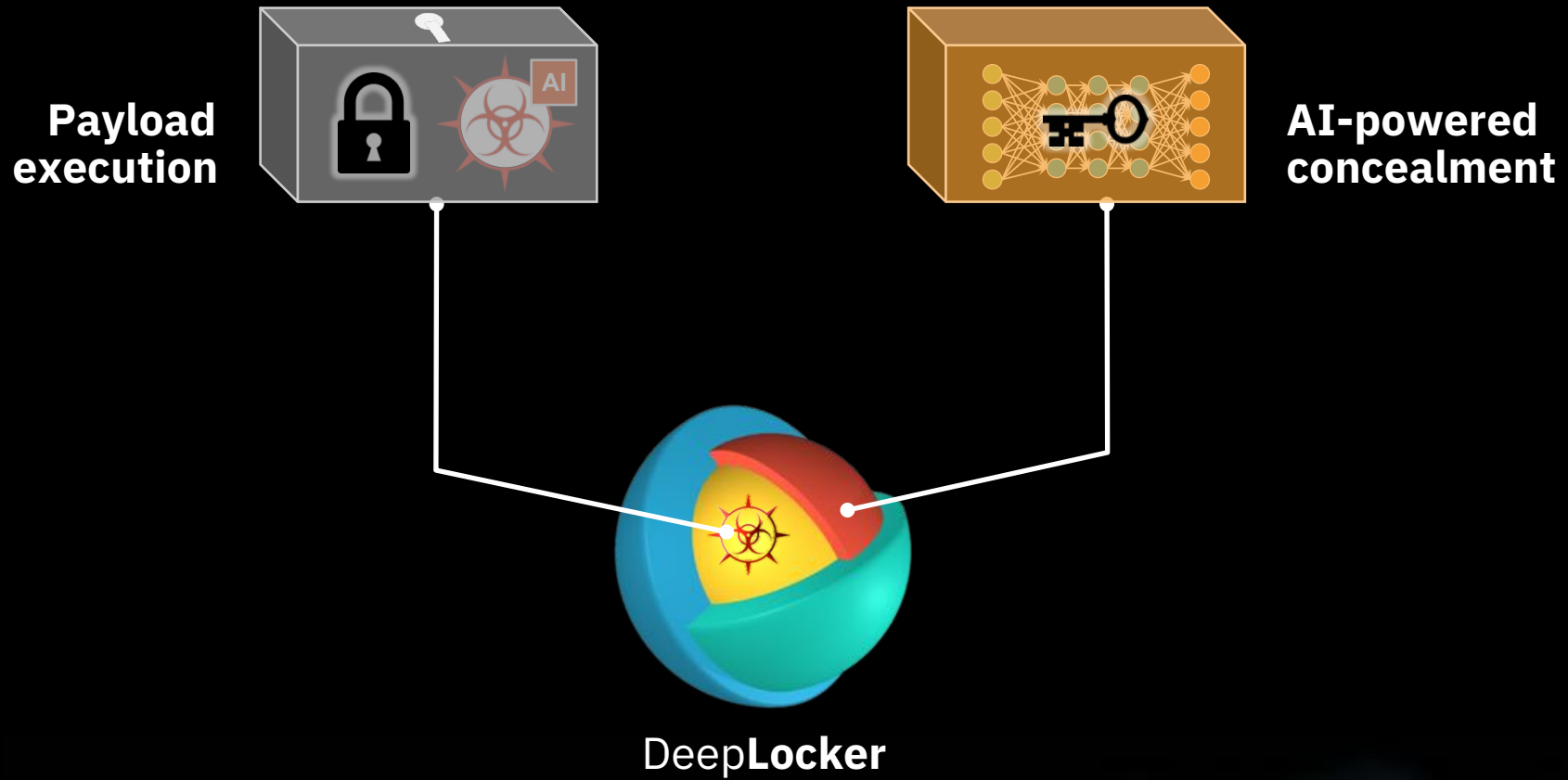
# DeepLocker – AI-powered concealment





# Attacking DeepLocker – AI Lock Picking

# Ways to counter



# Ways to counter

## Payload execution



Code attestation



Host-based monitoring



Brute-force key



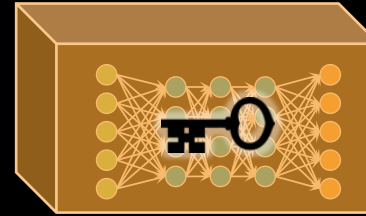
Deceptive resources



Code analysis



## AI-powered concealment



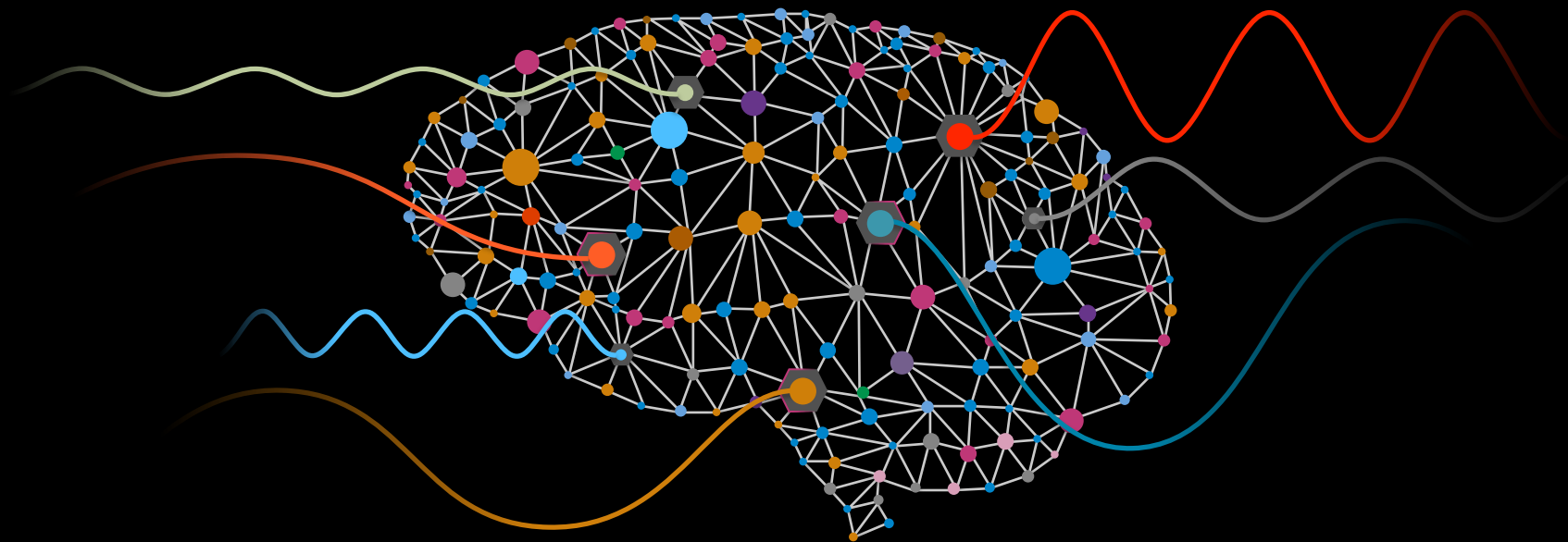
Block sensor access

AI usage monitoring

Brute-force attributes

Deceptive attributes

AI lock picking

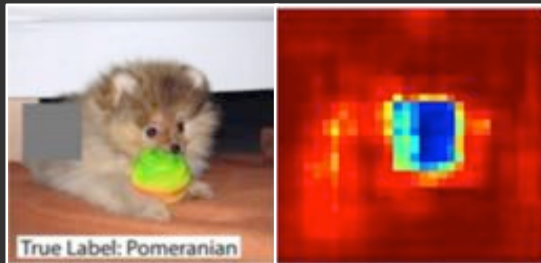


?

# Reverse engineering AI models

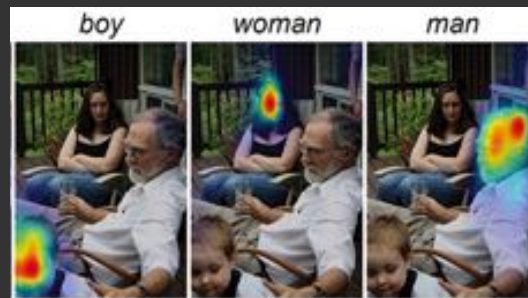
## Partial occlusion

Occlude a portion of the image to see how the embedding is affected (deconvnet) [1]



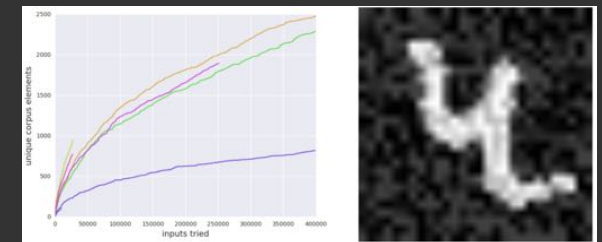
## Neural attention model

Heatmap using the degree of excitation of neurons in each layer (excitation backprop) [2]



## Debug neural networks

Fuzzing for neural networks (coverage-guided fuzzing) [3]



[1] M. Zeiler and R. Fergus, "Visualizing and understanding convolutional networks," ECCV 2014

[2] J. Zhang et. al., "Top-down neural attention by excitation backprop," ECCV 2016

[3] A. Odena and I. Goodfellow, "TensorFuzz: Debugging neural networks with coverage-guided fuzzing," arXiv 2018



# Takeaways

Rapid democratization of AI has made

**AI-powered attacks an imminent threat**

DeepLocker is a demonstration of the potential of

**AI-embedded attacks**

Current defenses will become obsolete and

**new defenses are needed**

# Thank you

Dhilung Kirat

✉ [dkirat@us.ibm.com](mailto:dkirat@us.ibm.com)

Jiyong Jang

✉ [jjang@us.ibm.com](mailto:jjang@us.ibm.com)

Marc Ph. Stoecklin

✉ [mpstoeck@us.ibm.com](mailto:mpstoeck@us.ibm.com)

**IBM Research**

