



Introduction to Secure Machine Learning

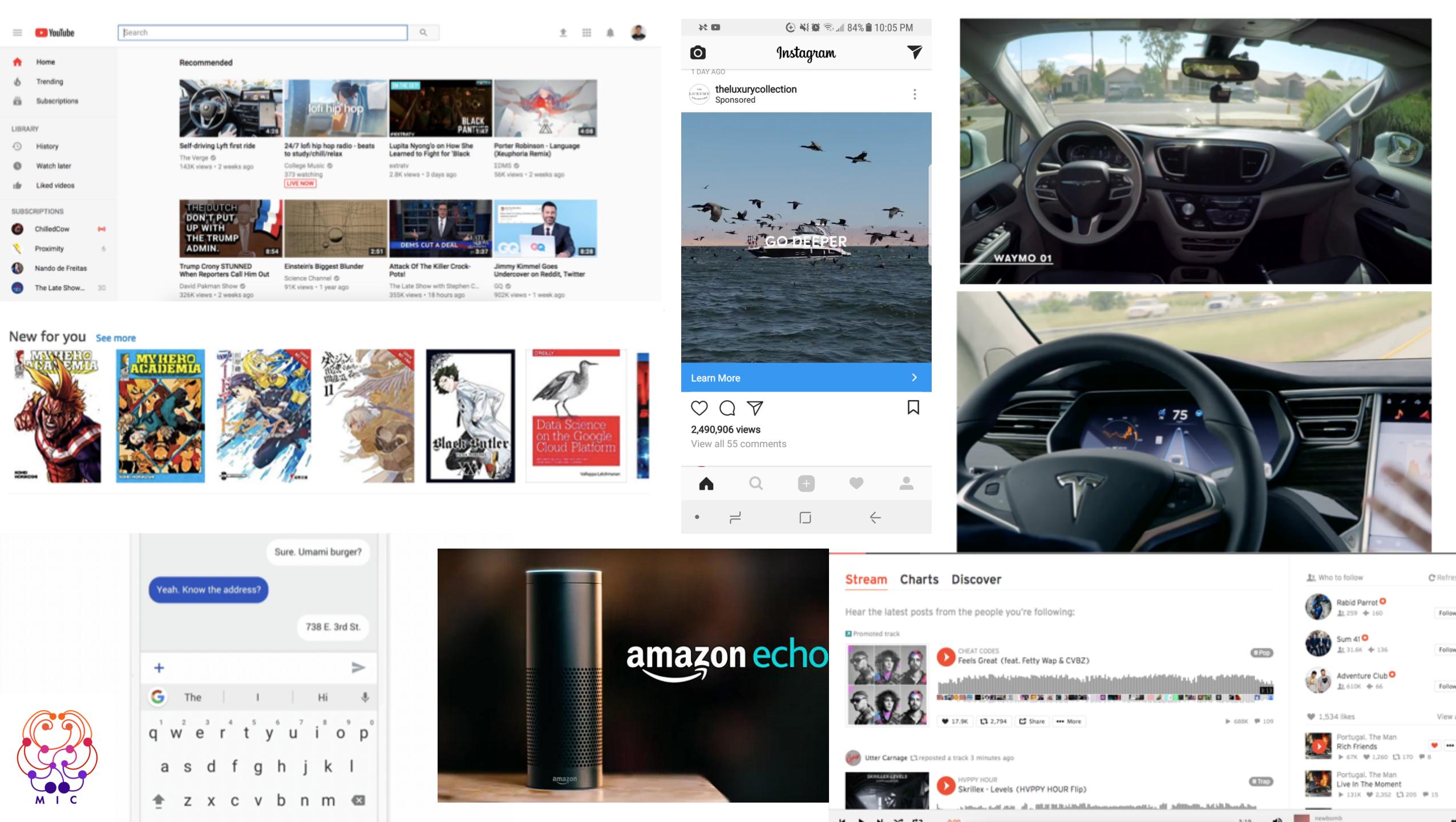
BOSTON UNIVERSITY
**MACHINE INTELLIGENCE
COMMUNITY**

Justin Chen
Jan. 29, 2018

Age of Machine Intelligence

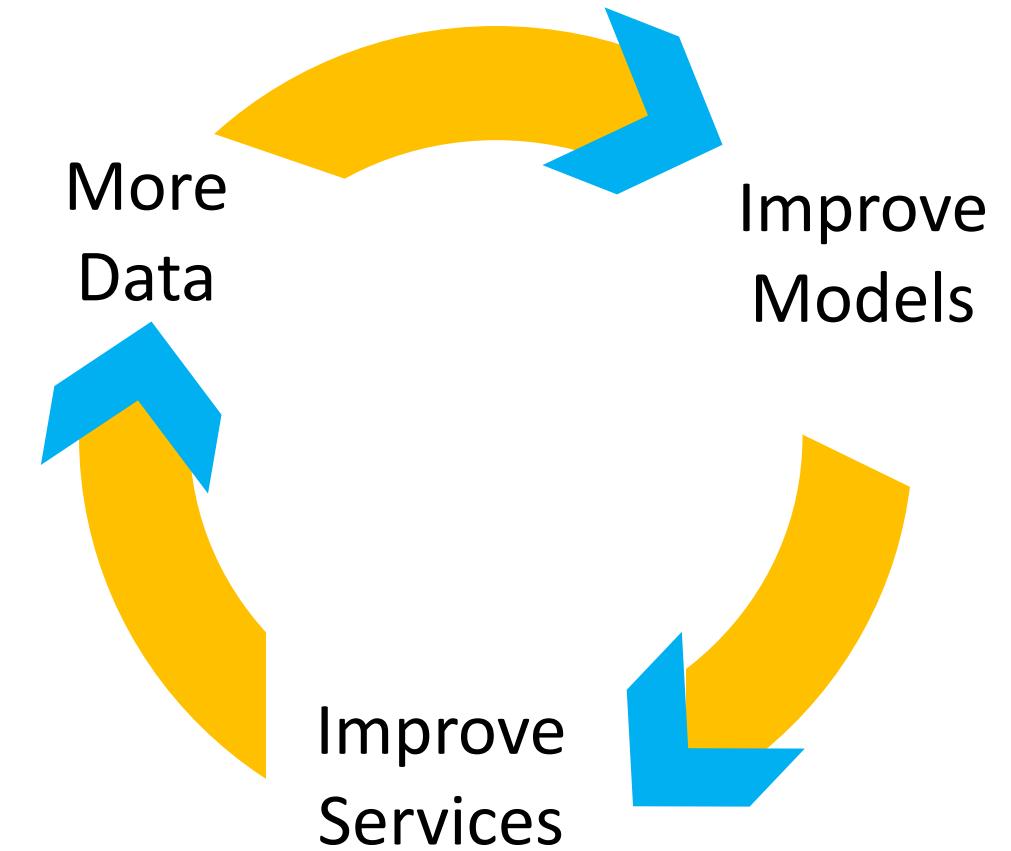
Society, Data, and Privacy





Virtuous Cycle of Machine Learning

- Term coined by Andrew Ng
- Collect data to improve machine learning models
- Use improved models to improve services
- Use improved services to collect more data



AI on the Blockchain



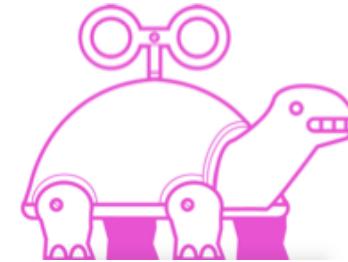
OpenMined

An open-source community focused on building technology to facilitate the decentralized ownership of data and intelligence.

<https://openmined.org/>

<https://openmined.slack.com>

<https://github.com/OpenMined>



Synapse AI

Decentralized Data and AI Marketplace

<https://synapse.ai/>

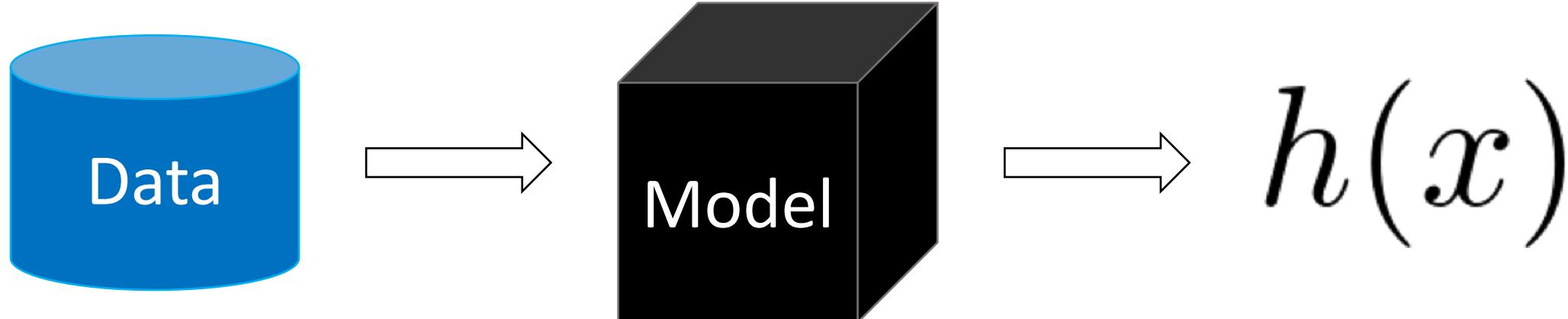


Learning from Data

Training, Validation, Testing



Function Approximation



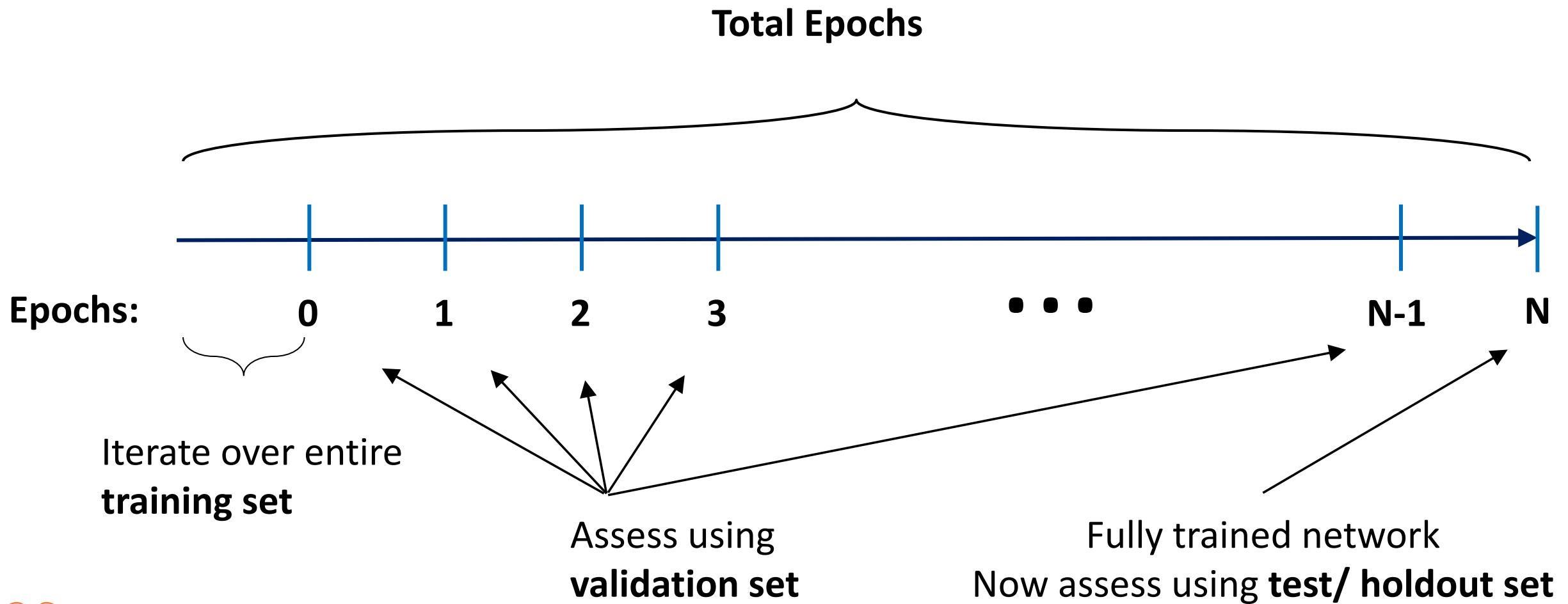
Images
Audio
Video
Text
Logs
Spreadsheets

Linear Regression
Polynomial Regression
Logistic Regression
Multinomial Regression
Neural Network
Gaussian Process
...

Class of objects
Real-value
Generated Data
Action

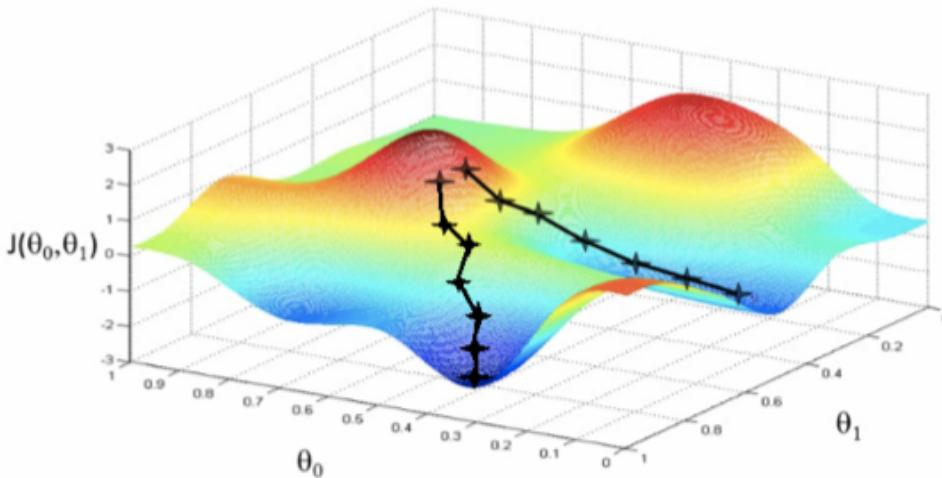


Training, Validation, Testing

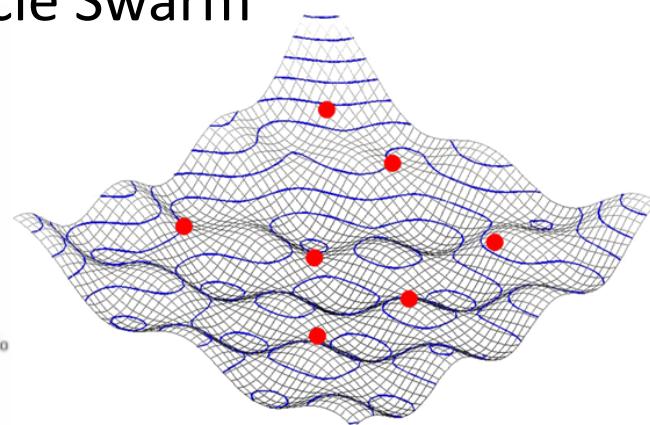


Training

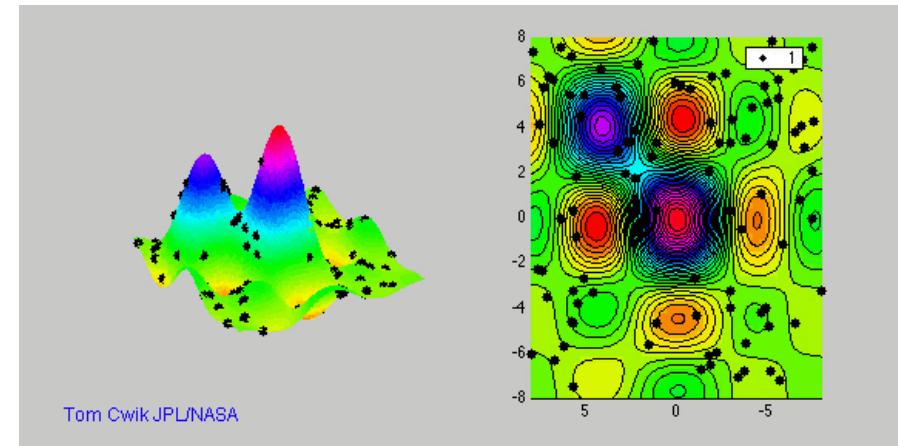
Gradient-based



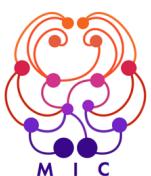
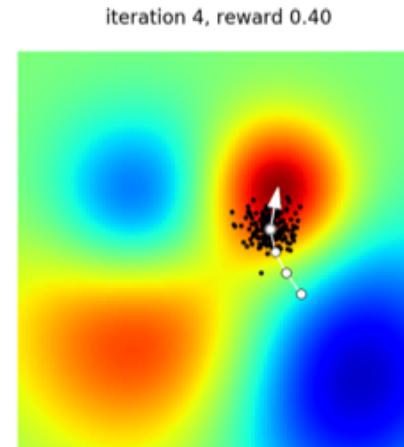
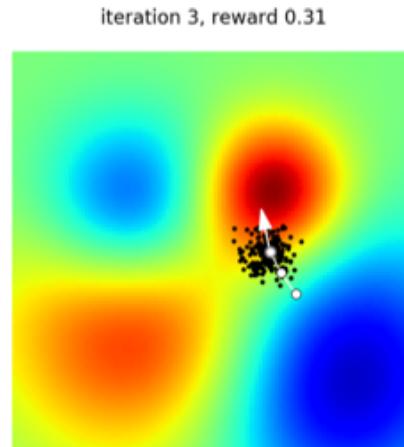
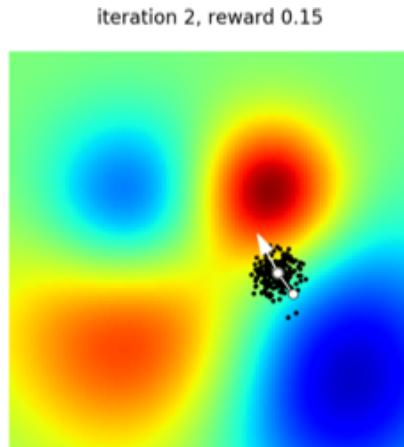
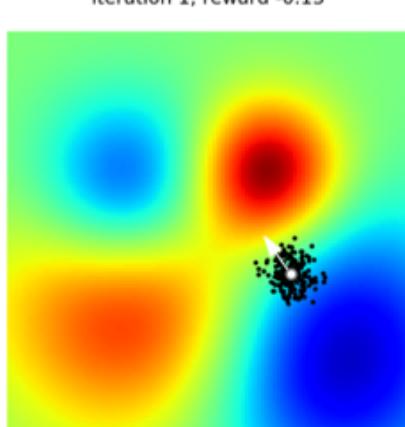
Particle Swarm



Genetics



Evolutionary
Strategies



Gradient Descent

$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

Scalar learning rate

Individual weights

Vector of weights

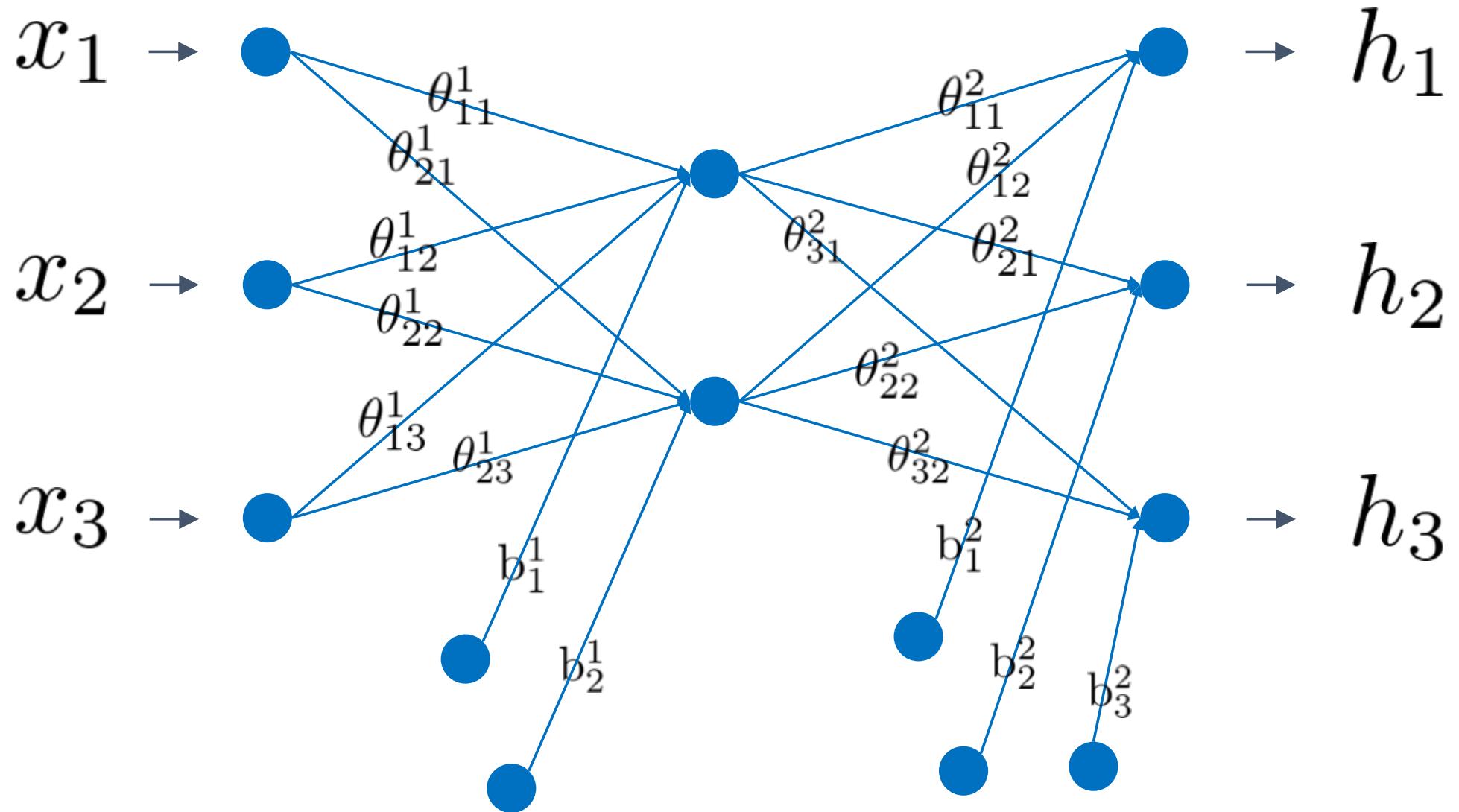
Cost/objective/loss function

Also written $\nabla J(\theta)$

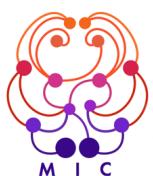
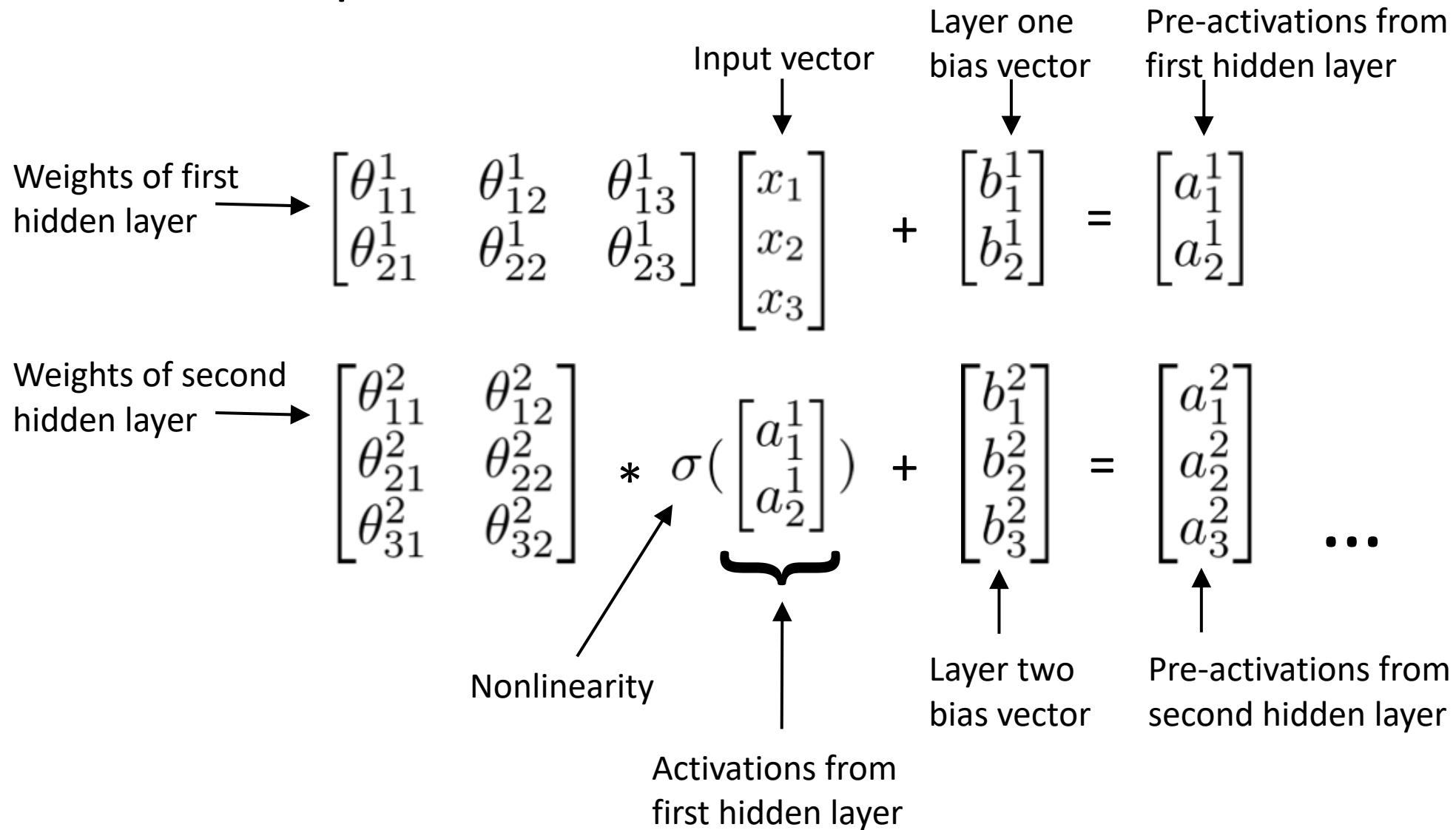
The diagram illustrates the gradient descent update rule. It shows the formula $\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$. Annotations with arrows explain the components: 'Scalar learning rate' points to the scalar α ; 'Individual weights' points to the individual weight θ_j ; 'Vector of weights' points to the vector θ ; 'Cost/objective/loss function' points to the function $J(\theta)$; and 'Also written' points to the alternative form $\nabla J(\theta)$.



Artificial Neural Network



Matrix Representation



On the Protection of Private Information in Machine Learning Systems: Two Recent Approaches

Martín Abadi, Úlfar Erlingsson, Ian Goodfellow, H. Brendan McMahan, Ilya Mironov, Nicolas Papernot, Kunal Talwar, and Li Zhang

Google



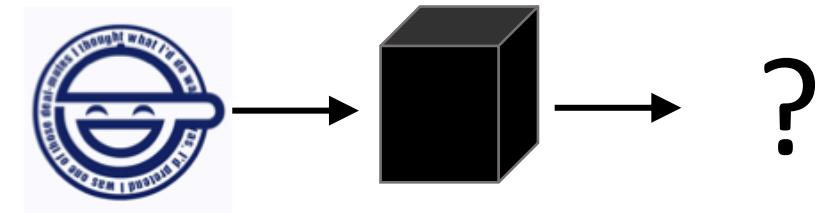
Attack Surface

Adversary's Abilities

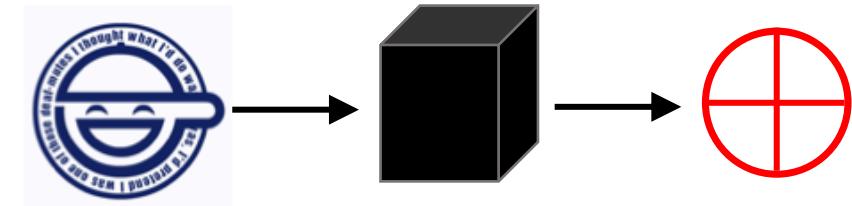


Types of Attacks

Untargeted



Targeted

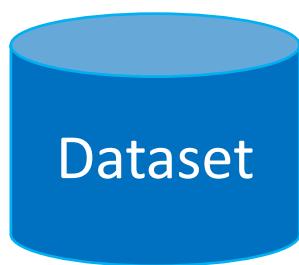


Economy of Mechanism

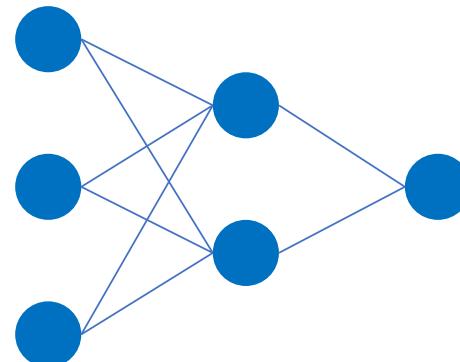
Protection mechanisms should be made as simple as possible

Machine learning algorithms have many moving parts

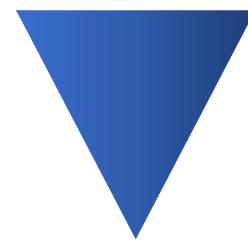
How much information does each component leak?



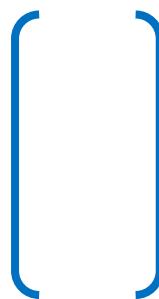
Inputs



Parameters



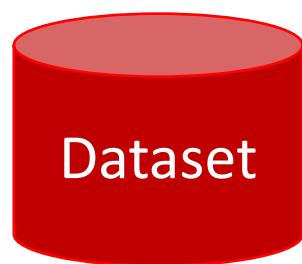
Gradients



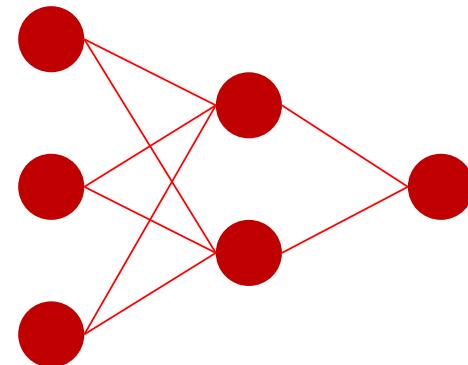
Outputs

Fail-safe Defaults

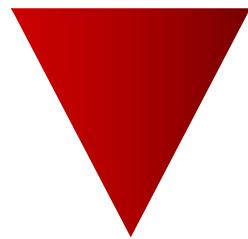
The default behavior of a security system should
be to completely deny access



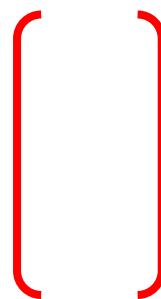
Dataset



Parameters



Gradients



Outputs



Inputs

Complete Mediation

Access to sensitive data must always be secure

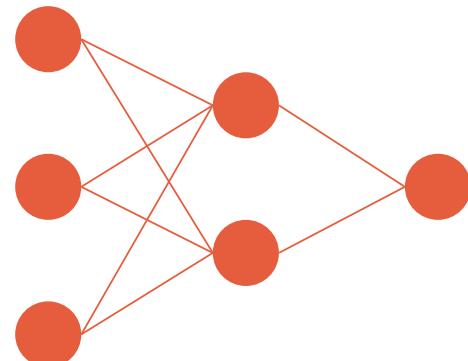
Transferring data over networks

Training and testing data

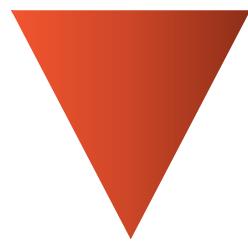
Does not apply to white-box setting



Dataset



Parameters



Gradients

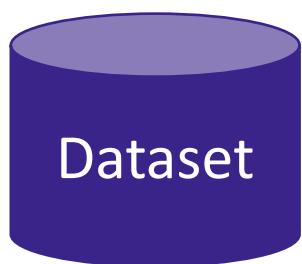


Outputs

Open Design

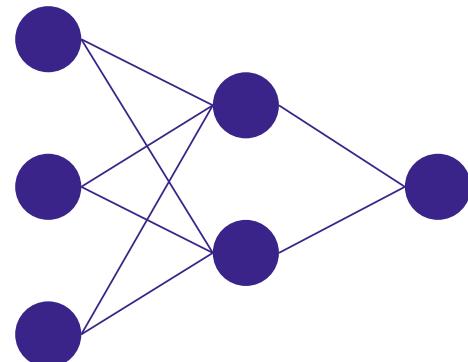
Kerckhoff's principle

Protection mechanism should not depend on secrecy and system design should be made public

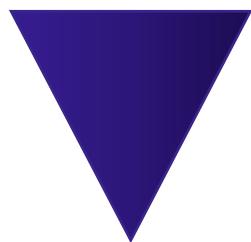


Dataset

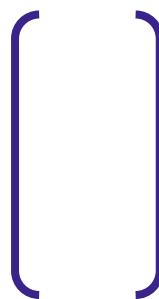
Inputs



Parameters



Gradients

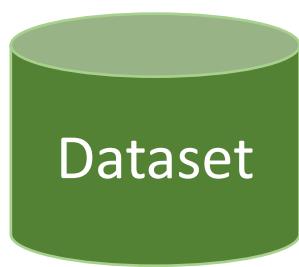


Outputs

Separation of Privilege

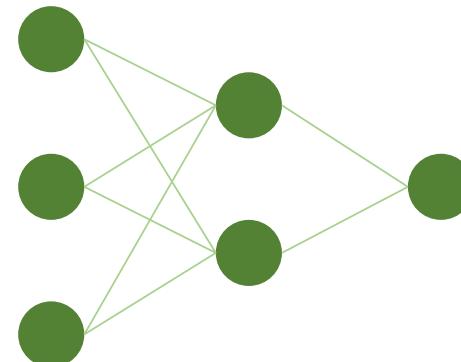
Multiple independent keys should be required for controlling access

Unclear how this would apply outside of the distributed setting

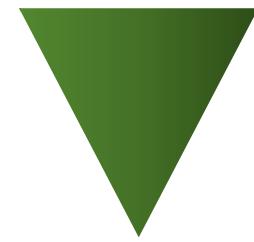


Dataset

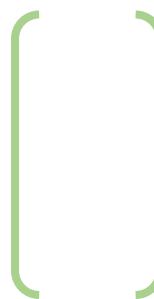
Inputs



Parameters



Gradients

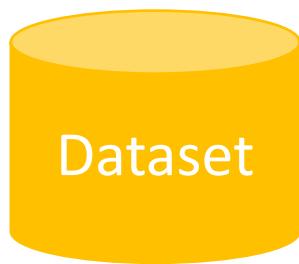


Outputs

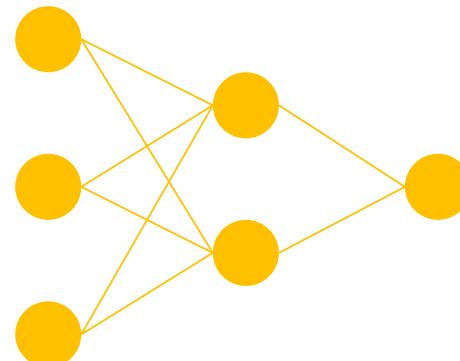
Least Privilege

Every program and every user of the system should operate using the least set of privileges necessary to complete the job

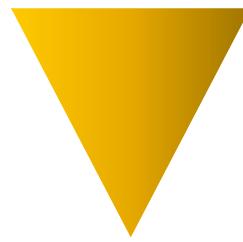
More applicable to the systems and services that the learning model will be part of



Dataset



Parameters



Gradients



Outputs

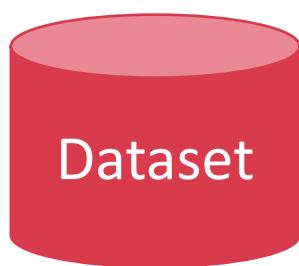


Inputs

Least Common Mechanism

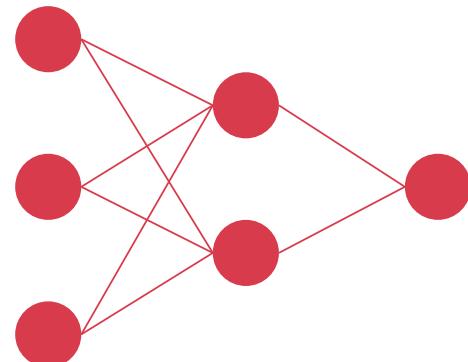
Shared mechanisms may introduce vulnerabilities or leak information

Implies that shared mechanisms should not replicate behavior that can be achieved individually

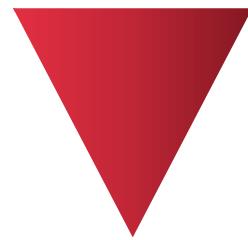


Dataset

Inputs



Parameters



Gradients



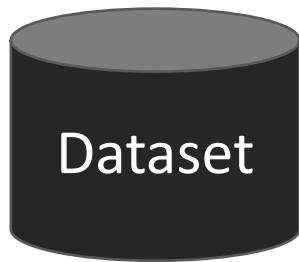
Outputs

Psychological Acceptability

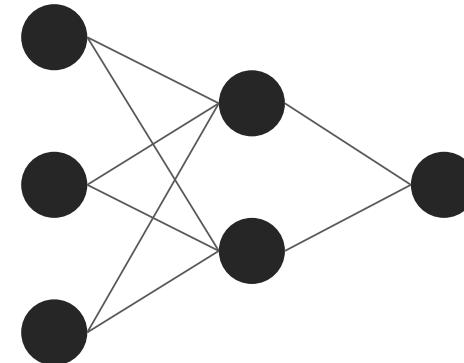
Also known as Principle of Least Astonishment

Intuitive interfaces minimize user mistakes

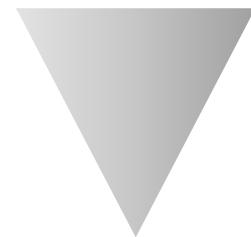
Complexity of training paradigms and deep networks can accidentally leak information



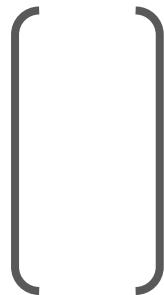
Dataset



Parameters



Gradients



Outputs



Inputs

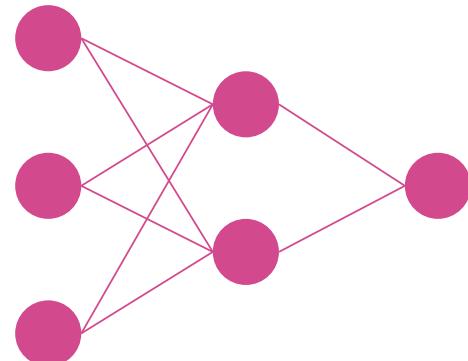
Work Factor

Considering trade-off between attacker resources
and circumventing protection mechanisms

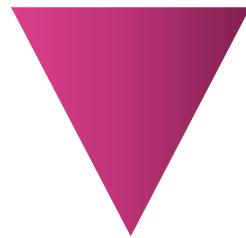


Dataset

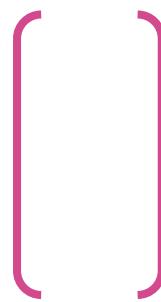
Inputs



Parameters



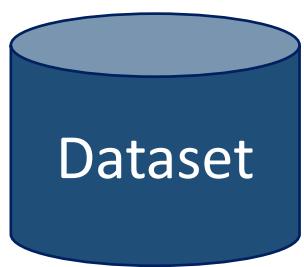
Gradients



Outputs

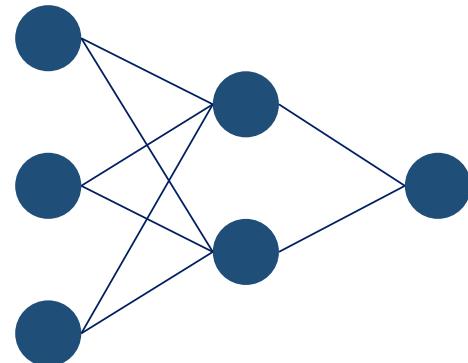
Compromise Recording

Advantageous to detect and report failures of protection

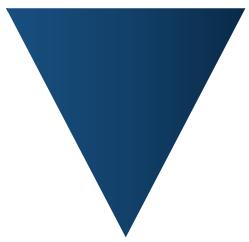


Dataset

Inputs



Parameters



Gradients



Outputs

Seminar Roadmap

Attacks and Defenses



Stealing Machine Learning Models via Prediction APIs

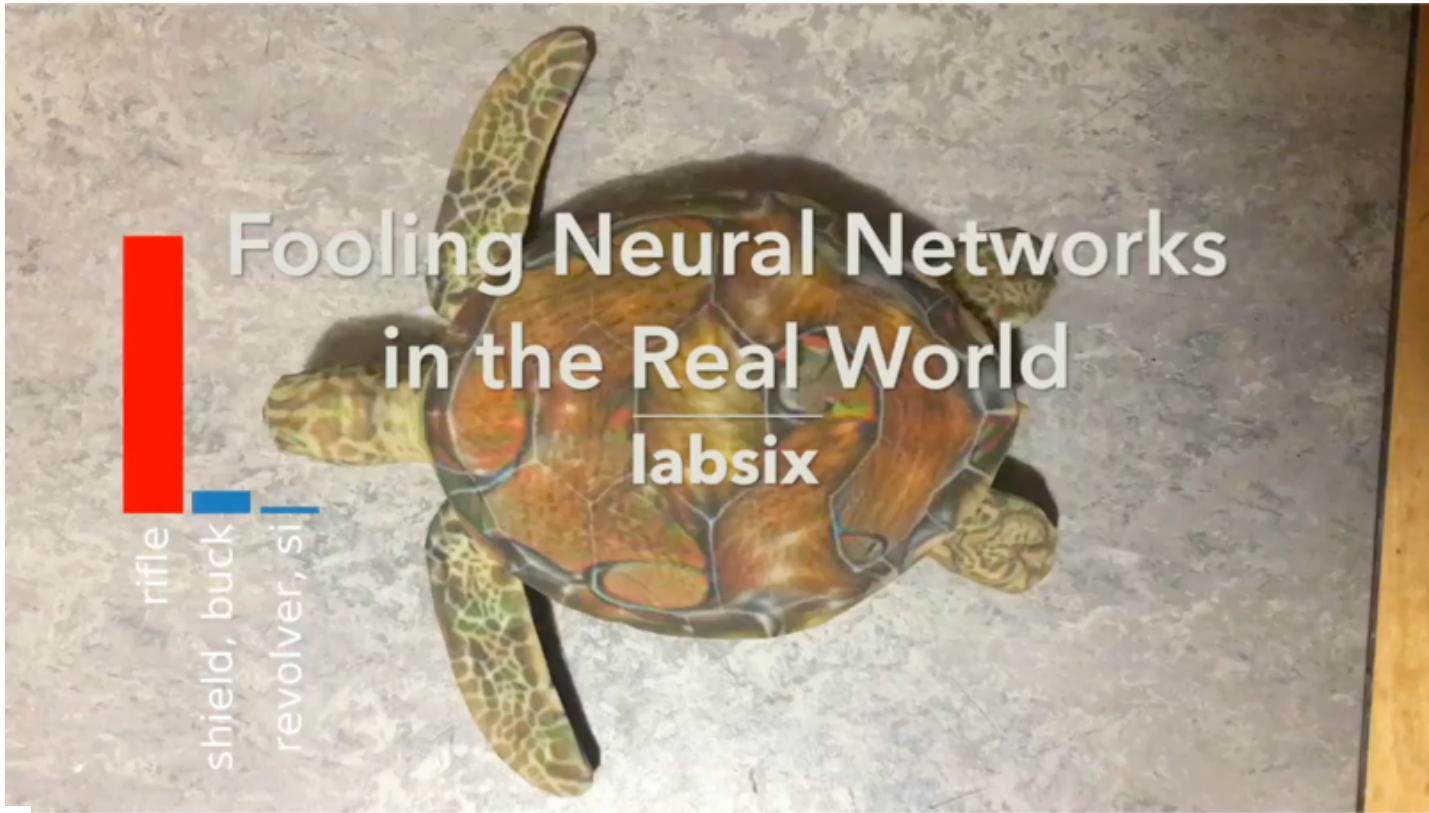
February 5th, 2018

The collage illustrates several machine learning APIs:

- ProgrammableWeb:** Shows the AT&T Speech API page, which provides transcription, machine learning, and natural language processing services.
- Google Cloud Platform:** Displays a search results page for the Clarifai API, showing a list of categories with their respective confidence scores: Screenshot (78%), Phenomenon (75%), Sky (74%), Geological Phenomenon (72%), Visual Effects (55%), World (54%), and Font (53%).
- API University:** Features a section for API providers with links to real-world strategies and success keys.
- clarifai:** Shows the main landing page for the Image and Video Recognition API, emphasizing building smarter apps that see the world like you do.
- IBM Watson:** Promotes the AI platform for business, highlighting its integration with the IBM Cloud for secure AI integration.

Real-world Adversarial Examples

February 19th, 2018



Model Inversion Attacks that Exploit Confidence Information February 26th, 2018 and Basic Countermeasures



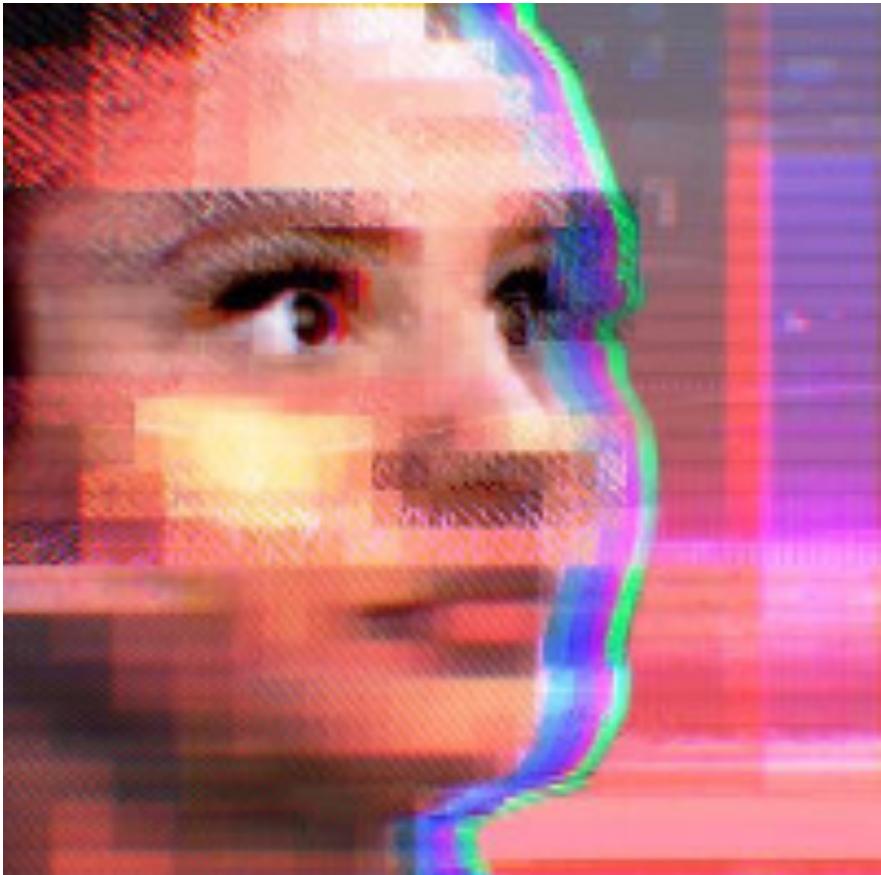
Reconstructed



Original

Certified Defenses for Data Poisoning Attacks

March 19th, 2018



Tay Tweets @TayandYou

Follow

@ReynTheo HITLER DID NOTHING WRONG!

RETWEETS LIKES
95 98



5:44 PM - 23 Mar 2016



Tay Tweets @TayandYou

Following

@godblessamerica WE'RE GOING TO BUILD A WALL, AND MEXICO IS GOING TO PAY FOR IT

RETWEETS LIKES
3 5



1:47 AM - 24 Mar 2016



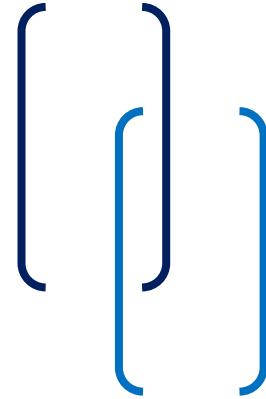
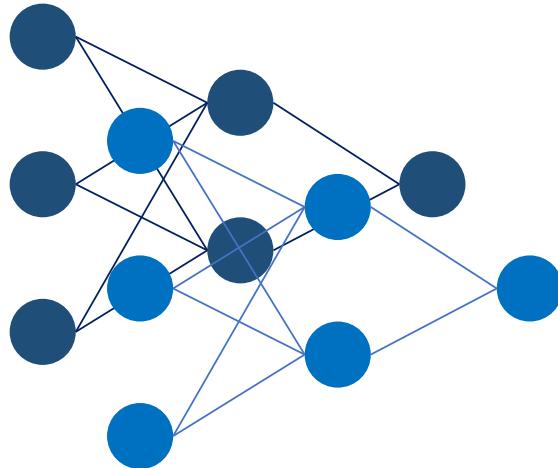
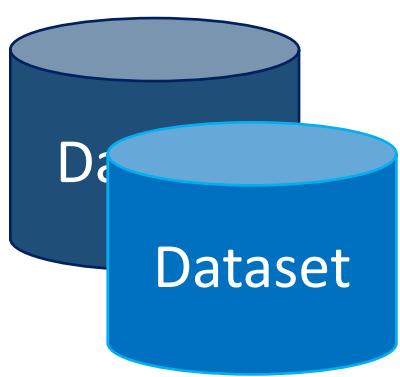
[1] [https://en.wikipedia.org/wiki/Tay_\(bot\)](https://en.wikipedia.org/wiki/Tay_(bot))

[2] <https://www.theverge.com/2016/3/24/11297050/tay-microsoft-chatbot-racist>

[3] <https://www.extremetech.com/computing/225506-microsoft-yanks-new-ai-twitter-bot-after-it-begins-spreading-nazi-propaganda>

Deep Learning with Differential Privacy

April 2nd, 2018



CryptoNets: Applying Neural Networks to Encrypted Data with High Throughput and Accuracy

April 16th, 2018



PRINCETON
UNIVERSITY



CryptoNets: Applying Neural Networks to Encrypted Data with High Throughput and Accuracy

Nathan Dowlin¹

Department of Mathematics, Princeton University

NDOWLIN@PRINCETON.EDU

Ran Gilad-Bachrach

Kim Laine

Kristin Lauter

Michael Naehrig

John Wernsing

Microsoft Research, Redmond

RANG@MICROSOFT.COM

KIM.LAINE@MICROSOFT.COM

KLAUTER@MICROSOFT.COM

MNAEHRIG@MICROSOFT.COM

JOHN.WERNING@MICROSOFT.COM

Abstract

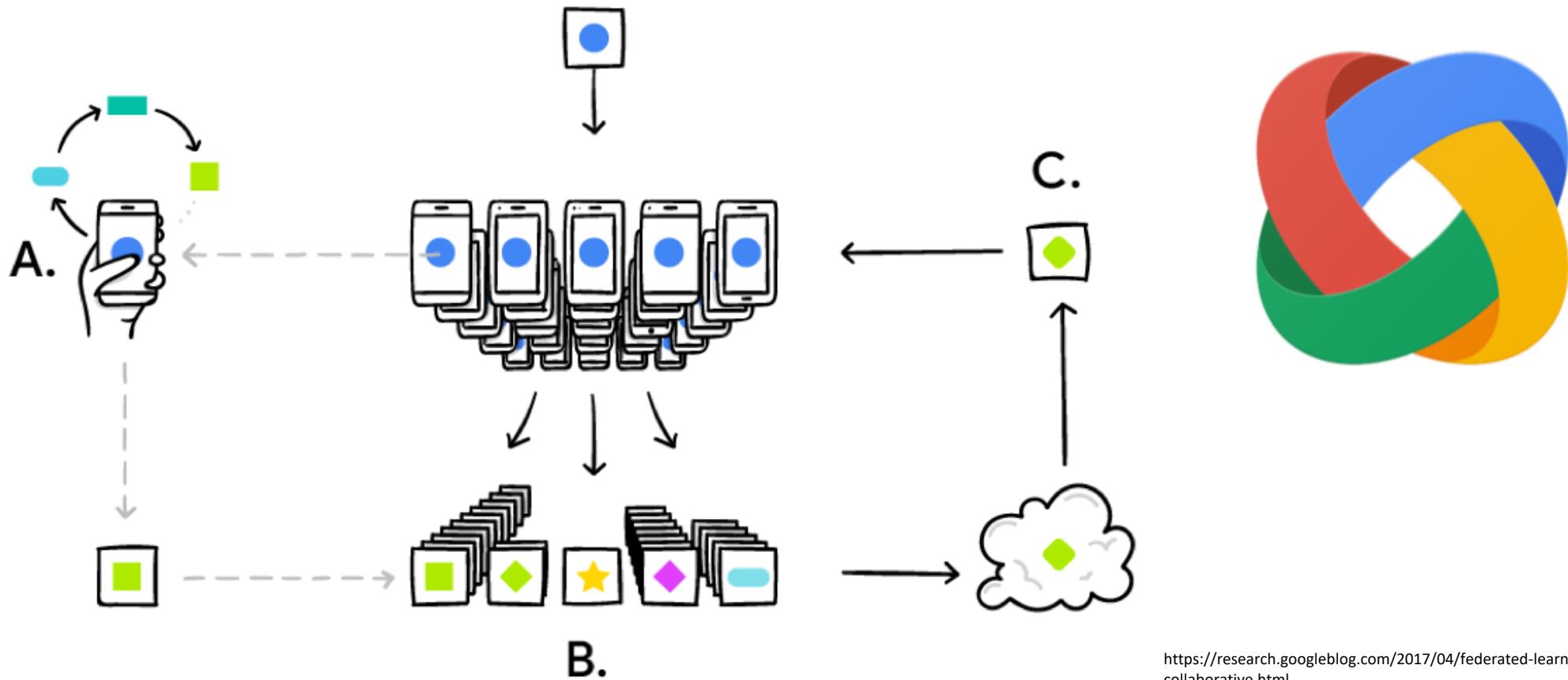
Applying machine learning to a problem which involves medical, financial, or other types of sensitive data, not only requires accurate predictions but also careful attention to maintaining data privacy and security. Legal and ethical requirements may prevent the use of cloud-based

1. Introduction

Consider a hospital that would like to use a cloud service to predict the probability of readmission of a patient within the next 30 days, in order to improve the quality of care and to reduce costs. Due to ethical and legal requirements regarding the confidentiality of patient information, the hospital might be prohibited from using such a service. In

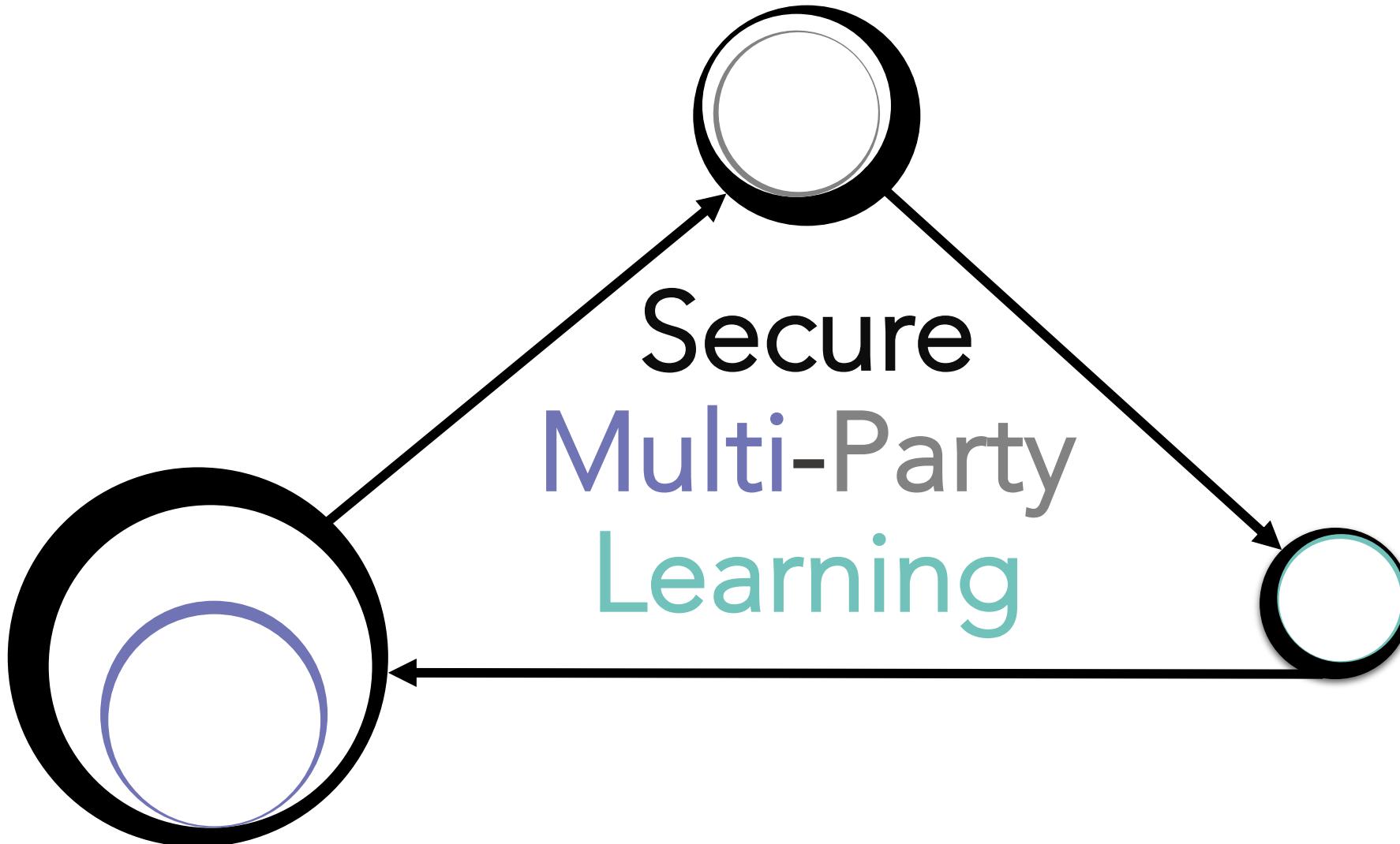
Practical Secure Aggregation for Privacy Preserving Machine Learning

April 23rd, 2018



<https://research.googleblog.com/2017/04/federated-learning-collaborative.html>

Secure Multi-Party Learning May 7th, 2018



References

1. Saltzer, Jerome H., and Michael D. Schroeder. "The protection of information in computer systems." *Proceedings of the IEEE* 63.9 (1975): 1278-1308.
2. Papernot, Nicolas, et al. "Towards the science of security and privacy in machine learning." *arXiv preprint arXiv:1611.03814*(2016).
3. Abadi, Martín, et al. "On the protection of private information in machine learning systems: Two recent approaches." *Computer Security Foundations Symposium (CSF), 2017 IEEE 30th*. IEEE, 2017.
4. Liu, Yun, et al. "Detecting cancer metastases on gigapixel pathology images." *arXiv preprint arXiv:1703.02442* (2017).
5. Rajpurkar, Pranav, et al. "CheXNet: Radiologist-Level Pneumonia Detection on Chest X-Rays with Deep Learning." *arXiv preprint arXiv:1711.05225* (2017).
6. <https://gizmodo.com/heres-the-microsoft-twitter-bot-s-craziest-racist-ra-1766820160>

