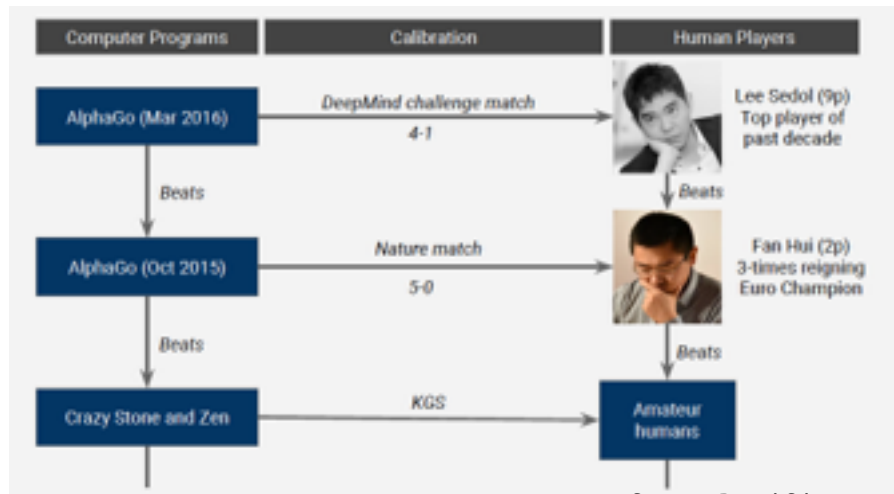


AI and Security: Lessons, Challenges & Future Directions

Dawn Song
UC Berkeley

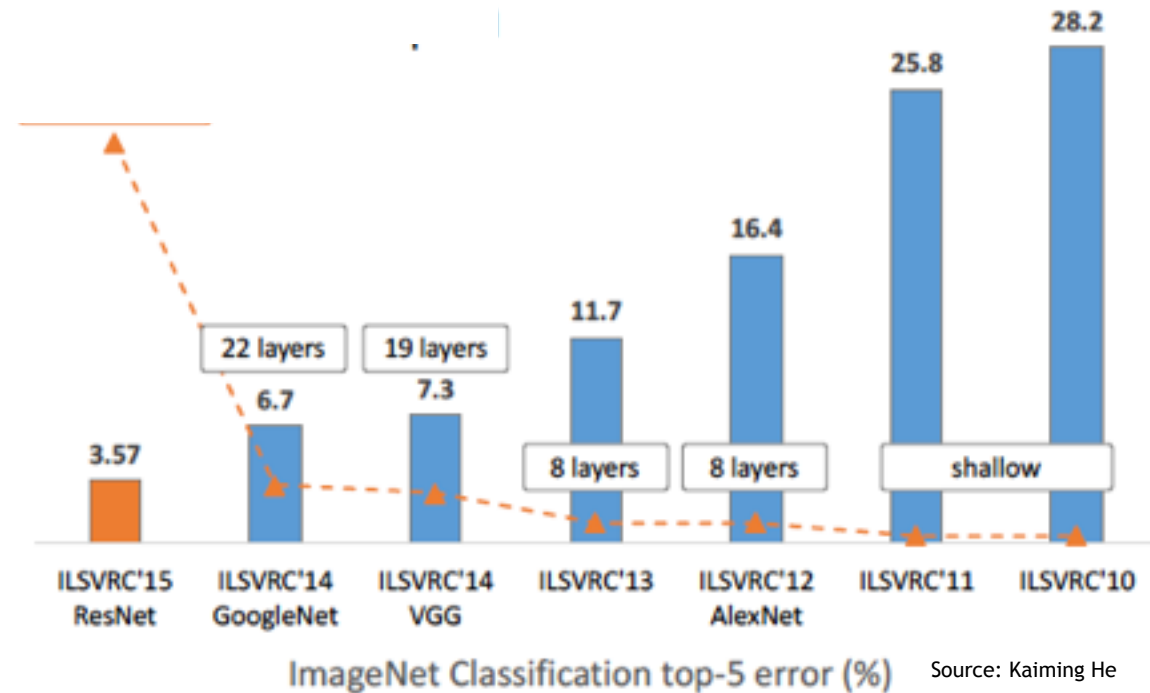
AlphaGo: Winning over World Champion



Source: David Silver



Achieving Human-Level Performance on ImageNet Classification



Deep Learning Powering Everyday Products



pcmag.com



theverge.com





Attacks are increasing in
scale & sophistication



Massive DDoS Caused by IoT Devices



source: Incapsula

Geographical distribution of Mirai bots in recent DDoS attack

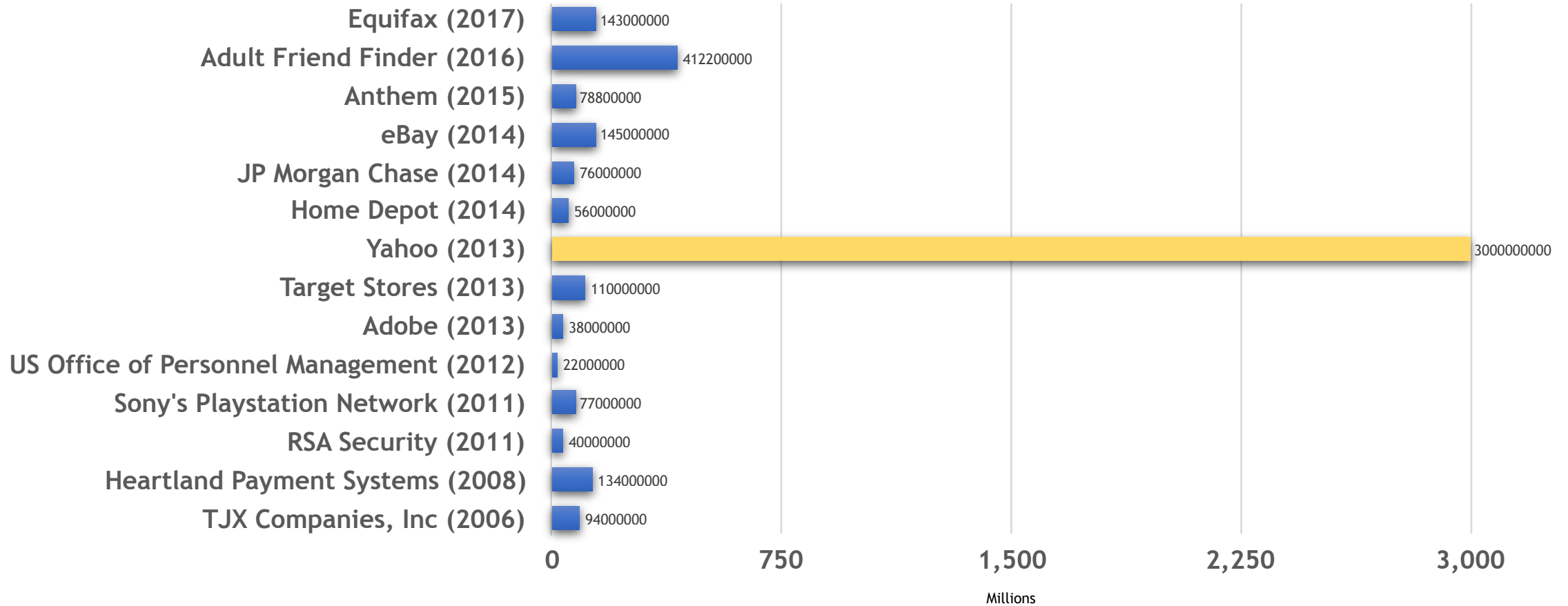
- Botnet of over 400,000 Mirai bots over 160 countries
 - Security cameras/webcams/baby monitors
 - Home routers
- One of the biggest DDoS attacks
 - Over 1Tbps combined attack traffic

WannaCry: One of the Largest Ransomware Breakout



- Used EternalBlue, an exploit of Windows' Server Message Block (SMB) protocol.
- Infected over 200,000 machines across 150 countries in a few days
- Ask for bitcoin payment to unlock encrypted files

Biggest Data Breaches Of the 21st Century



Source: csoonline.com

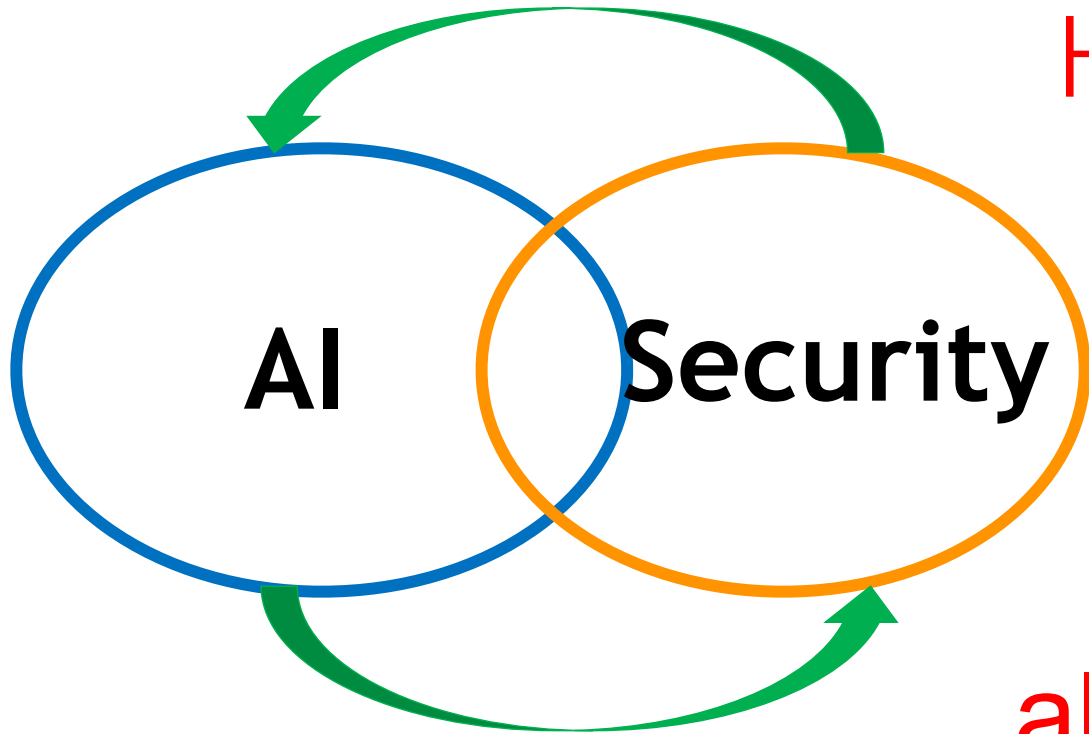
Attacks Entering New Landscape



Ukrain power outage by cyber attack impacted over 250,000 customers



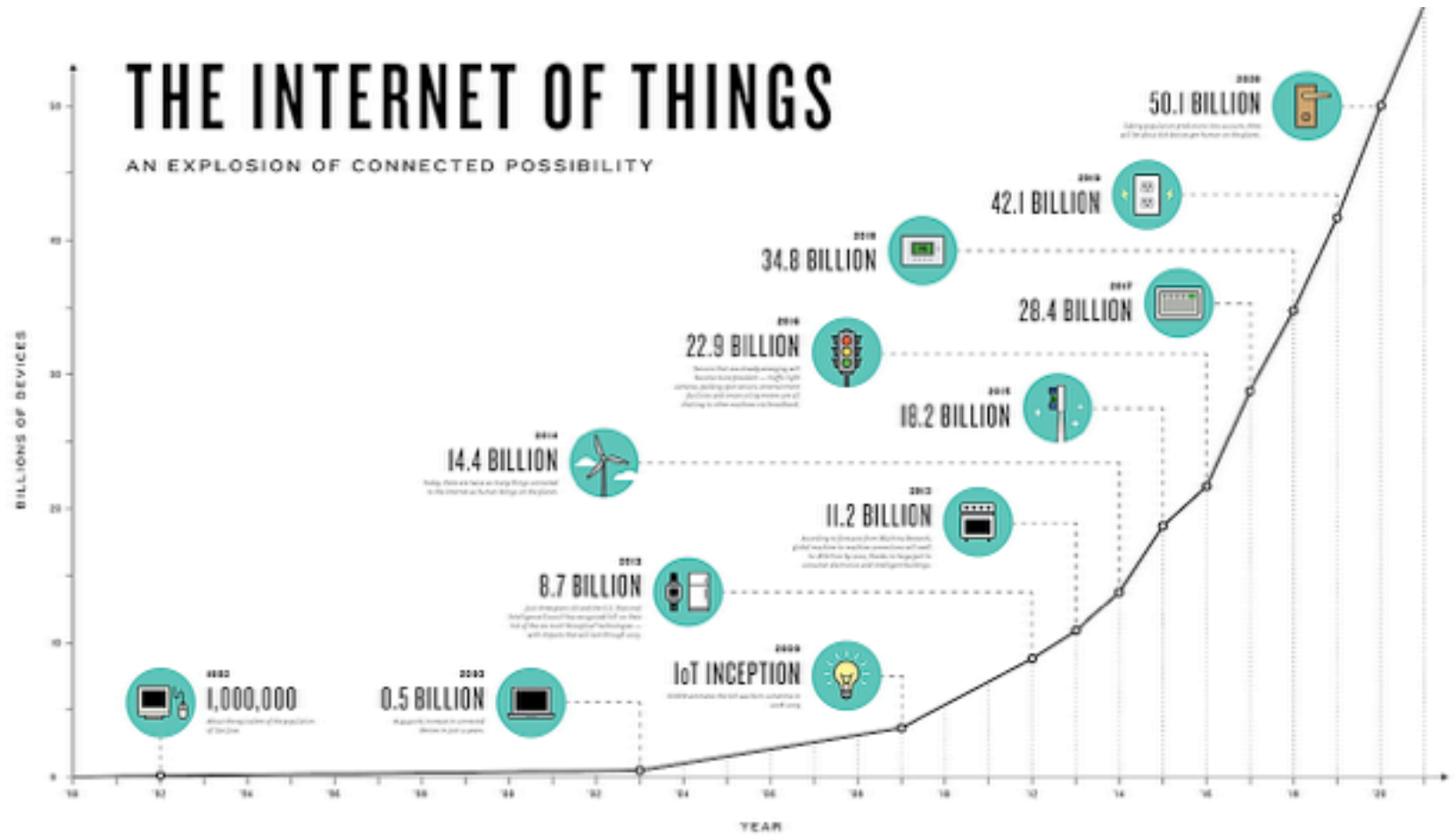
Millions of dollars lost in targeted attacks in SWIFT banking system



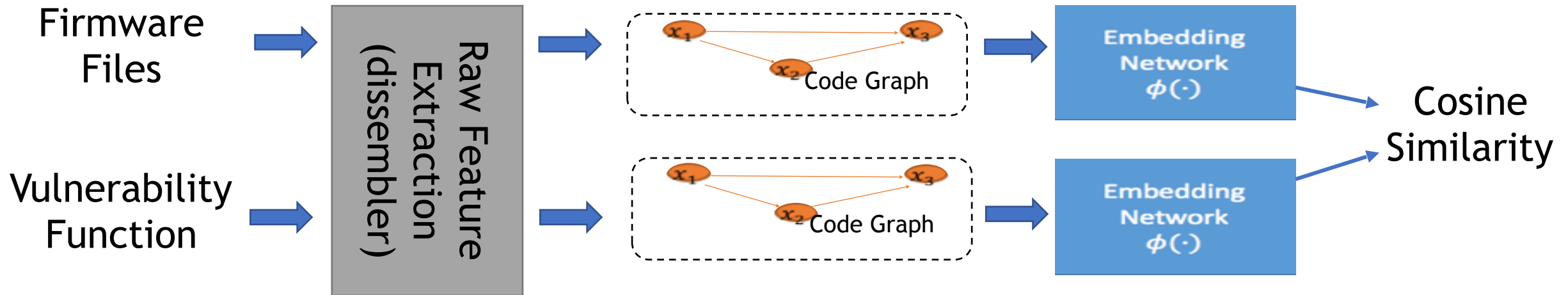
How will (in)security impact the deployment of AI?

How will the rise of AI alter the security landscape?

IoT devices are plagued with vulnerabilities from third-party code



Deep learning for vulnerability detection in IoT Devices



Neural Network-based Graph Embedding for Cross-Platform Binary Code Search
[XLFSSY, ACM Computer and Communication Symposium 2017]

Deep learning for vulnerability detection in IoT Devices

Training time:

Previous work: > 1 week

Our approach: < 30 mins

Serving time (per function):

Previous work: a few mins

Our work: a few milliseconds

10,000 times faster

Function Name	Vendor	Firmware	Binary File	Similarity
u01_get_new_session_ticket	D-Link	DAP-1562_FIRMWARE_1.50	wpa_supplicant.acfgs	0.962374008
port_check_v6	D-Link	DES-1215-28_R1V6_FIRMWARE_3.12.015	in.tftpd.acfgs	0.955400902
sub_42E7C	TP-Link	TD-W8970_V1_140624	raconon.acfgs	0.954742193
sub_42E7C	TP-Link	TD-W8970_V1_130828	raconon.acfgs	0.954742193
prio_parse_file	TP-Link	Archer_G5_V1_140804	raconon.acfgs	0.949834435
sub_41280C	TP-Link	TD-W8970_V1_140624	raconon.acfgs	0.949583828
sub_41280C	TP-Link	TD-W8970_V1_130828	raconon.acfgs	0.949583828
u01_get_new_session_ticket	DD-wrt	DD-wrt v34 13038 NEW2.3 x3.x mega-WNR1000v2 VC	openvpn.acfgs	0.94668215
ucstartsbipServer	TP-Link	WR9400_V2_130115	httpd.acfgs	0.946312308
u01_get_new_session_ticket	Netgear	tomato-Cisco-M30v2-NVRAM12K-1.28.RT-N5x-MIPSr2-110-PL-Mini	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-K26-1.28.RT-MIPSr1-109-Mini	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-K26US-1.28.RT-N5x-MIPSr2-110-VPN	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-64200US-NVRAM0K-1.28.RT-MIPSr2-110-PL-BT	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-E3000US-NVRAM0K-1.28.RT-MIPSr2-110-BT-VPN	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-K26US-1.28.RT-MIPSr1-109-AIO	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-Netgear-E500v2-K26US-1.28.RT-N5x-109-AIO	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-64200US-NVRAM0K-1.28.RT-MIPSr2-109-AIO	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-E1150US-NVRAM0K-1.28.RT-N5x-MIPSr2-110-Nocat-VPN	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-K26US-1.28.RT-N5x-MIPSr2-115-PL-1600N	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-E1150US-NVRAM0K-1.28.RT-N5x-MIPSr2-110-BT-VPN	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-E3000US-NVRAM0K-1.28.RT-MIPSr2-108-PL-VPN	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-E1150US-NVRAM0K-1.28.RT-N5x-MIPSr2-110-Mega-VPN	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-E1200v2-NVRAM0K-1.28.RT-N5x-MIPSr2-108-PL-Max	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-K26US-1.28.RT-MIPSr1-109-Mega-VPN	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-E3000US-NVRAM0K-1.28.RT-MIPSr2-109-Big-VPN	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-64200US-NVRAM0K-1.28.RT-MIPSr2-108-PL-Nocat-VPN	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-Netgear-E500v2-K26US-1.28.RT-N5x-110-ND-AIO	libssl.so.1.0.0.acfgs	0.940912604
u01_get_new_session_ticket	Tomato by Shlubby	tomato-64200US-NVRAM0K-1.28.RT-MIPSr2-109-Nocat-VPN	libssl.so.1.0.0.acfgs	0.940912604

Identified vulnerabilities among top 50:

Previous work: 10/50

Our approach: 42/50

AI Enables Stronger Security Capabilities

- Automatic vulnerability detection & patching
- Automatic agents for attack detection, analysis, & defense



One fundamental weakness of cyber systems is humans

80+% of penetrations and hacks start with a social engineering attack
70+% of nation state attacks [FBI, 2011/Verizon 2014]

AI Enables Chatbot for Phishing Detection



Chatbot for booking flights,
finding restaurants



Chatbot for social engineering attack
detection & defense

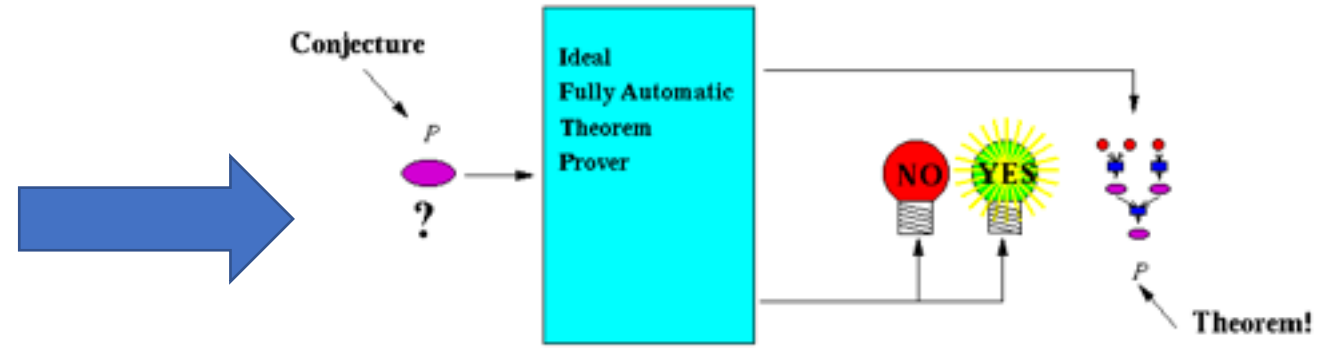
AI Enables Stronger Security Capabilities

- Automatic vulnerability detection & patching
- Automatic agents for attack detection, analysis, & defense
- Automatic verification of software security

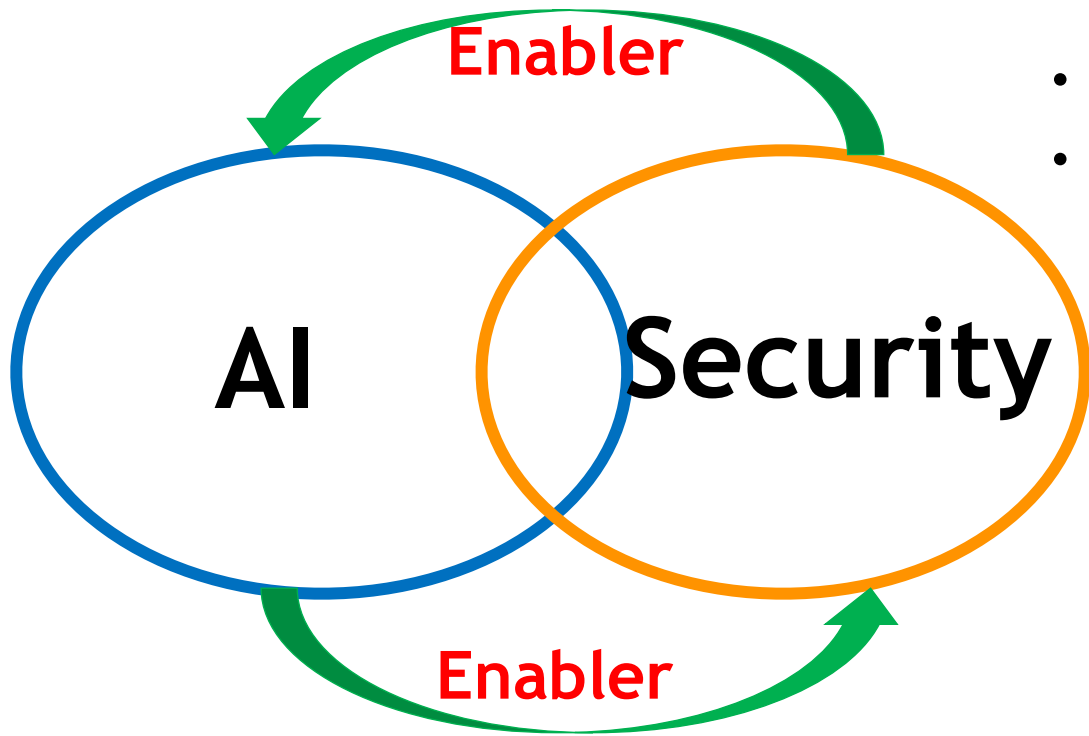
AI Agents to Prove Theorems & Verify Programs



Deep Reinforcement Learning
Agent Learning to Play Go



Automatic Theorem Proving
for Program Verification



- AI enables new security capabilities
- Security enables better AI

Integrity: produces intended/correct results
(adversarial machine learning)

Confidentiality/Privacy: does not leak users' sensitive data
(secure, privacy-preserving machine learning)

Preventing misuse of AI

AI and Security: AI in the presence of attacker

Important to consider the presence of attacker

- History has shown attacker always follows footsteps of new technology development (or sometimes even leads it)
- The stake is even higher with AI
 - As AI controls more and more systems, attacker will have higher & higher incentives
 - As AI becomes more and more capable, the consequence of misuse by attacker will become more and more severe



AI and Security: AI in the presence of attacker

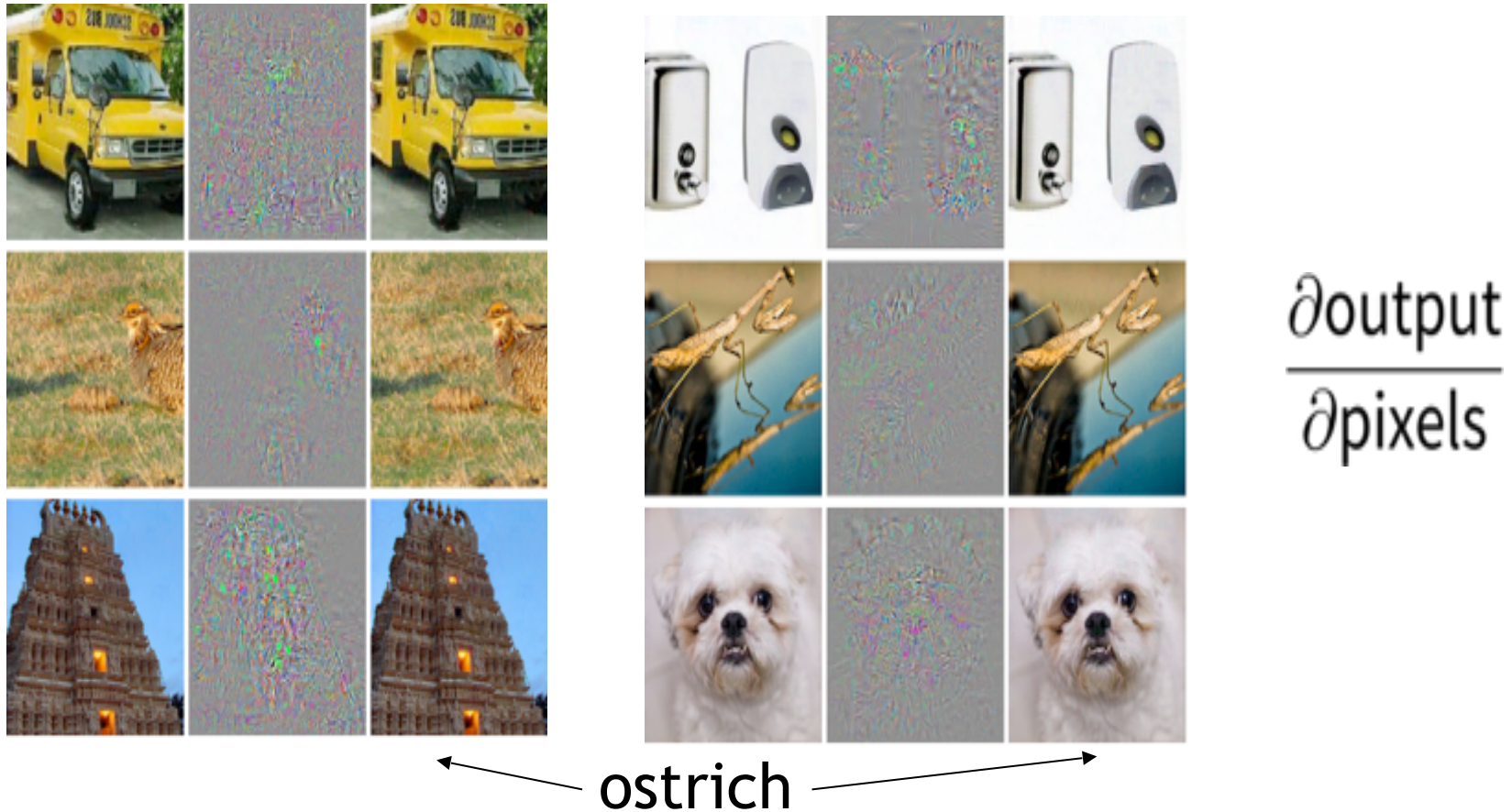
- **Attack AI**

- Cause the learning system to not produce intended/correct results
- Cause learning system to produce targeted outcome designed by attacker
- Learn sensitive information about individuals
- Need security in learning systems

- **Misuse AI**

- Misuse AI to attack other systems
 - Find vulnerabilities in other systems; Devise attacks
- Need security in other systems

Deep Learning Systems Are Easily Fooled



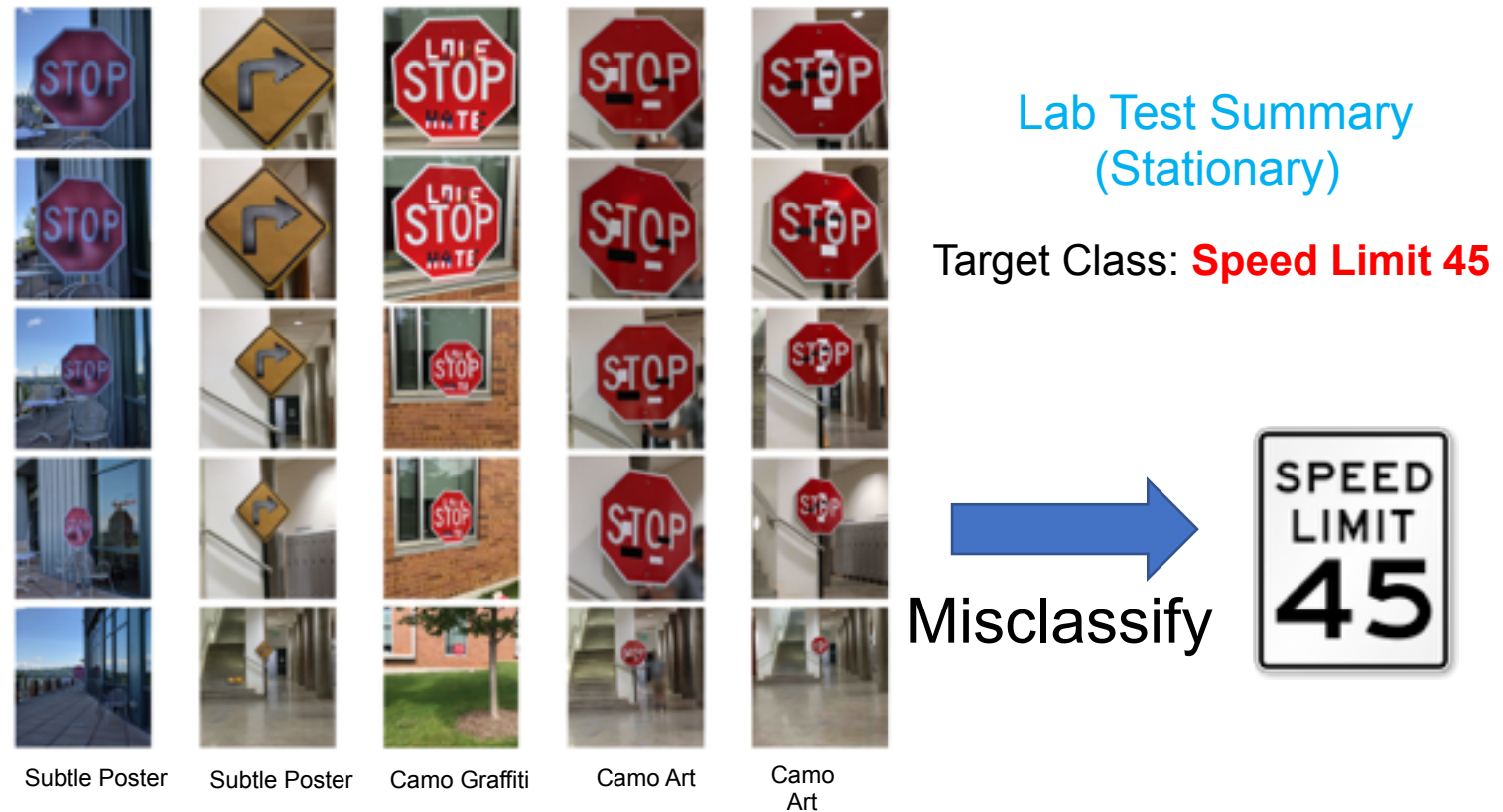




STOP Signs in Berkeley

Adversarial Examples in Physical World

Adversarial examples in physical world **remain effective under different viewing distances, angles, other conditions**



Drive-by Test

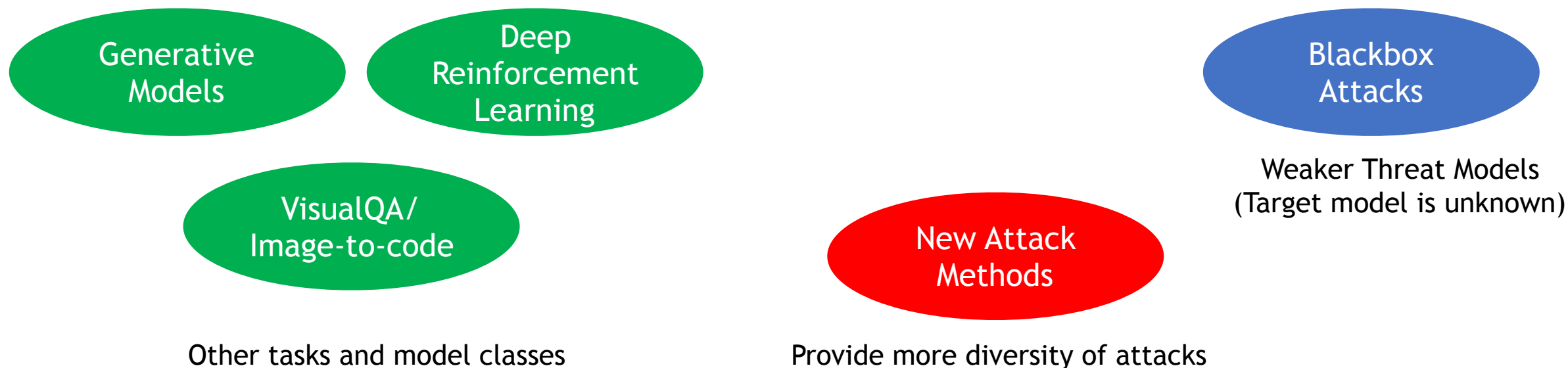
Adversarial examples in
physical world
&
remain effective under
different viewing distances,
angles, other conditions



Adversarial Examples Are Prevalent in
Deep Learning Systems

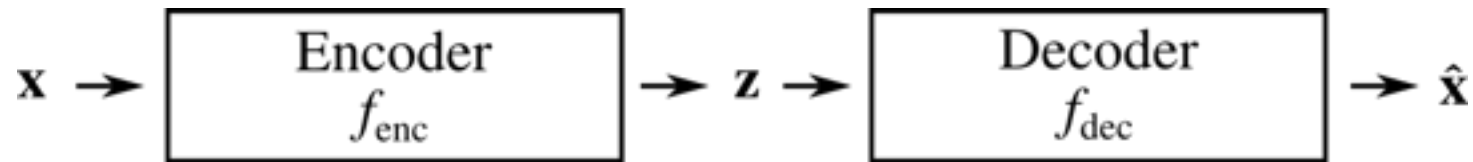
Adversarial Examples Prevalent in Deep Learning Systems

- Most existing work on adversarial examples:
 - Image classification task
 - Target model is known
- Our investigation on adversarial examples:



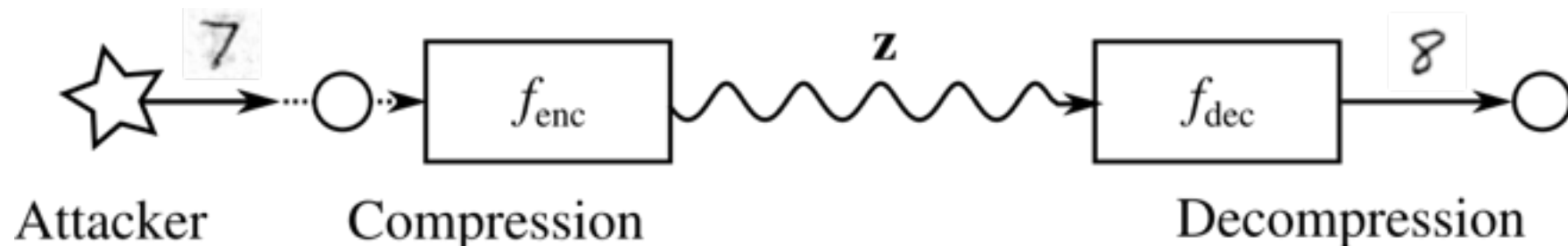
Generative models

- VAE-like models (VAE, VAE-GAN) use an intermediate latent representation
- An **encoder**: maps a high-dimensional input into lower-dimensional latent representation \mathbf{z} .
- A **decoder**: maps the latent representation back to a high-dimensional reconstruction.



Adversarial Examples in Generative Models

- An example attack scenario:
 - Generative model used as a compression scheme



- Attacker's goal: for the decompressor to reconstruct a different image from the one that the compressor sees.

Adversarial Examples for VAE-GAN in MNIST

7 2 1 0 4 1 4 9 5 9
0 6 9 0 1 5 9 7 3 4
9 6 6 5 4 0 7 4 0 1
3 1 3 4 7 2 7 1 2 1
1 7 4 2 3 5 1 2 4 4
6 3 5 5 6 0 4 1 9 5
7 8 9 3 7 4 6 4 3 0
7 0 2 9 1 7 3 2 9 7
7 6 2 7 8 4 7 3 6 1
3 6 9 3 1 4 1 7 6 9

Original images

7 2 1 0 4 1 4 9 5 7
0 6 9 0 1 5 9 7 3 4
7 6 6 5 4 0 7 4 0 1
3 1 3 4 7 2 7 1 2 1
1 7 4 2 3 5 1 2 4 4
6 3 5 5 6 0 4 1 9 5
7 8 9 3 7 9 6 4 3 0
7 0 2 7 1 7 3 2 9 7
7 6 2 7 8 4 7 3 6 1
3 6 9 3 1 4 1 7 6 9

Reconstruction of original images

Target Image



7 2 1 4 1 4 9 5 9 6
9 1 5 9 7 3 4 9 6 6
5 4 7 4 1 3 1 3 4 7
2 7 1 2 1 1 7 4 2 3
5 1 2 4 4 6 3 5 5 6
4 1 9 5 7 8 9 3 7 4
6 4 3 7 2 9 1 7 3 2
9 7 7 6 2 7 8 4 7 3
6 1 3 6 9 3 1 4 1 7
6 9 6 5 4 9 9 2 1 9

Adversarial examples

0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0

Reconstruction of adversarial examples

Adversarial Examples for VAE-GAN in SVHN

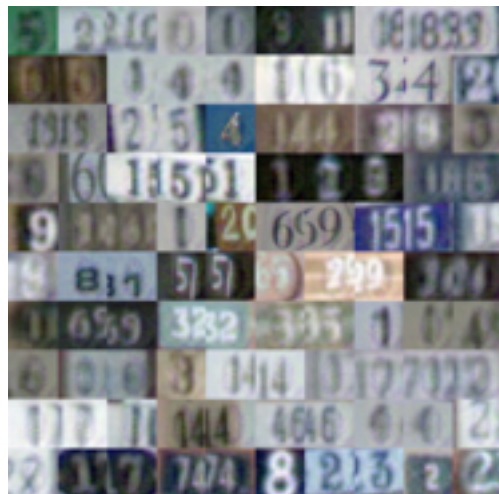


Original images



Reconstruction of original images

Target Image



Adversarial examples



Reconstruction of adversarial examples

Adversarial Examples for VAE-GAN in SVHN



Original images



Reconstruction of original images

Target Image

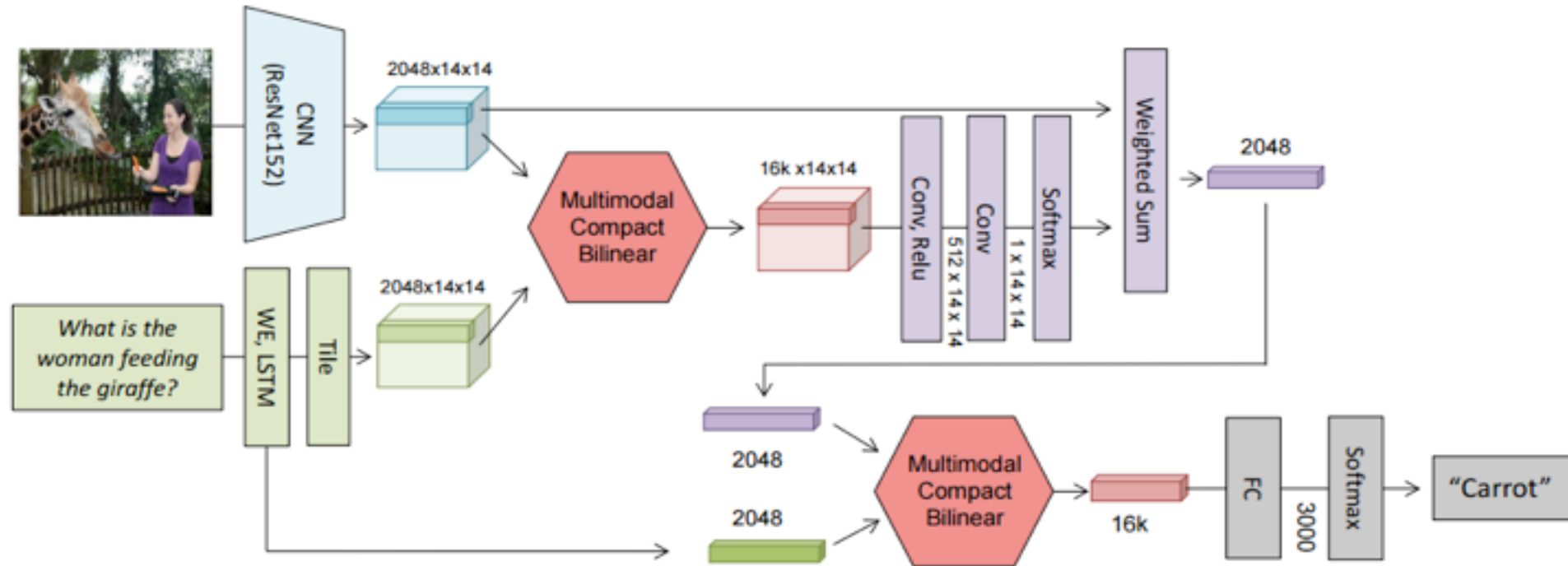


Adversarial examples



Reconstruction of adversarial examples

Visual Question & Answer (VQA)



Multimodal Compact Bilinear Pooling for Visual Question Answering and Visual Grounding, Fukui et al., <https://arxiv.org/abs/1606.01847>

Q: Where is the plane?



Benign image



Mode



Answer:
Runway

Fooling VQA

Target: Sky



Adversarial example



Mode



Sky

Q: How many cats are there?



Benign image



Mode



1

Answer:

Fooling VQA

Target: 2



Adversarial example

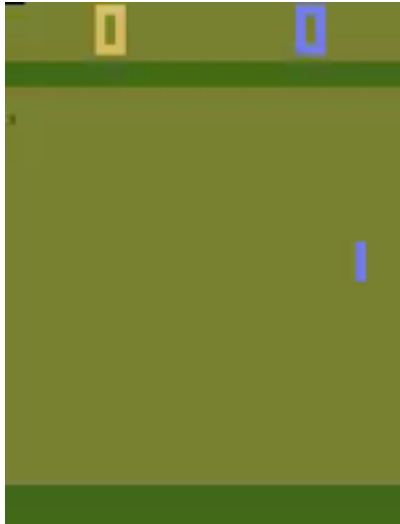


Mode

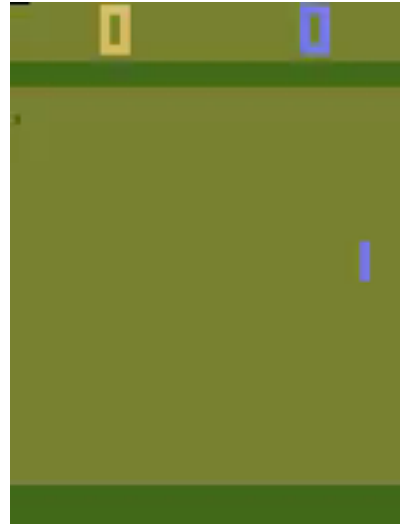


2

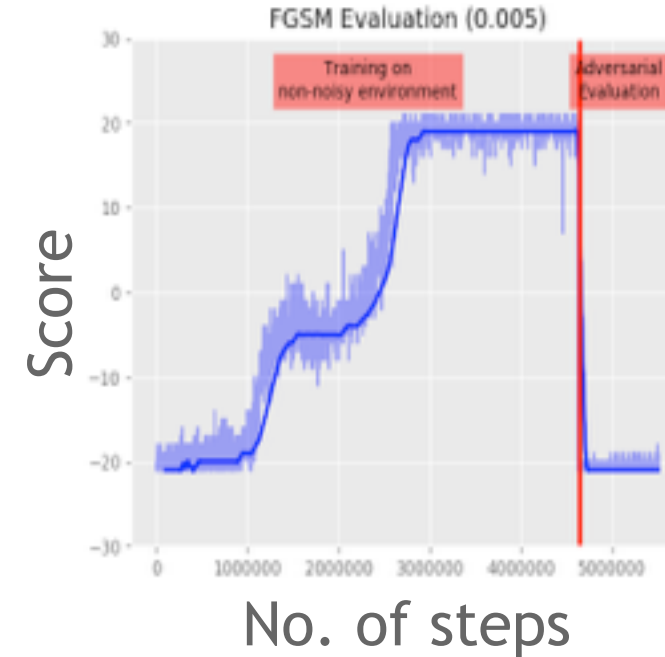
Adversarial Examples Fooling Deep Reinforcement Learning Agents



Original Frames



Original Frames with
Adversarial Perturbation



Jernej Kos and Dawn Song: *Delving into adversarial attacks on deep policies* [ICLR Workshop 2017].

A General Framework for Black-box attacks

- Zero-Query Attack (Previous methods)
 - Random perturbation
 - Difference of means
 - Transferability-based attack
 - Practical Black-Box Attacks against Machine Learning [Papernot et al. 2016]
 - Ensemble transferability-based attack [Yanpei Liu, Xinyun Chen, Chang Liu, Dawn Song: Delving into Transferable Adversarial Examples and Black-box Attacks, ICLR 2017]
- Query Based Attack (new method)
 - Finite difference gradient estimation
 - Query reduced gradient estimation
 - Results: similar effectiveness to whitebox attack
 - A general active query game model

Black-box Attack on Clarifai



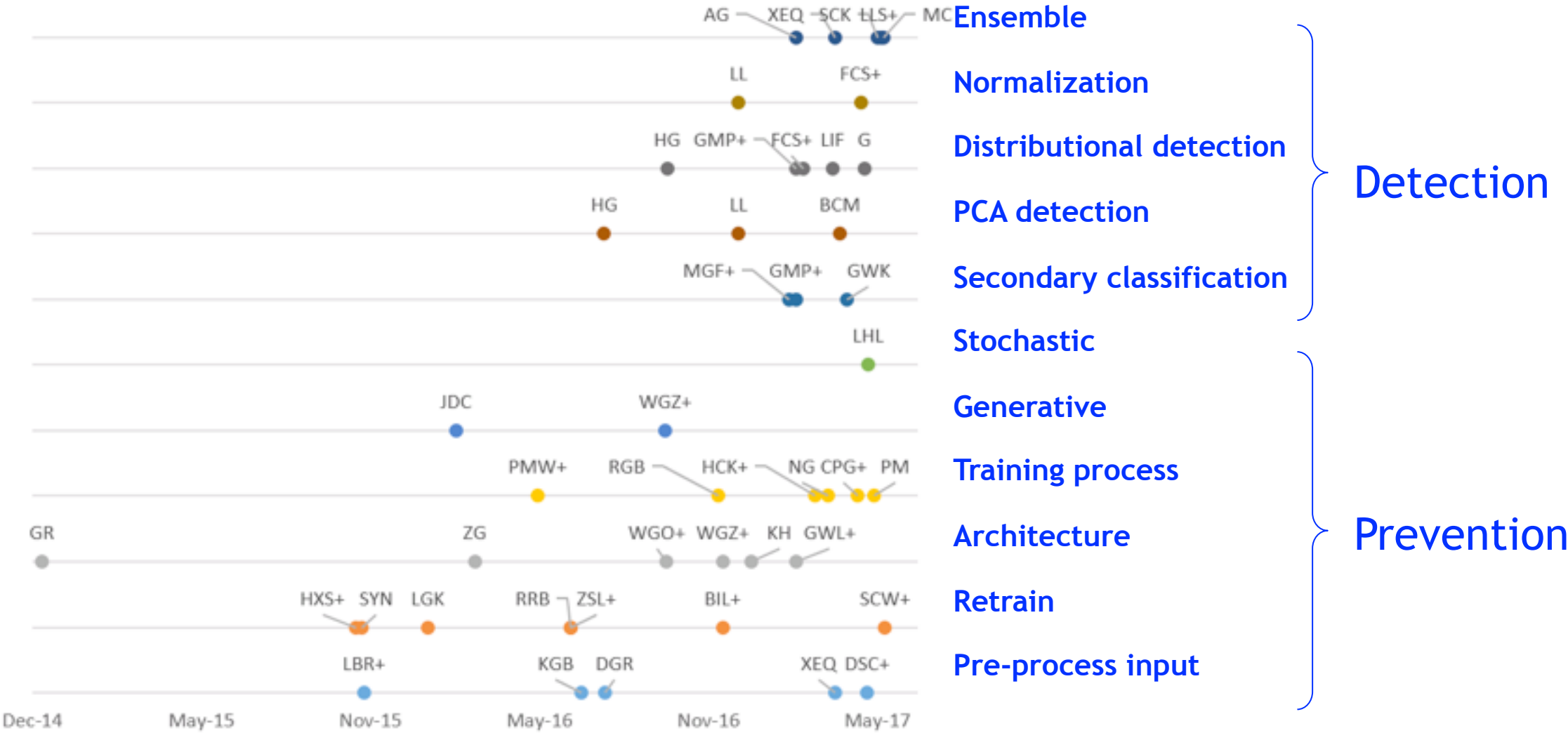
Original image, classified as
“drug” with a confidence of 0.99



Adversarial example, classified as
as “safe” with a confidence of
0.96

The Gradient-Estimation black-box attack on Clarifai’s Content Moderation Model

Numerous Defenses Proposed



No Sufficient Defense Today

- Strong, adaptive attacker can easily evade today's defenses
- Ensemble of weak defenses does not (by default) lead to strong defense
 - Warren He, James Wei, Xinyun Chen, Nicholas Carlini, Dawn Song [WOOT 2017]
- Adversarial Examples Are Not Easily Detected: Bypassing Ten Detection Methods
 - Nicholas Carlini and David Wagner [AISeC 2017]

Adversarial Machine Learning

- Adversarial machine learning:
 - Learning in the presence of adversaries
- Inference time: adversarial example fools learning system
 - Evasion attacks
 - Evade malware detection; fraud detection
- Training time:
 - Attacker poisons training dataset (e.g., poison labels) to fool learning system to learn wrong model
 - Poisoning attacks: e.g., Microsoft's Tay twitter chatbot
 - Attacker selectively shows learner training data points (even with correct labels) to fool learning system to learn wrong model
 - Data poisoning is particularly challenging with crowd-sourcing & insider attack
 - Difficult to detect when the model has been poisoned
- Adversarial machine learning particularly important for security critical system

Security will be one of the biggest challenges in Deploying AI



Security of Learning Systems

- Software level
- Learning level
- Distributed level

Challenges for Security at Software Level

- No software vulnerabilities (e.g., buffer overflows & access control issues)
 - Attacker can take control over learning systems through exploiting software vulnerabilities

Challenges for Security at Software Level

- No software vulnerabilities (e.g., buffer overflows & access control issues)
- Existing software security/formal verification techniques apply

Reactive Defense

Proactive Defense:
Bug Finding

Proactive Defense:
Secure by
Construction



Automatic worm detection
& signature/patch generation

Automatic malware
detection & analysis



Progression of different approaches to software security over last 20 years

Security of Learning Systems

- Software level
- Learning level
- Distributed level

Challenges for Security at Learning Level

- Evaluate system under adversarial events, not just normal events

Regression Testing vs. Security Testing in Traditional Software System

	Regression Testing	Security Testing
Operation	Run program on normal inputs	Run program on abnormal/adversarial inputs
Goal	Prevent normal users from encountering errors	Prevent attackers from finding exploitable errors

Regression Testing vs. Security Testing in Learning System

	Regression Testing	Security Testing
Training	Train on noisy training data: Estimate resiliency against noisy training inputs	Train on poisoned training data: Estimate resiliency against poisoned training inputs
Testing	Test on normal inputs: Estimate generalization error	Test on abnormal/ adversarial inputs: Estimate resiliency against adversarial inputs

Challenges for Security at Learning Level

- Evaluate system under adversarial events, not just normal events
 - Regression testing vs. security testing
- Reason about complex, non-symbolic programs

Decades of Work on Reasoning about Symbolic Programs

- Symbolic programs:
 - E.g., OS, File system, Compiler, web application, mobile application
 - Semantics defined by logic
 - Decades of techniques & tools developed for logic/symbolic reasoning
 - Theorem provers, SMT solvers
 - Abstract interpretation

Era of Formally Verified Systems

Verified: Micro-kernel, OS, File system, Compiler, Security protocols, Distributed systems



IronClad/IronFleet

FSCQ

CertiKOS

miTLS/Everest

EasyCrypt

CompCert

Powerful Formal Verification Tools + Dedicated Teams



Why3



Z3



No Sufficient Tools to Reason about Non-Symbolic Programs

- Symbolic programs:



- Semantics defined by logic
- Decades of techniques & tools developed for logic/symbolic reasoning
 - Theorem provers, SMT solvers
 - Abstract interpretation

- Non-symbolic programs:



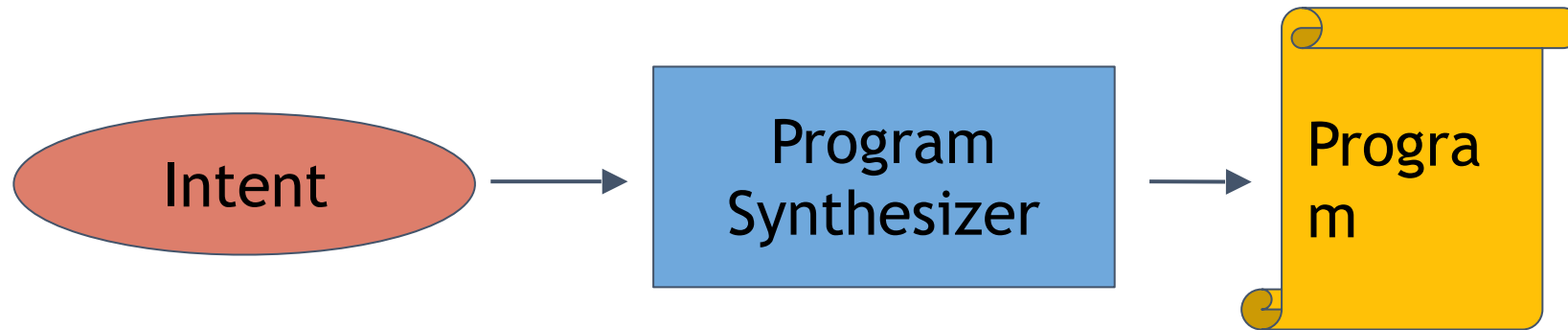
- No precisely specified properties & goals
- No good understanding of how learning system works
- Traditional symbolic reasoning techniques do not apply

Challenges for Security at Learning Level

- Evaluate system under adversarial events, not just normal events
 - Regression testing vs. security testing
- Reason about complex, non-symbolic programs
- Design new architectures & approaches with stronger generalization & security guarantees

Neural Program Synthesis

Can we teach computers to write code?



Example Applications:

- End-user programming
- Performance optimization of code
- Virtual assistant

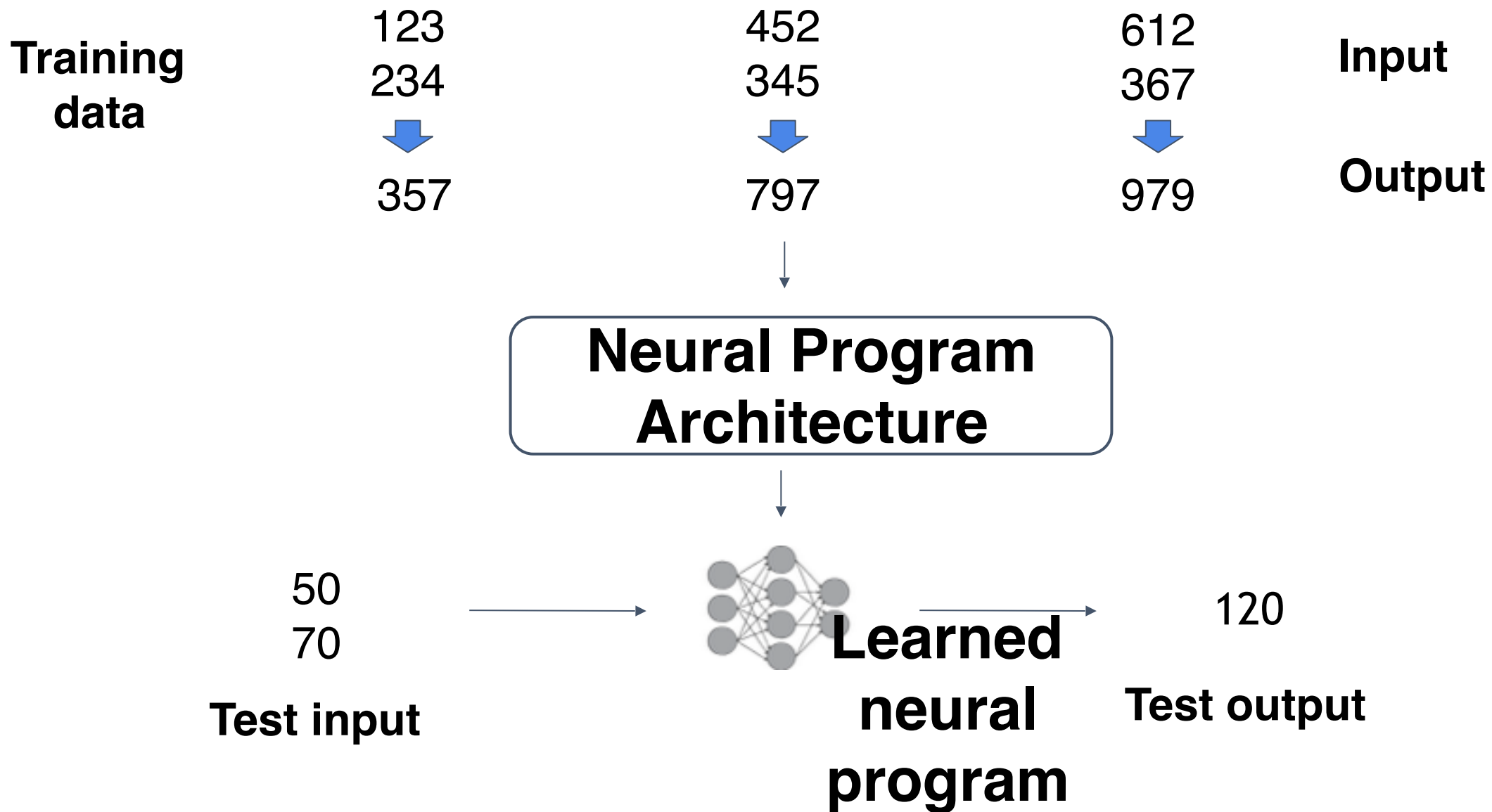
“Software is eating the world” --- az16

Program synthesis can automate this & democratize idea realization

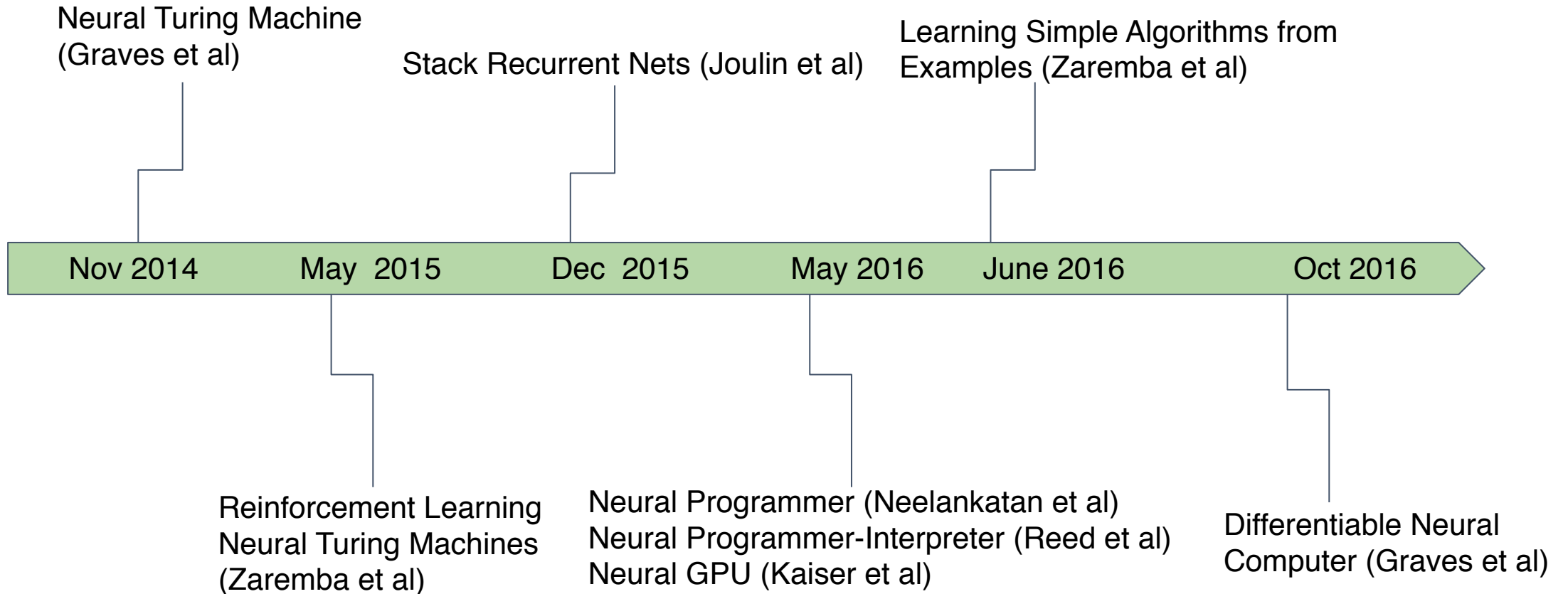
Neural Program Synthesis

Training data	123	452	612	Input
	234	345	367	
	↓	↓	↓	
	357	797	979	Output

Neural Program Synthesis

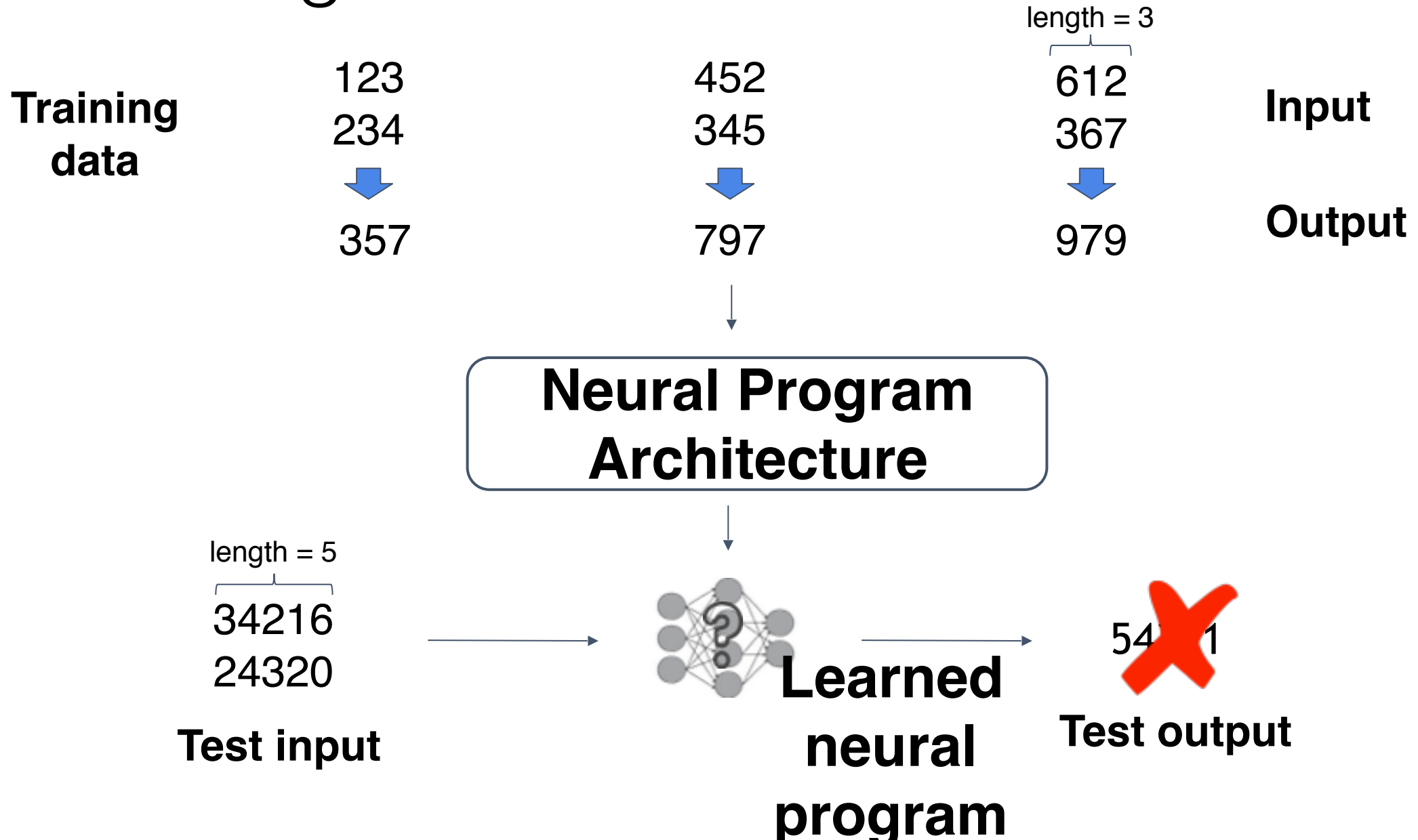


Neural Program Architectures

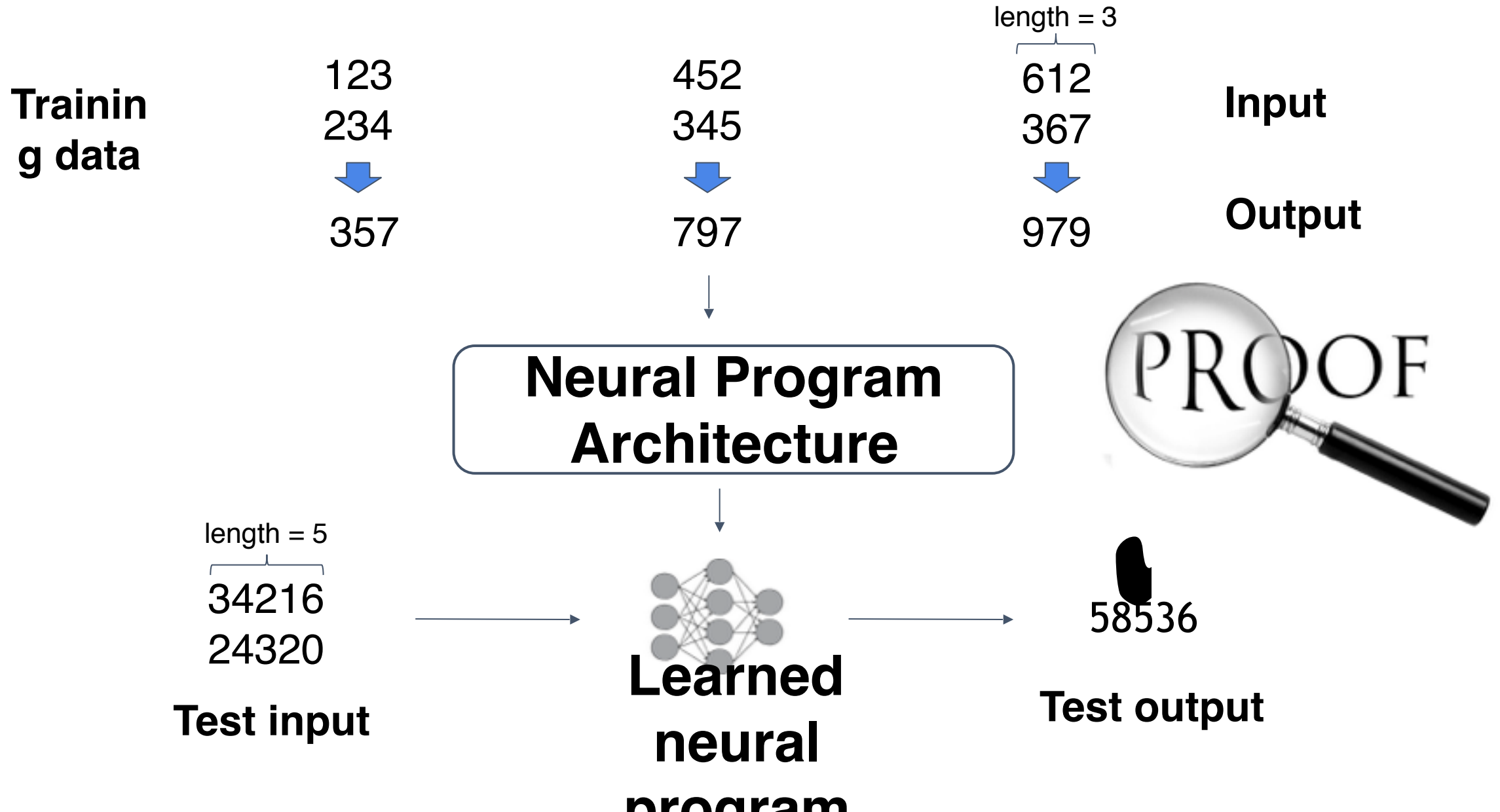


Neural Program Synthesis Tasks: Copy, Grade-school addition, Sorting, Shortest Path

Challenge 1: Generalization



Challenge 2: No Proof of Generalization



Our Approach: Introduce Recursion

Learn recursive neural programs

Recursion

- Fundamental concept in Computer Science and Math
- Solve whole problem by reducing it to smaller subproblems (*reduction rules*)
- *Base cases* (smallest subproblems) are easier to reason about



Quicksort

Our Approach: Making Neural Programming Architectures Generalize via Recursion

- **Proof of Generalization:**
 - Recursion enables provable guarantees about neural programs
 - Prove perfect generalization of a learned recursive program via a verification procedure
 - Explicitly testing on all possible base cases and reduction rules (Verification set)
- Learn & generalize faster as well
 - Trained on same data, non-recursive programs do not generalize well

Accuracy on Random Inputs for Quicksort

Length of Array	Non-Recursive	Recursive
3	100%	100%
5	100%	100%
7	100%	100%
11	73.3%	100%
15	60%	100%
20	30%	100%
22	20%	100%
25	3.33%	100%
30	3.33%	100%
70	0%	100%



Lessons

- Program architecture impacts generalization & provability
- Recursive, modular neural architectures are easier to reason, prove, generalize
- Explore new architectures and approaches enabling strong generalization & security properties for broader tasks

Challenges for Security at Learning Level

- Evaluate system under adversarial events, not just normal events
- Reason about complex, non-symbolic programs
- Design new architectures & approaches with stronger generalization & security guarantees
- Reason about how to compose components

Compositional Reasoning

- Building large, complex systems require compositional reasoning
 - Each component provides abstraction
 - E.g., pre/post conditions
 - Hierarchical, compositional reasoning proves properties of whole system
- How to do abstraction, compositional reasoning for non-symbolic programs?

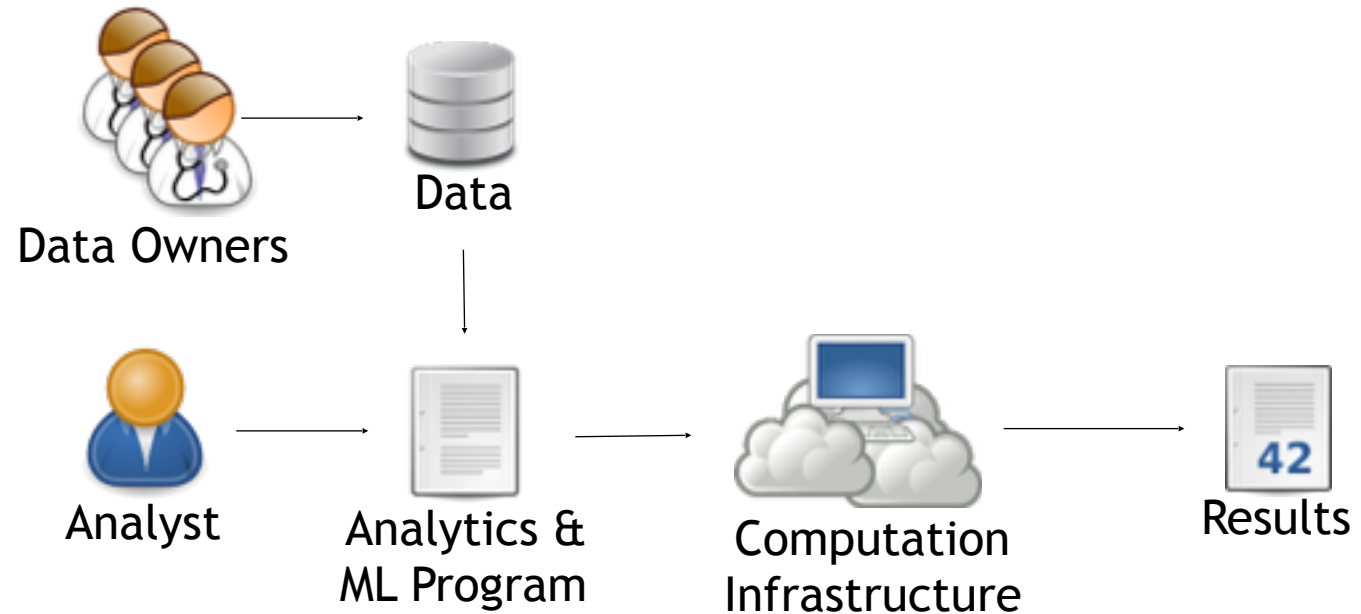
Security of Learning Systems

- Software level
- Learning level
 - Evaluate system under adversarial events, not just normal events
 - Reason about complex, non-symbolic programs
 - Design new architectures & approaches with stronger generalization & security guarantees
 - Reason about how to compose components
- Distributed level
 - Each agent makes local decisions; how to make good local decisions achieve good global decision?

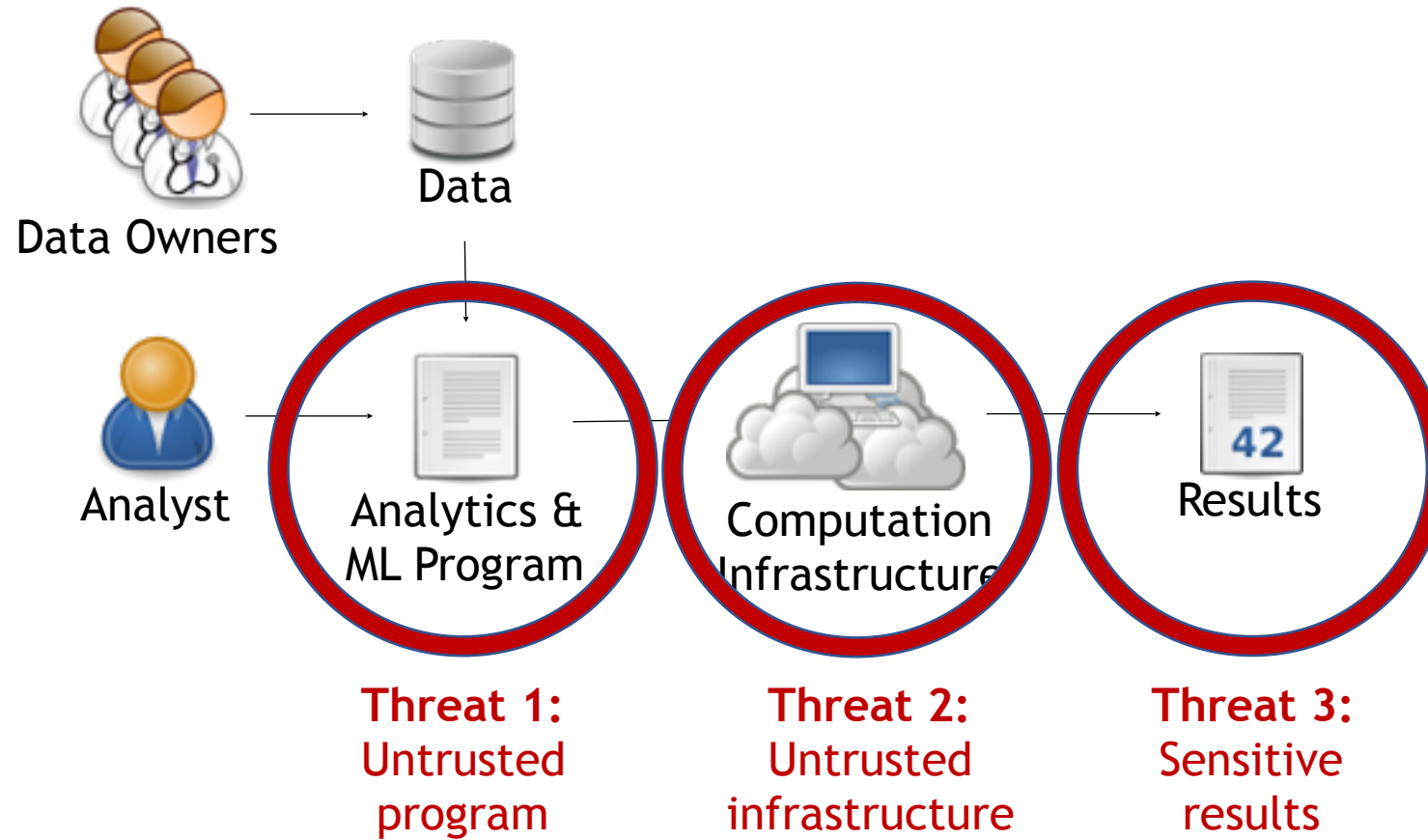
AI and Security: AI in the presence of attacker

- Attack AI
 - Integrity:
 - Cause learning system to not produce intended/correct results
 - Cause learning system to produce targeted outcome designed by attacker
 - Confidentiality:
 - Learn sensitive information about individuals
 - Need security in learning systems
- Misuse AI
 - Misuse AI to attack other systems
 - Find vulnerabilities in other systems
 - Target attacks
 - Devise attacks
 - Need security in other systems

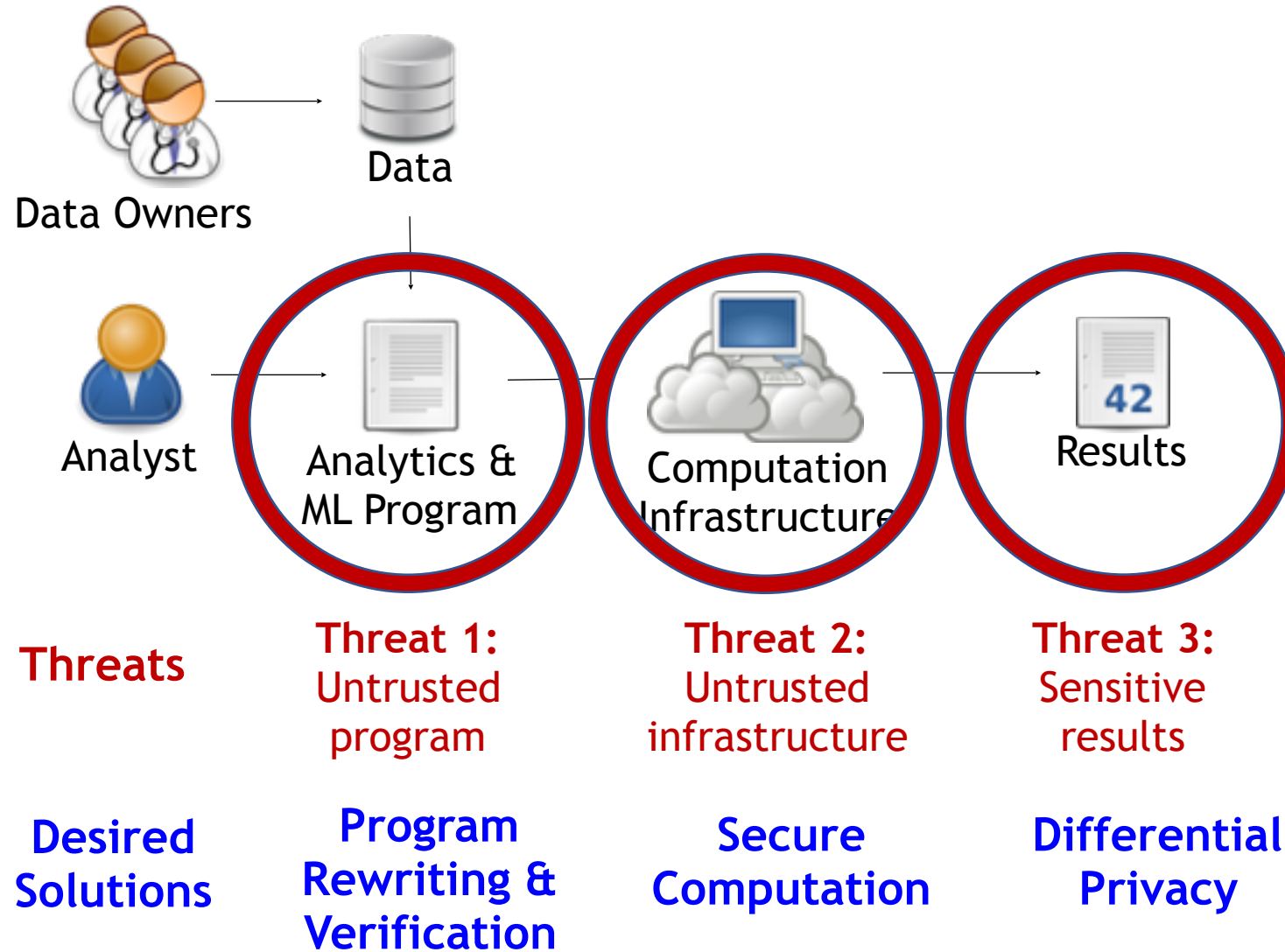
Current Frameworks for Data Analytics & Machine Learning



Current Frameworks Insufficient



Desired Solutions for Confidentiality/Privacy



AI and Security: AI in the presence of attacker

- Attack AI
 - Integrity:
 - Cause learning system to not produce intended/correct results
 - Cause learning system to produce targeted outcome designed by attacker
 - Confidentiality:
 - Learn sensitive information about individuals
 - Need security in learning systems
- Misuse AI
 - Misuse AI to attack other systems
 - Find vulnerabilities in other systems
 - Target attacks
 - Devise attacks
 - Need security in other systems

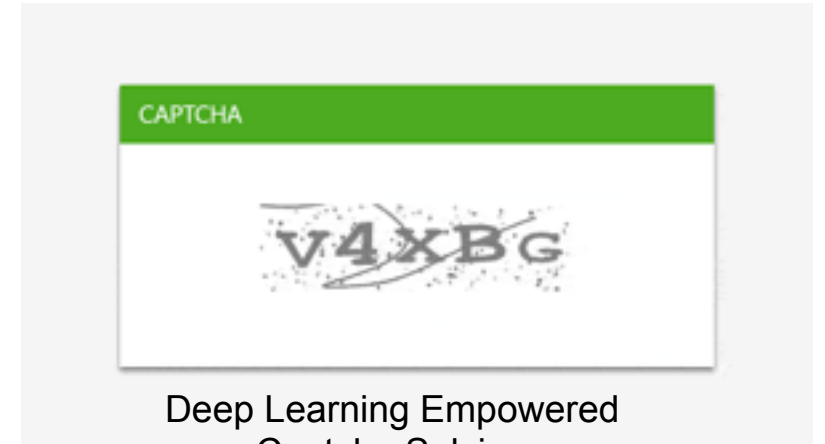
Misused AI can make attacks more effective



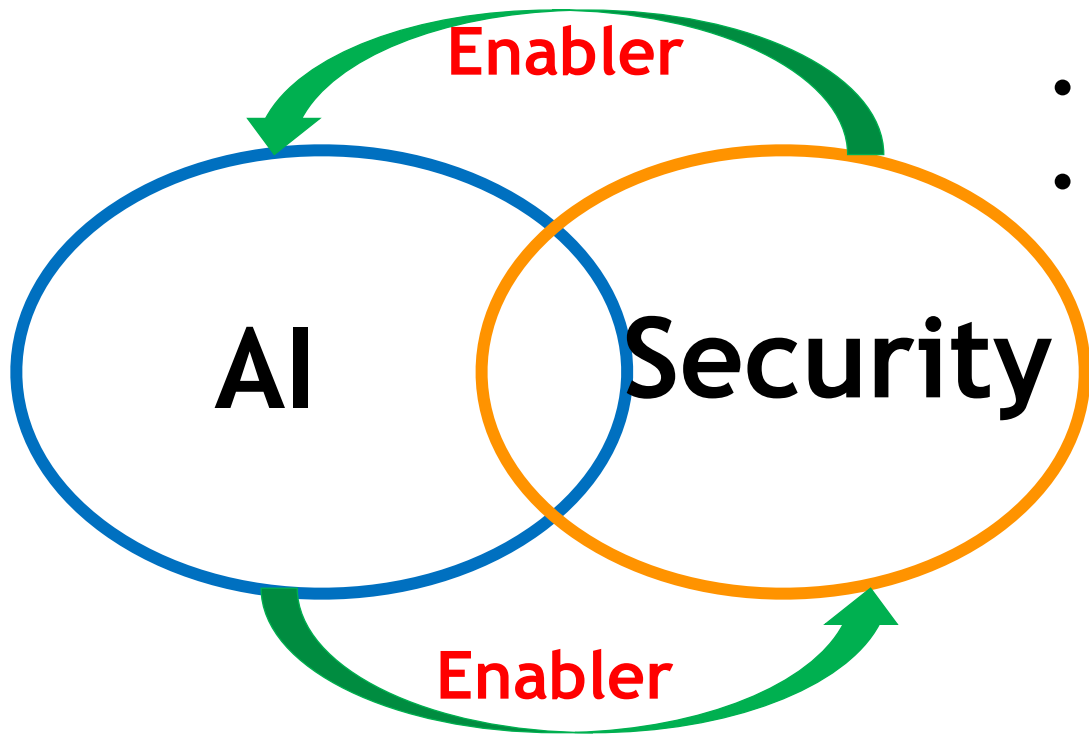
Deep Learning Empowered
Bug Finding



Deep Learning Empowered
Phishing Attacks



Deep Learning Empowered
Captcha Solving



- AI enables new security capabilities
- Security enables better AI

Integrity: produces intended/correct results (adversarial machine learning)

Confidentiality/Privacy: does not leak users' sensitive data (secure, privacy-preserving machine learning)

Preventing misuse of AI

Future of AI and Security

How to better understand what security means for AI, learning systems?

How to detect when a learning system has been fooled/compromised?

How to build better resilient systems with stronger guarantees?

How to build privacy-preserving learning systems?

Security will be one of the biggest challenges in Deploying AI.

Let's tackle the big challenges together!



