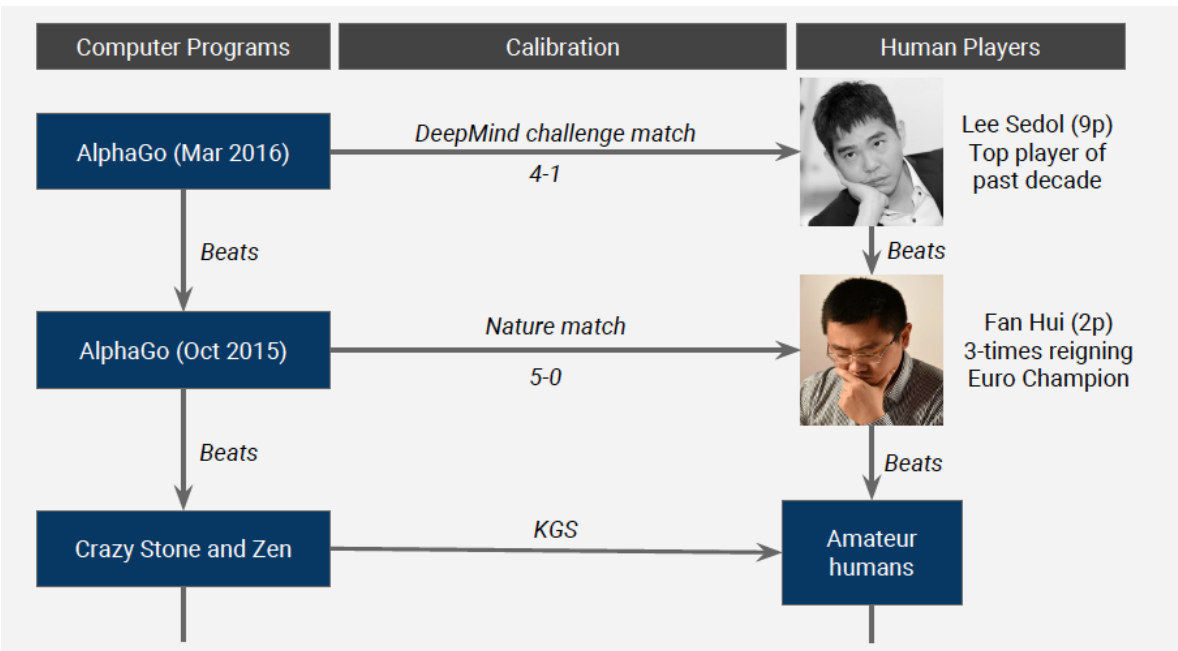


AI and Security in Cyber Physical Systems

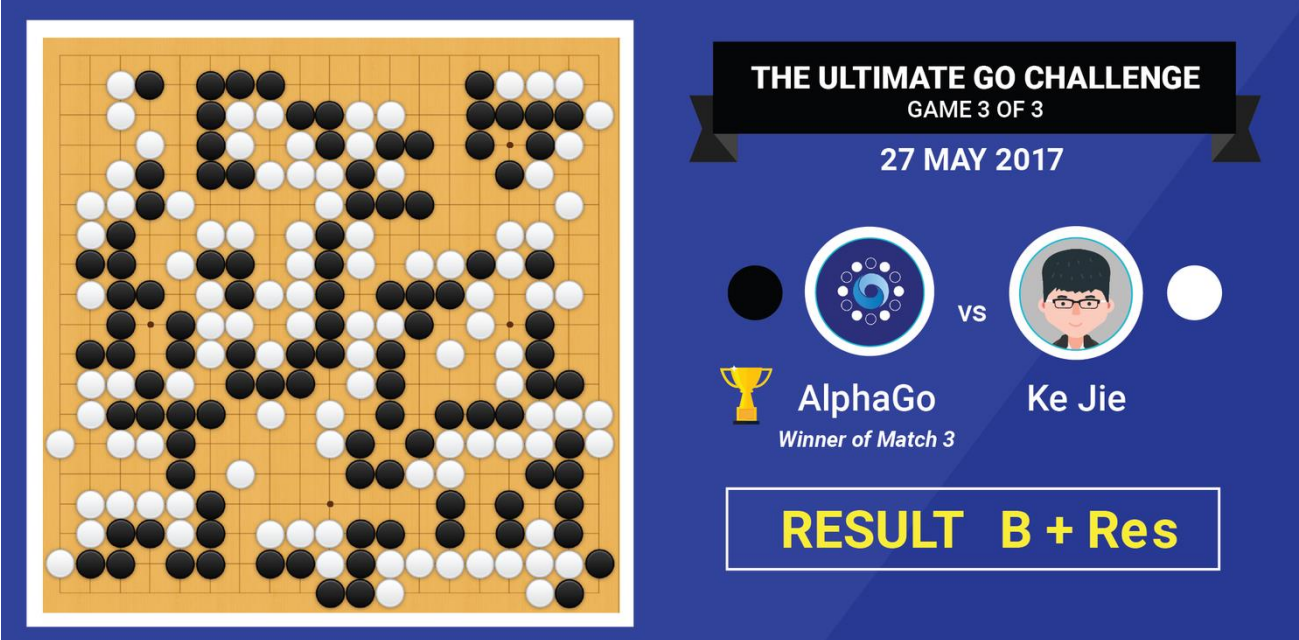
Dawn Song

UC Berkeley

AlphaGo: Winning over World Champion

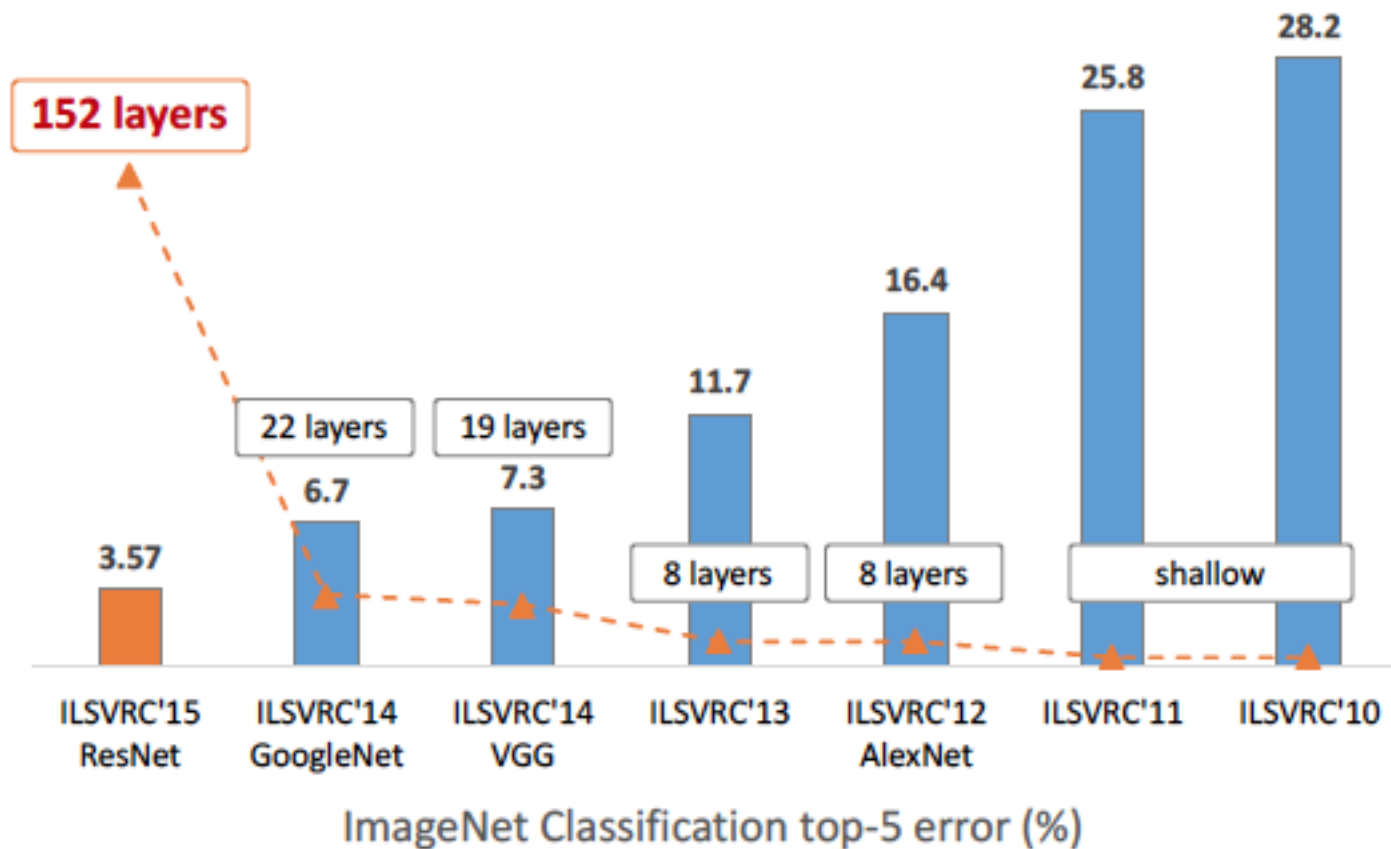


Source: David Silver



The image shows a Go board with a complex position of black and white stones. To the right is a graphic for 'THE ULTIMATE GO CHALLENGE', GAME 3 OF 3, dated 27 MAY 2017. It features the AlphaGo logo (a blue circle with white dots) and a cartoon of Ke Jie. A trophy icon is next to the text 'AlphaGo Winner of Match 3'. The result is displayed in a yellow box: 'RESULT B + Res'.

Achieving Human-Level Performance on ImageNet Classification



Deep Learning Powering Everyday Products



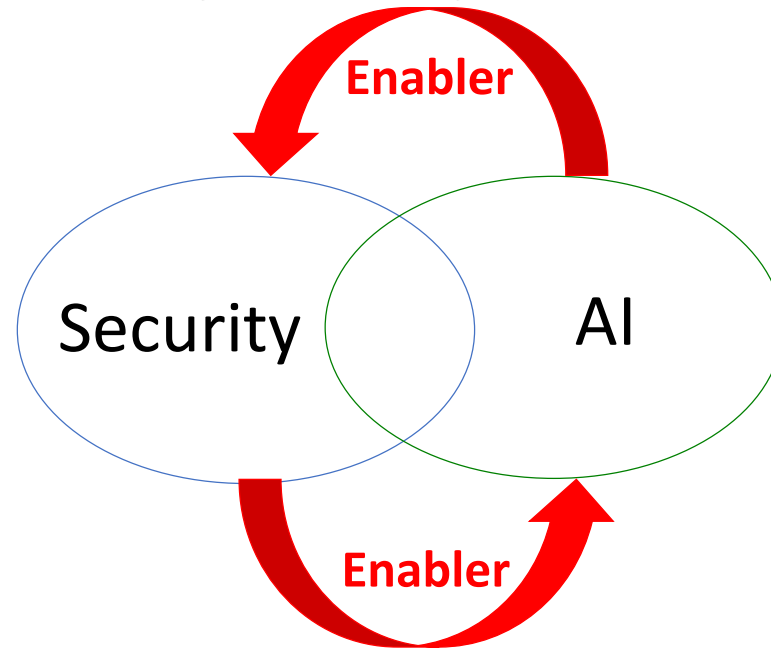
pcmag.com



theverge.com

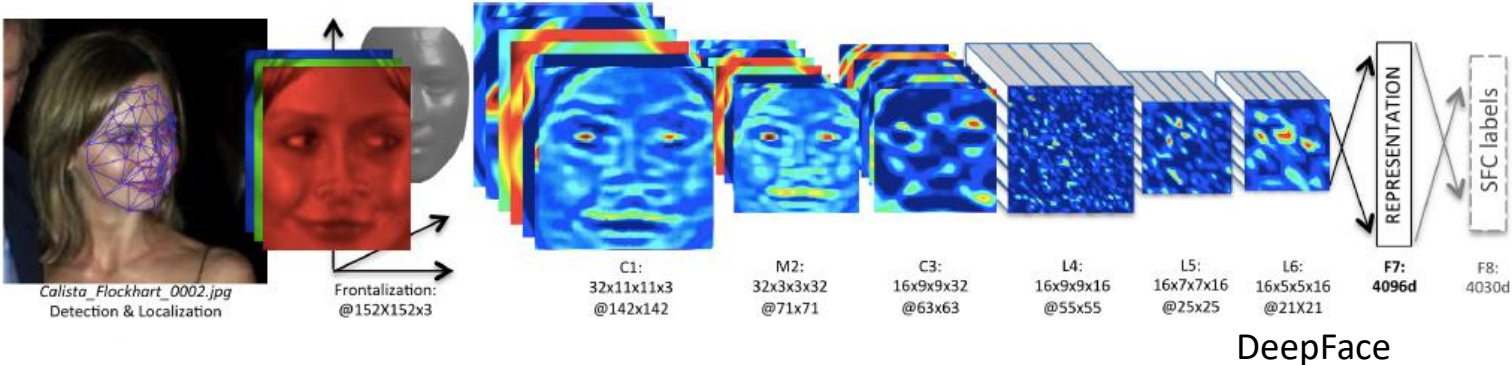


AI and Security in Cyber Physical Systems



- AI enables security applications
- 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**

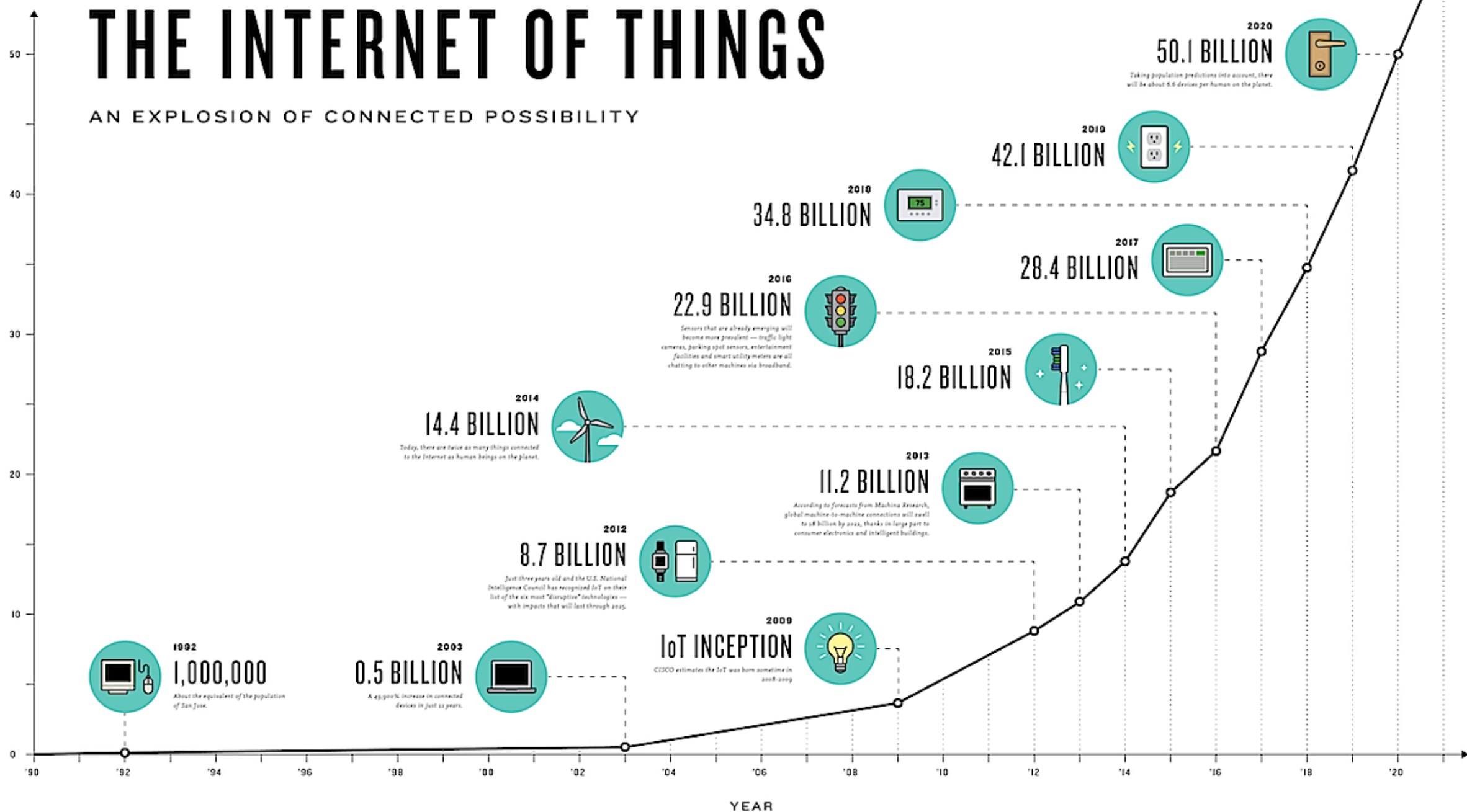
Deep Learning Improving Security Capabilities



THE INTERNET OF THINGS

AN EXPLOSION OF CONNECTED POSSIBILITY

BILLIONS OF DEVICES



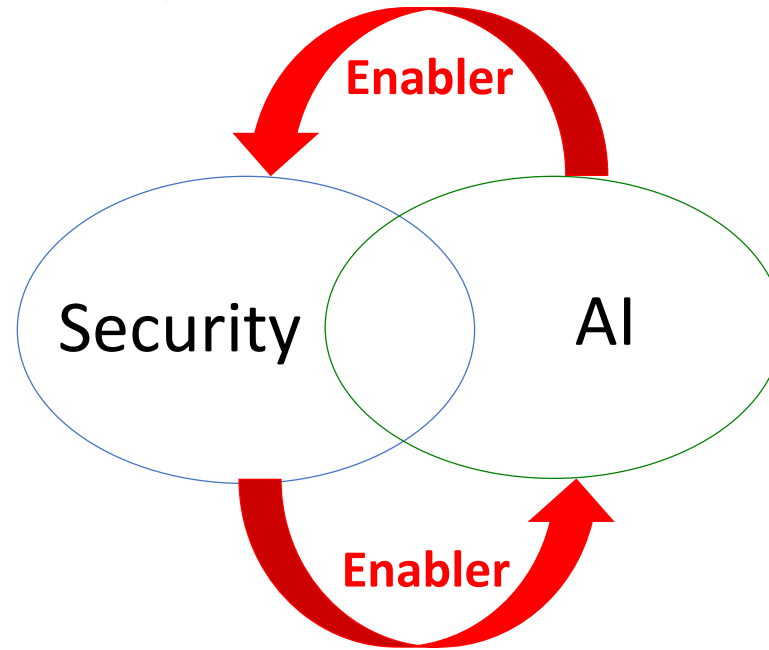
Firmware of IoT devices

- Binary code in various ISA
 - X86, MIPS, ARM, etc.
- Employ common open sourced code:
 - For example: OpenSSL
- Common vulnerability
 - Heartbleed

Deep Learning for IoT Vulnerability Detection

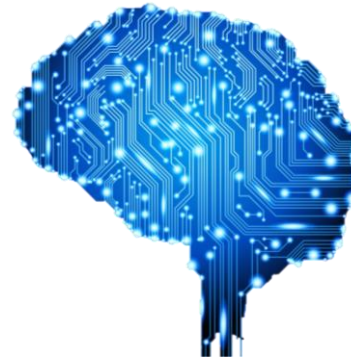
- Neural Network-based Graph Embedding for Cross-Platform Binary Code Search [XLFSY, ACM Computer and Communication Symposium 2017]
 - See talk by Chang Liu

AI and Security in Cyber Physical Systems



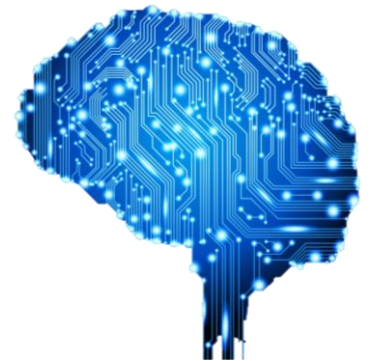
- AI enables security applications
- Security enables better AI
 - **Integrity**: produces intended/correct results (adversarial machine learning)
 - **Confidentiality/Privacy**: does not leak users' sensitive data (secure, privacy-preserving machine learning)
 - **Preventing misuse of AI**

AI and Security: AI in the presence of attacker



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 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
 - Target attacks
 - Devise attacks
 - Need security in other systems

AI and Security: AI in the presence of attacker

- Attack AI:
 - Cause 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
 - Target attacks
 - Devise attacks
 - Need security in other systems

Deep Learning Systems Are Easily Fooled



$$\frac{\partial \text{output}}{\partial \text{pixels}}$$

← ostrich →

Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., & Fergus, R. Intriguing properties of neural networks. ICLR 2014.

Adversarial examples fooling autonomous vehicles



Misclassified as
Speed Limit 60



Misclassified as
Speed Limit 75



Misclassified as
Speed Limit 75

Adversarial Examples in Physical World robust against viewpoint changes



Subtle Perturbations

Robust Physical-World Attacks on Machine Learning Models [EEFKLPRS, 2017]

Adversarial Examples in Physical World robust against viewpoint changes



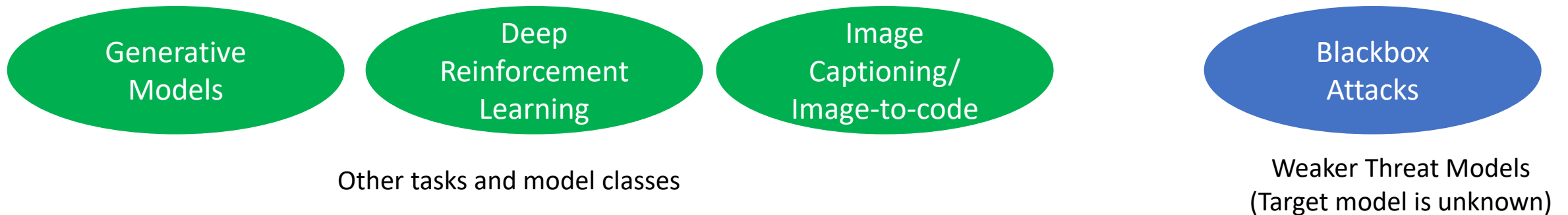
Camouflage Perturbations

Robust Physical-World Attacks on Machine Learning Models [EEFKLPRS, 2017]

Adversarial Examples Prevalent in Deep Learning Systems

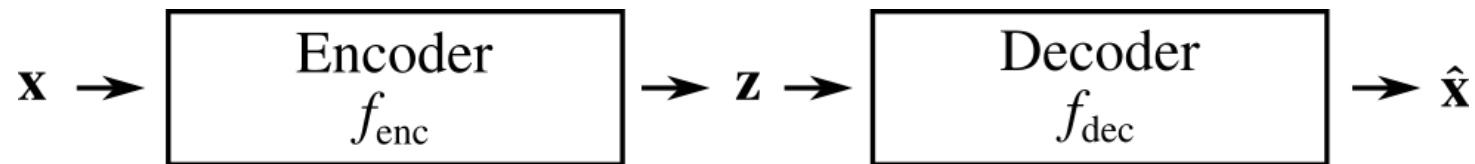
- Most existing work on adversarial examples:
 - Image classification task
 - Target model is known

- Our investigation on adversarial examples:



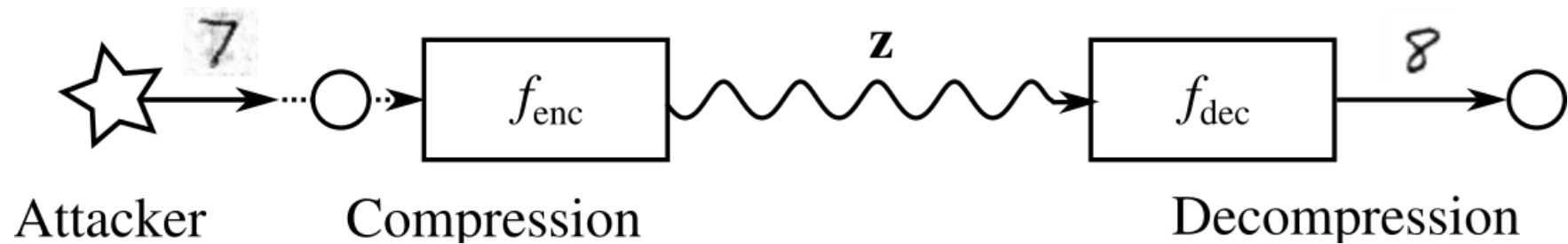
Generative models

- 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

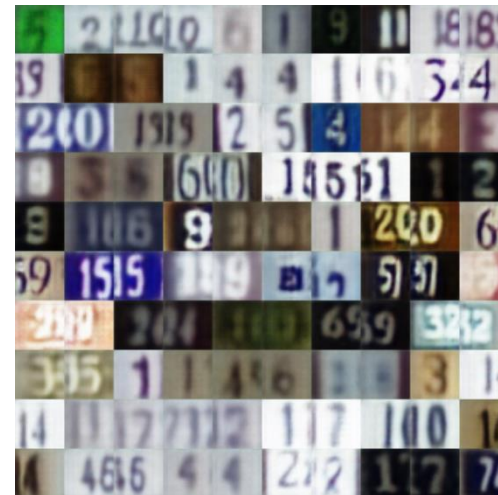
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0
0 0 0 0 0 0 0 0 0 0

Reconstruction of adversarial examples

Adversarial Examples for VAE-GAN in SVHN

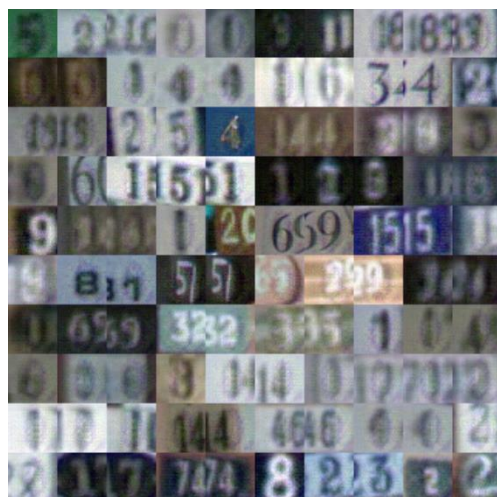
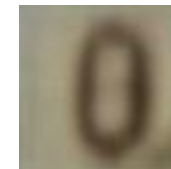


Original images



Reconstruction of original images

Target Image

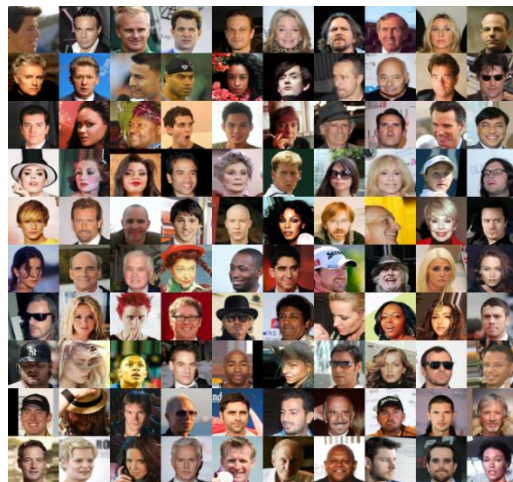


Adversarial examples



Reconstruction of adversarial examples

Adversarial Examples for VAE-GAN in SVHN



Original images

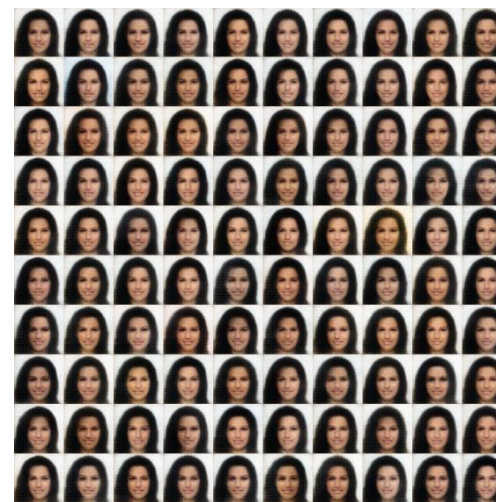


Reconstruction of original images

Target Image

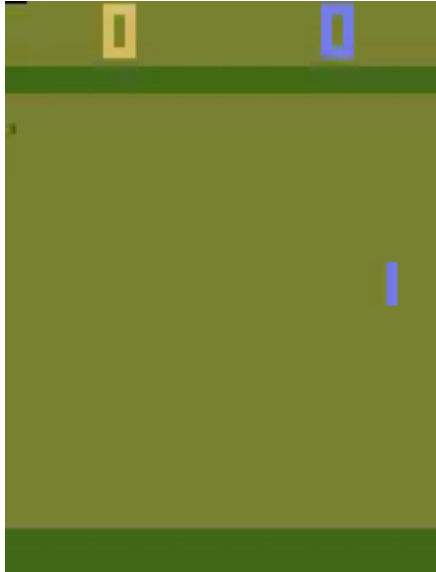


Adversarial examples



Reconstruction of adversarial examples

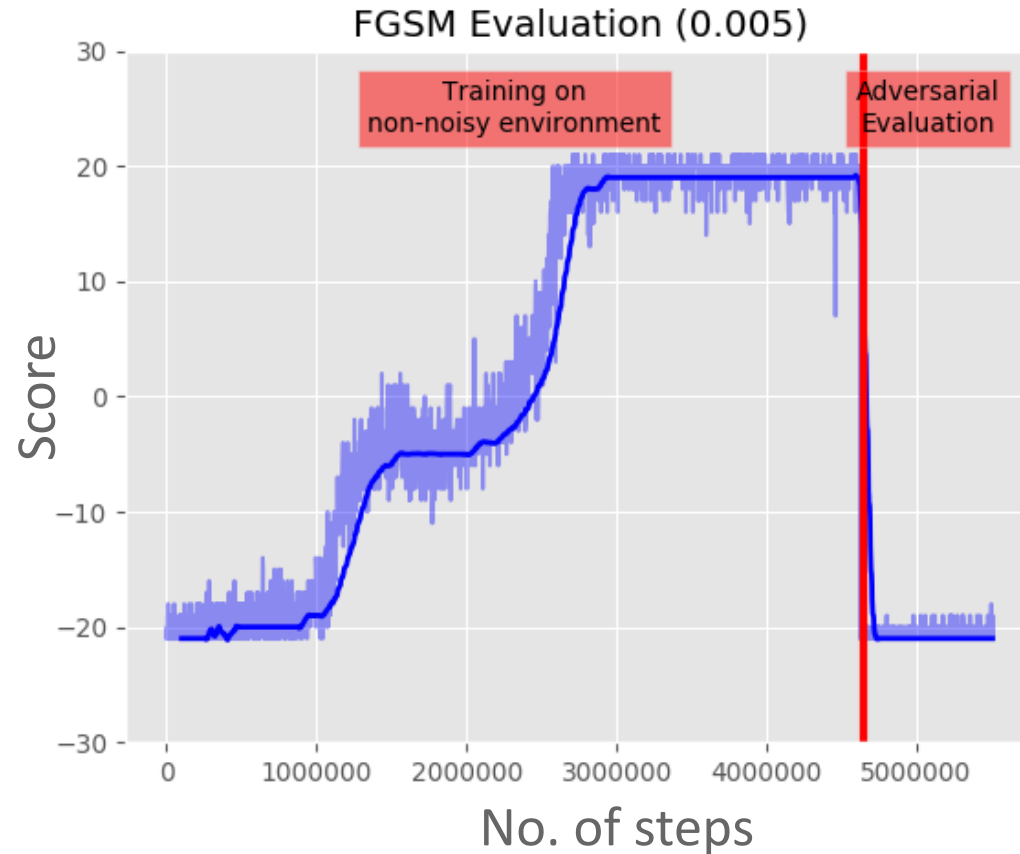
Deep Reinforcement Learning Agent (A3C) Playing Pong



Original Frames

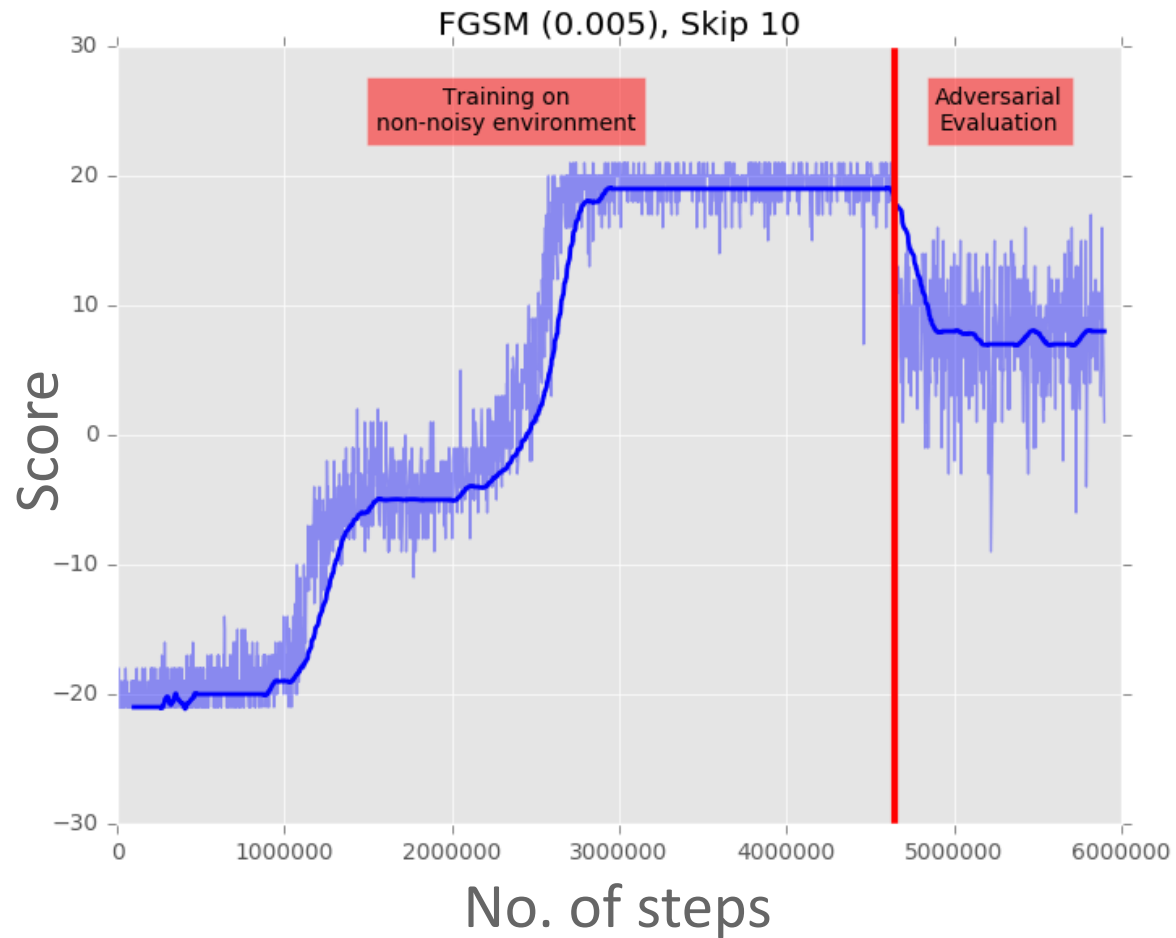
Jernej Kos and Dawn Song: Delving into adversarial attacks on deep policies [ICLR Workshop 2017].

Adversarial Examples on A3C Agent on Pong

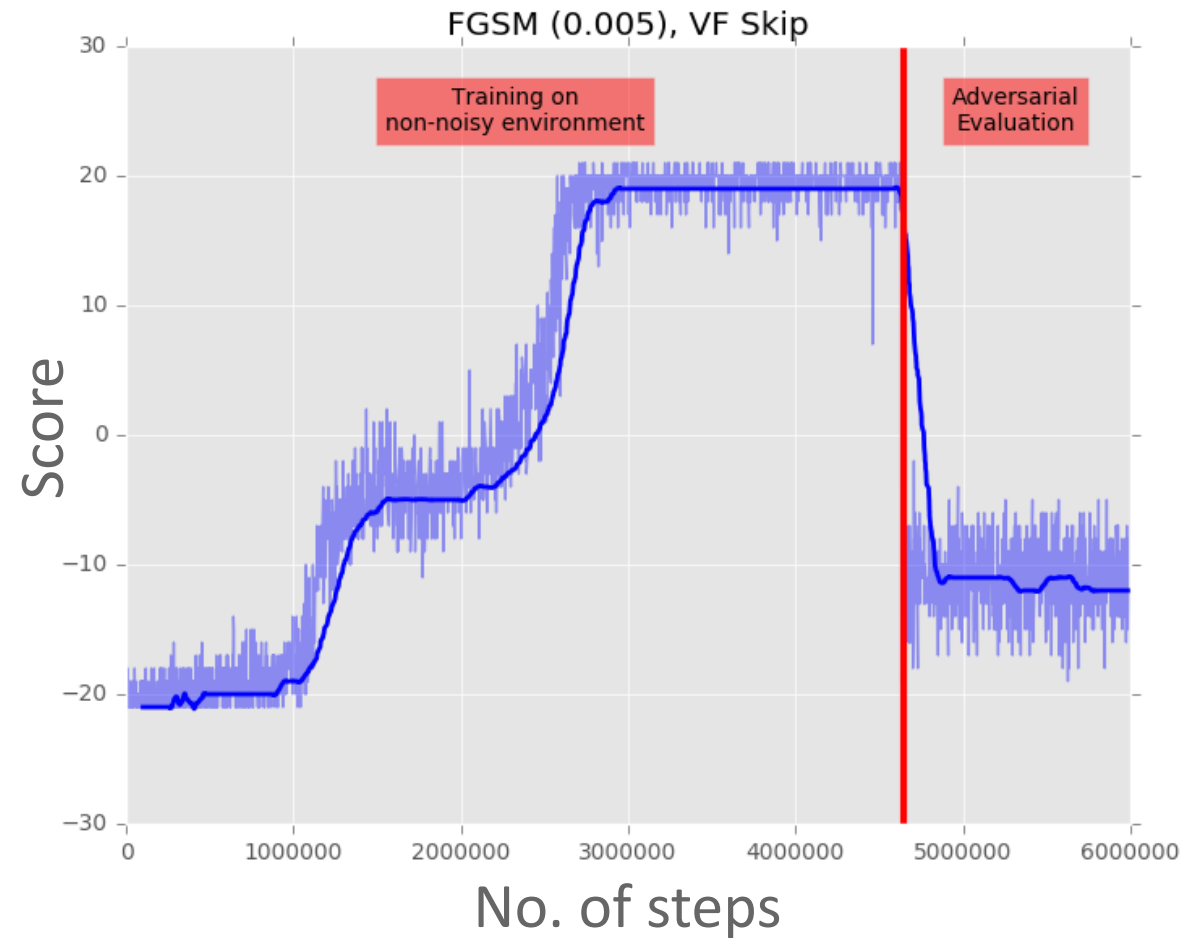


Jernej Kos and Dawn Song: Delving into adversarial attacks on deep policies [ICLR Workshop, 2017]

Attacks Guided by Value Function

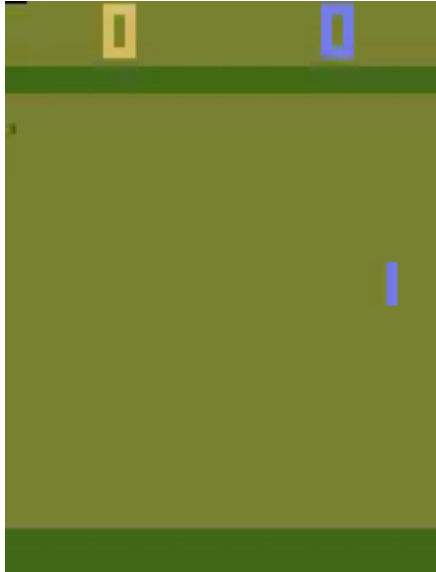


Blindly injecting adversarial perturbations every 10 frames.

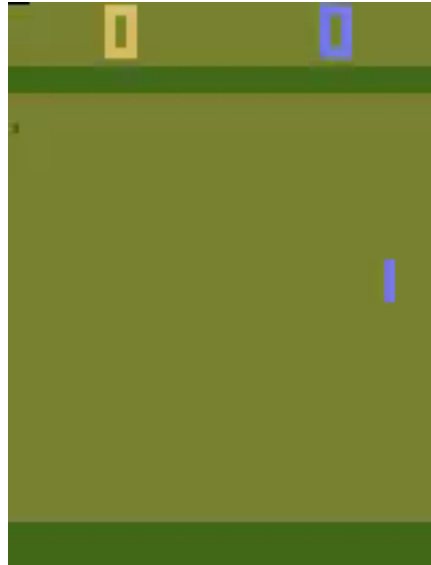


Injecting adversarial perturbations guided by the value function.

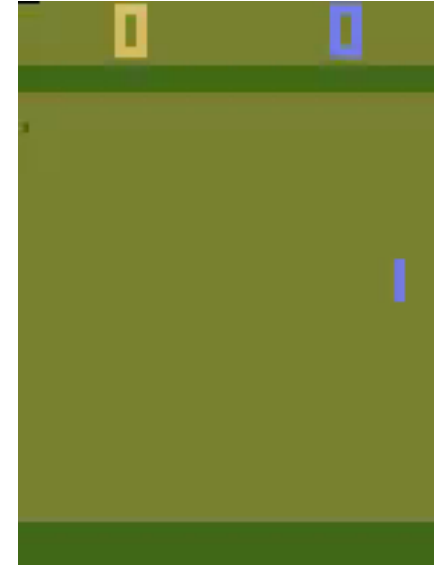
Agent in Action



Original Frames



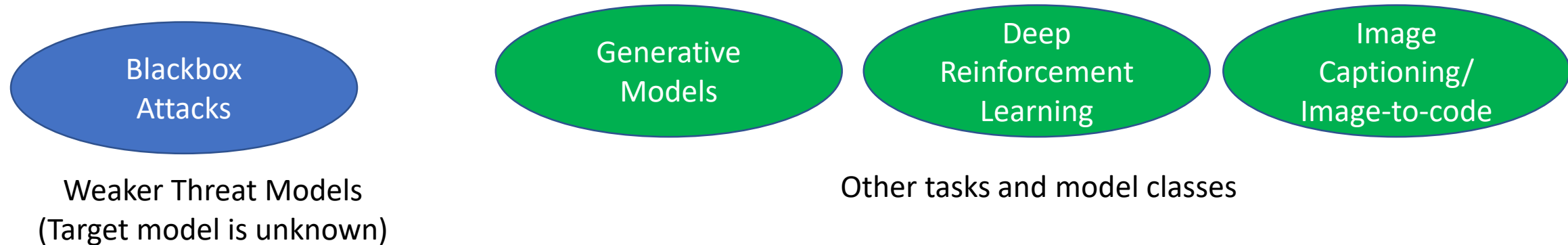
With FGSM perturbations
($\epsilon = 0.005$) inject in
every frame



With FGSM perturbations
($\epsilon = 0.005$) inject based
on value function

Adversarial Examples Prevalent in Deep Learning Systems

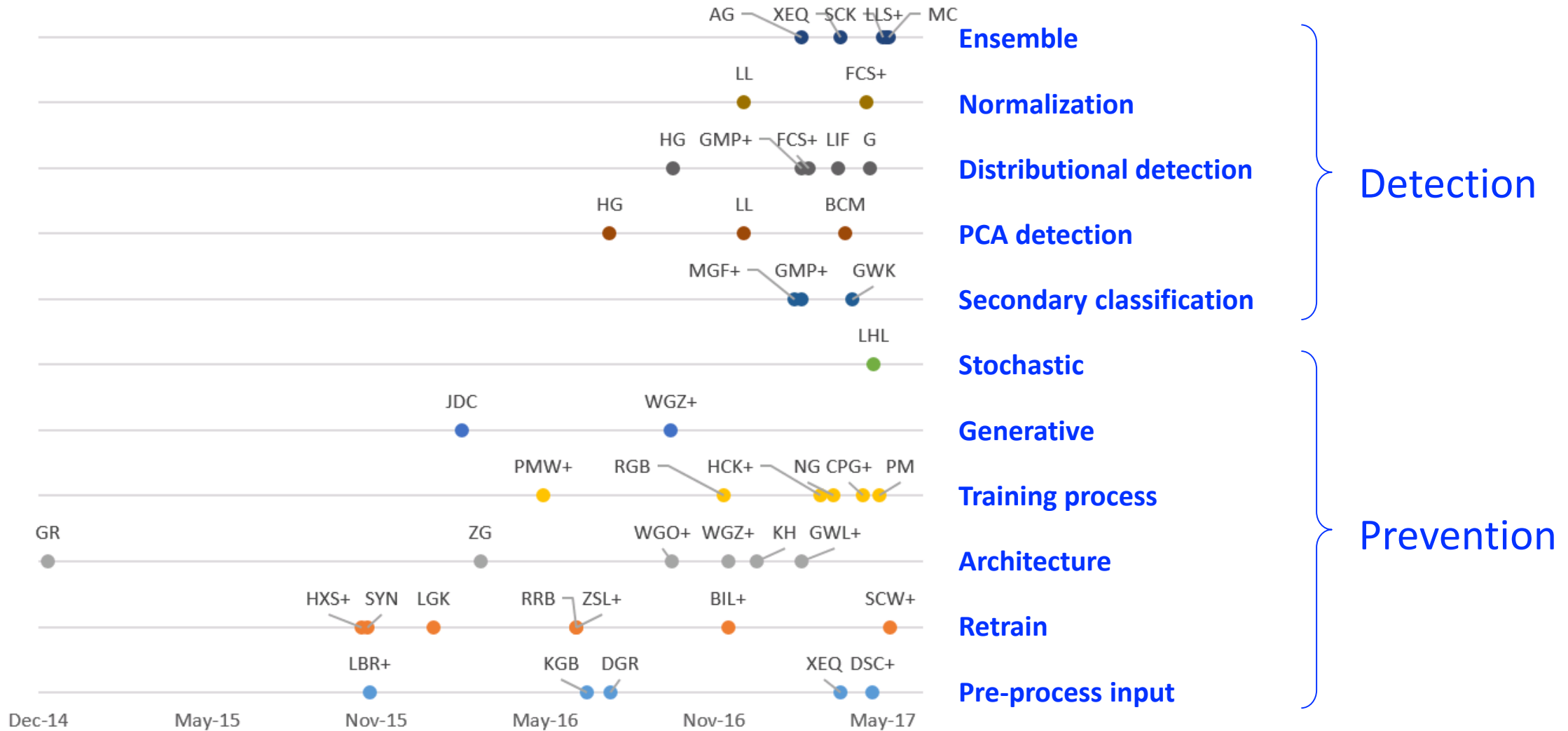
- Most existing work on adversarial examples:
 - Image classification task
 - Target model is known
- Our investigation on adversarial examples:



Numerous Defenses Proposed

- Input processing
 - Gaussian blur, median blur
 - Quantization
- Adversary re-training
 - Re-train with generated adversarial examples
- Detecting adversarial examples
 - Detecting anomalous high-frequency patterns in input
 - Detecting anomalous activations
 - Detecting low confidence output

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 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

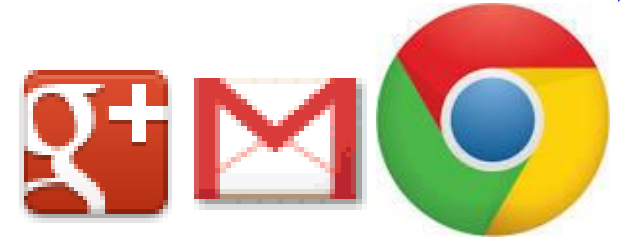
Reactive Defense

Proactive Defense:
Bug Finding

Proactive Defense:
Secure by Construction

Automatic worm detection
& signature/patch generation

Automatic malware
detection & analysis



Progression of my approach 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



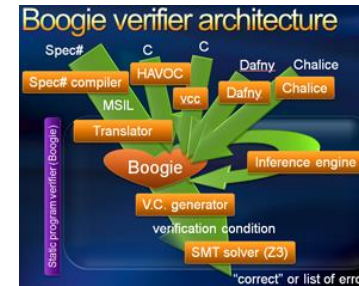
Coq



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:



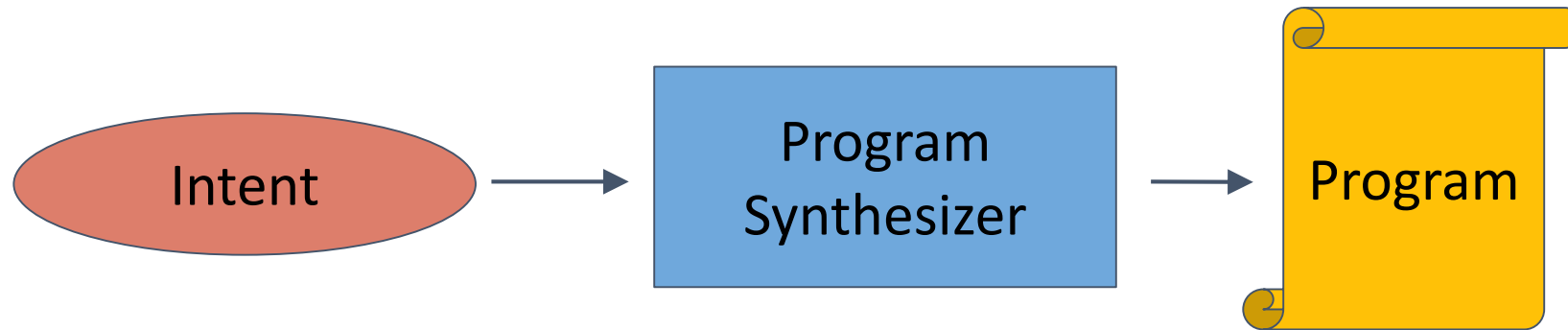
- No precisely specified properties & goals
- No good understanding of how learning system works
- Traditional symbolic reasoning techniques do not apply

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

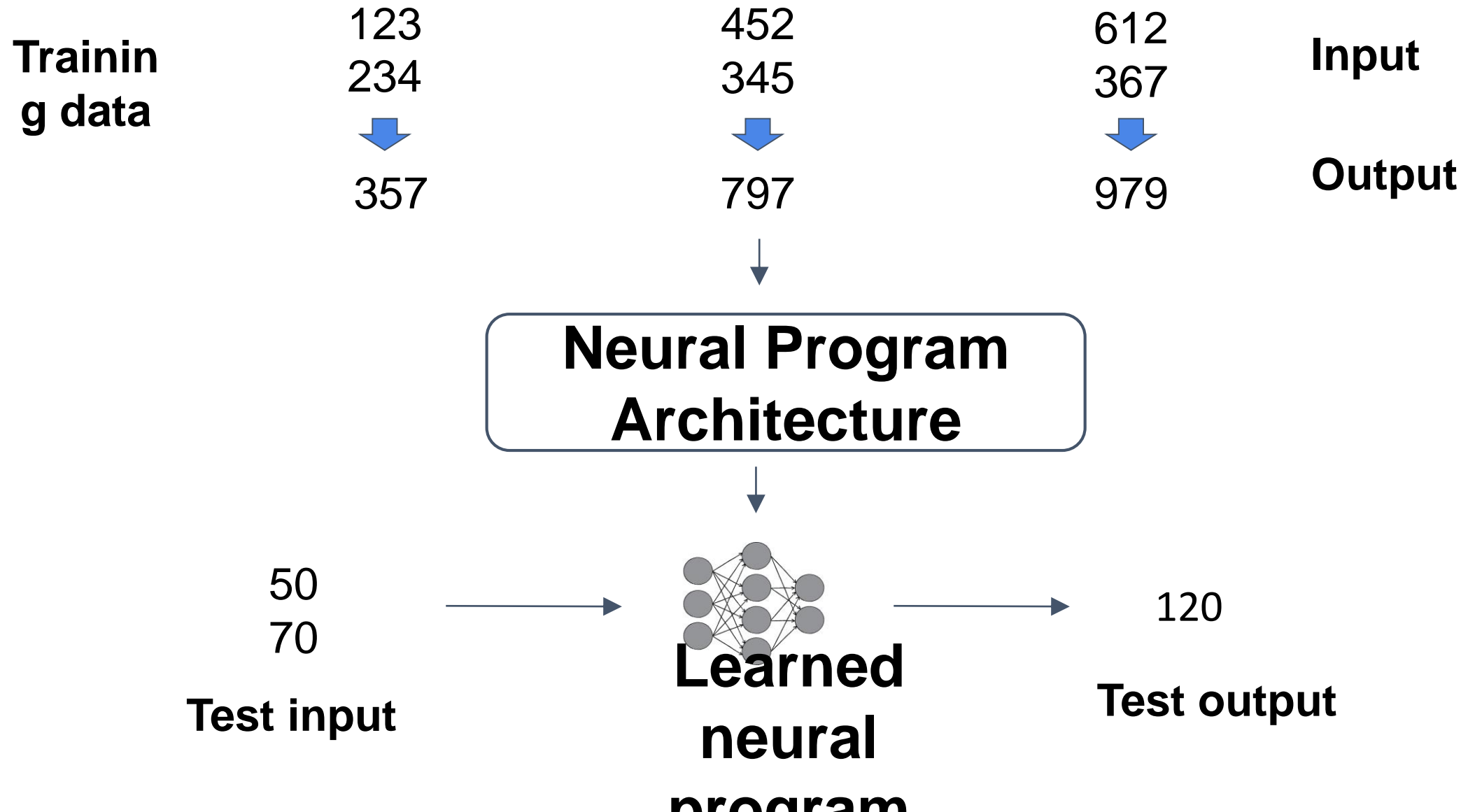
Program synthesis can automate this & democratize idea realization



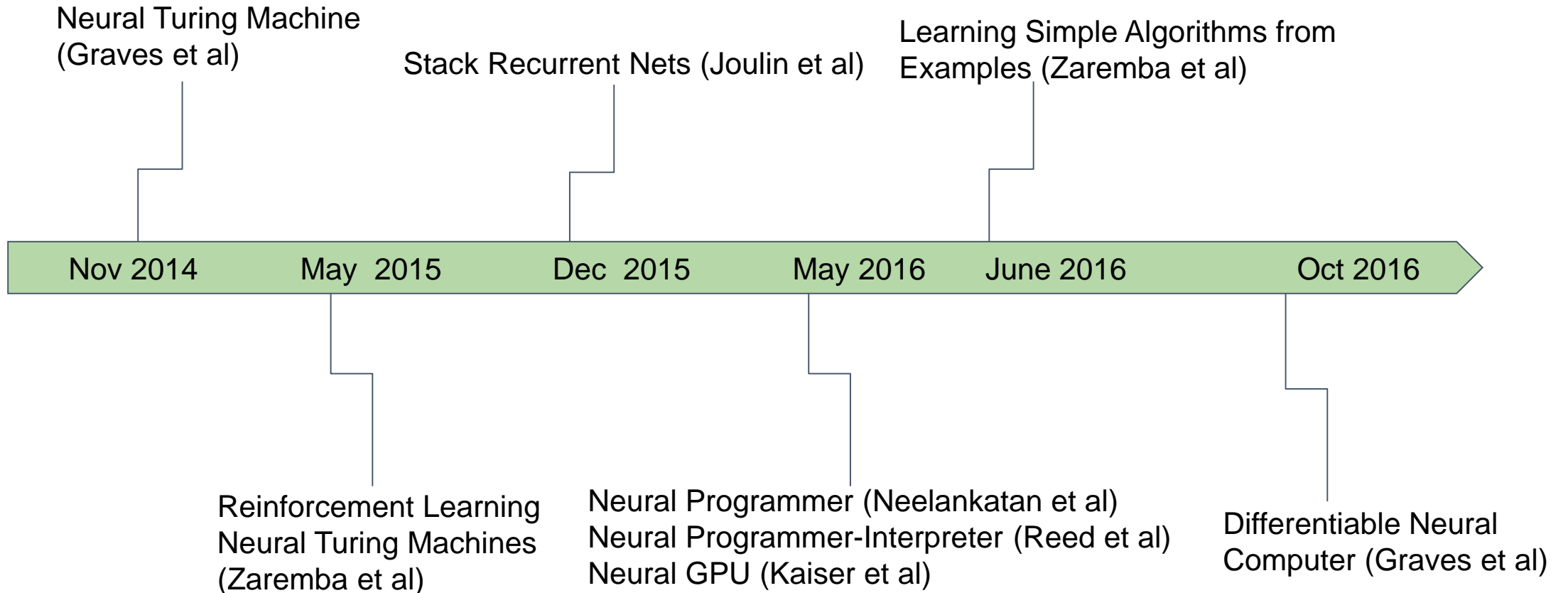
Neural Program Synthesis

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

Neural Program Synthesis

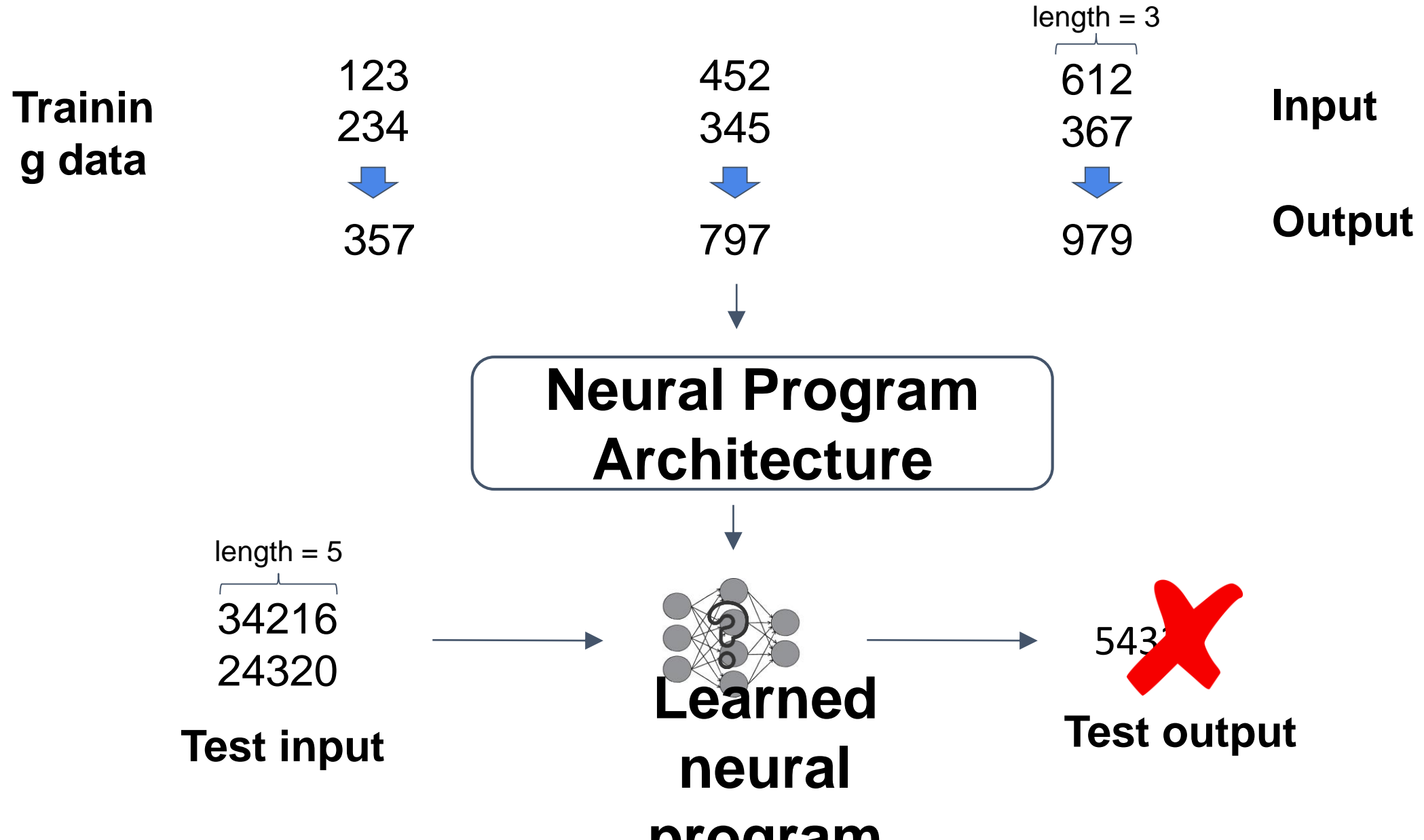


Neural Program Architectures

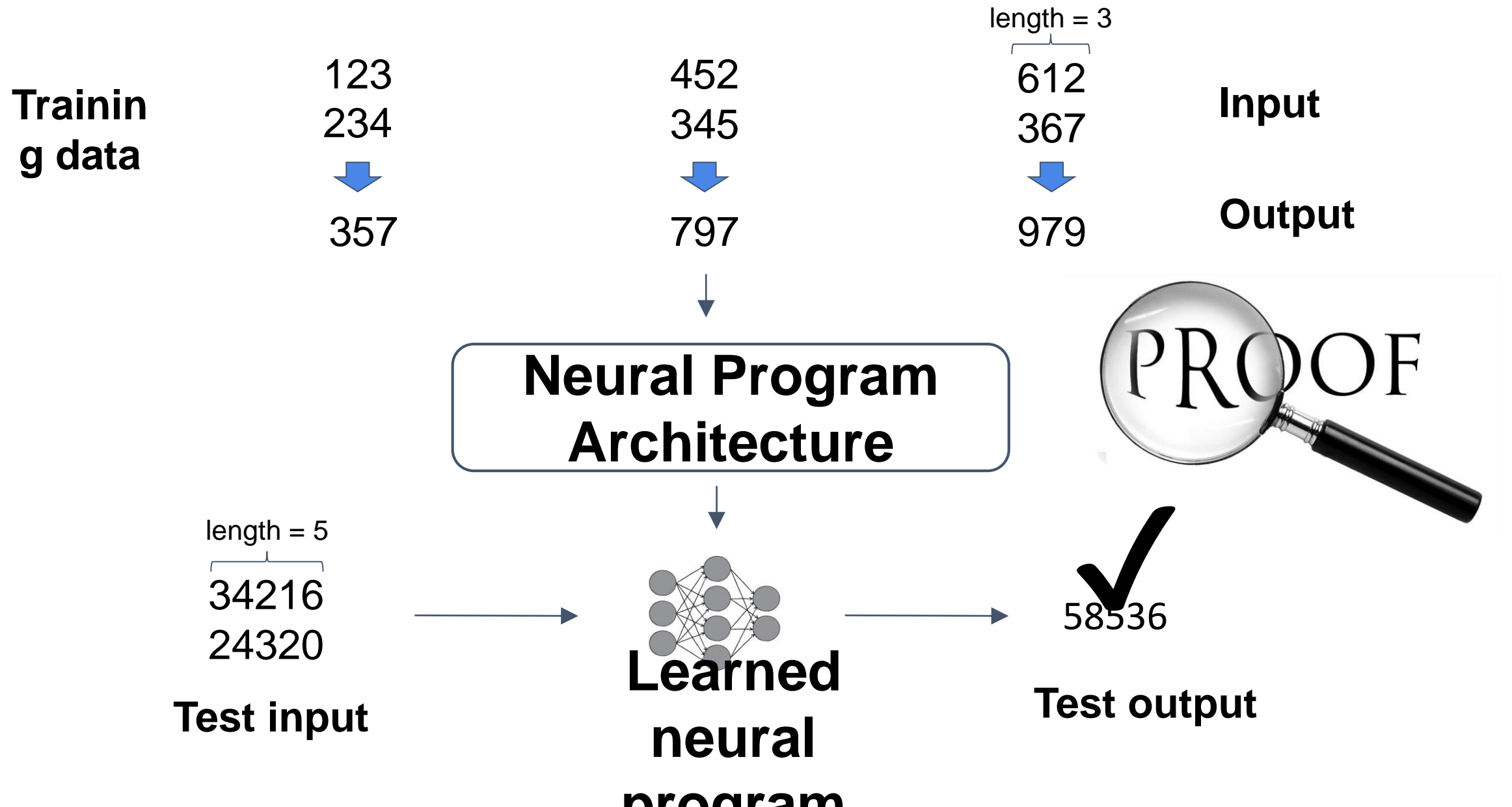


Neural Program Synthesis Tasks: Copy, Grade-school addition, Sorting, Shortest Path

Challenge 1: Generalization



Challenge 2: No Proof of Generalization

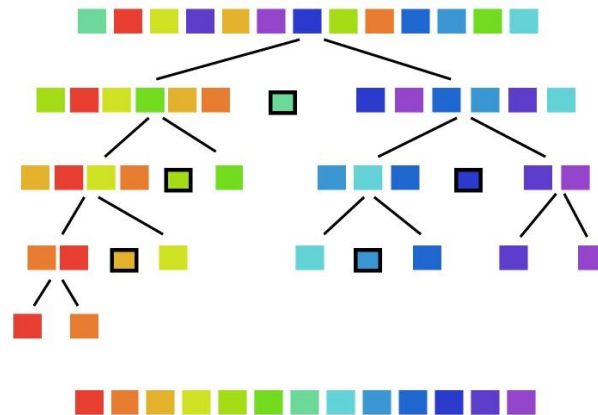


Our Approach: Introduce Recursion

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



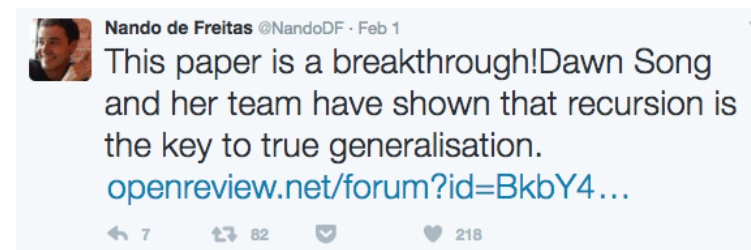
Quicksort

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

| <u>Length of Array</u> | <u>Non-Recursive</u> | <u>Recursive</u> |
|------------------------|----------------------|------------------|
| 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% |



Importance of Recursion in Neural Program Architectures

- We introduce recursion, for the first time, into neural program architectures, and learn recursive neural programs
- We address two main challenges using recursion:
 - Generalization to more complex inputs
 - Proof of generalization

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

Misused AI can make attacks more effective



Deep Learning Empowered
Bug Finding



Deep Learning Empowered
Phishing Attacks

Misused AI for large-scale, automated, targeted manipulation



Female
25-35 Years old
AMEX User



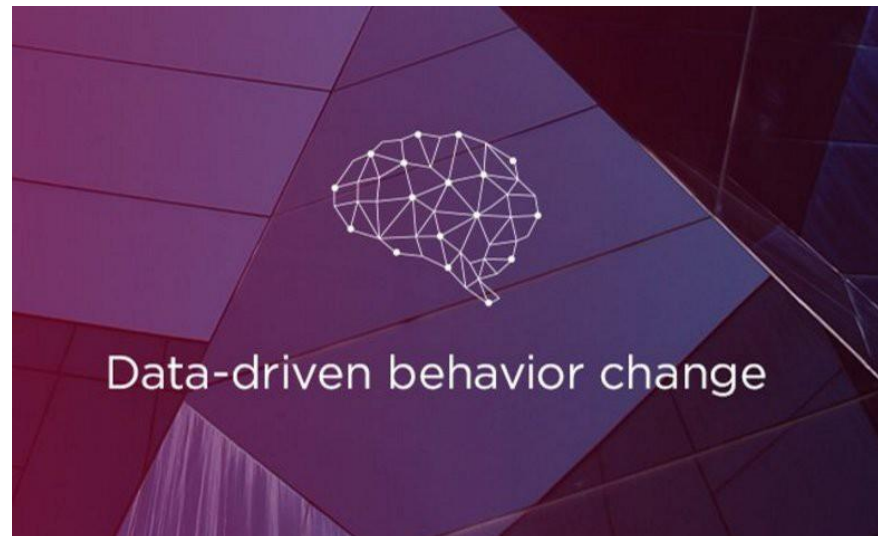
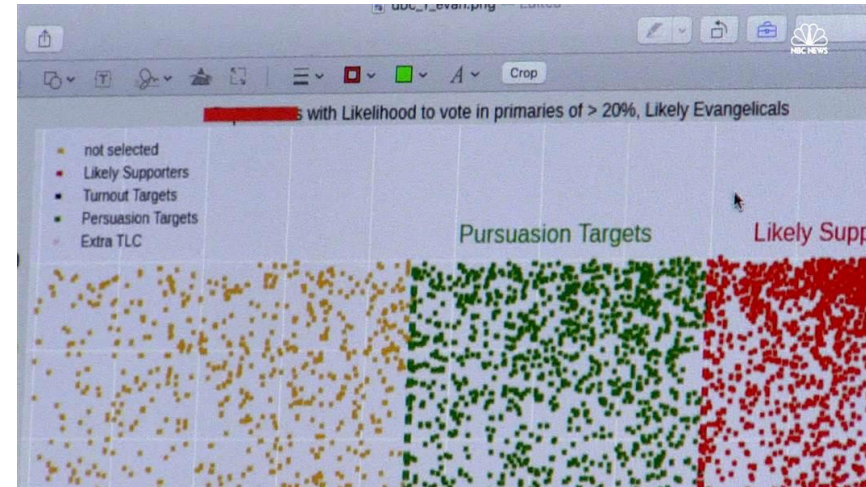
People with high openness and extraversion love new experiences they can share with lots of people.



Female
25-35 Years old
AMEX User



People with low openness and extraversion really value down time spent with their closest friends.



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



