BISHOP FOX®

# Weaponizing Machine Learning
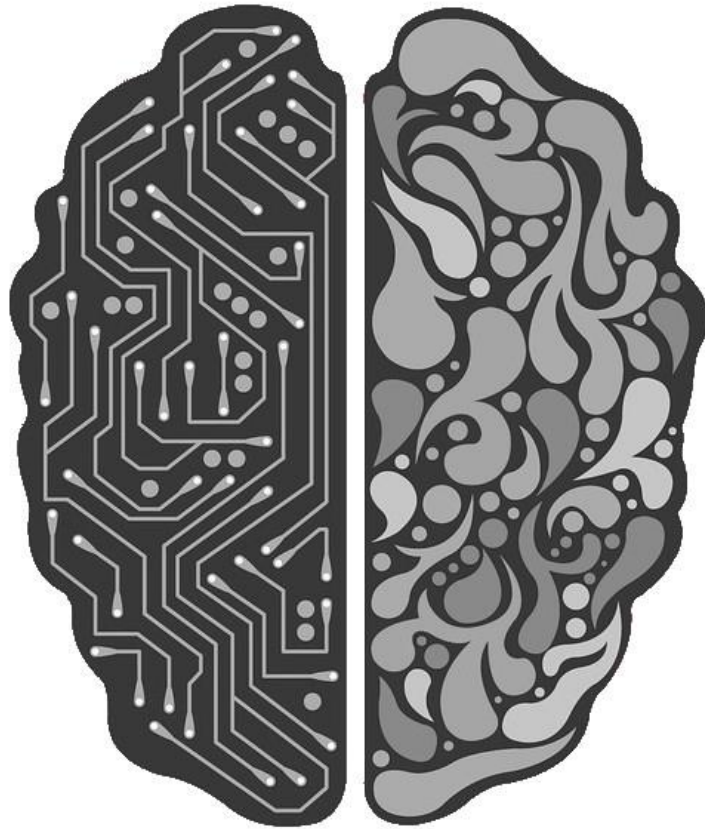
## HUMANITY WAS OVERRATED ANYWAY



July 30, 2017

# DeepHack

**Artificial Intelligence**



**Hacking**

# The Scenario
I'D LIKE TO OPEN AN ACCOUNT



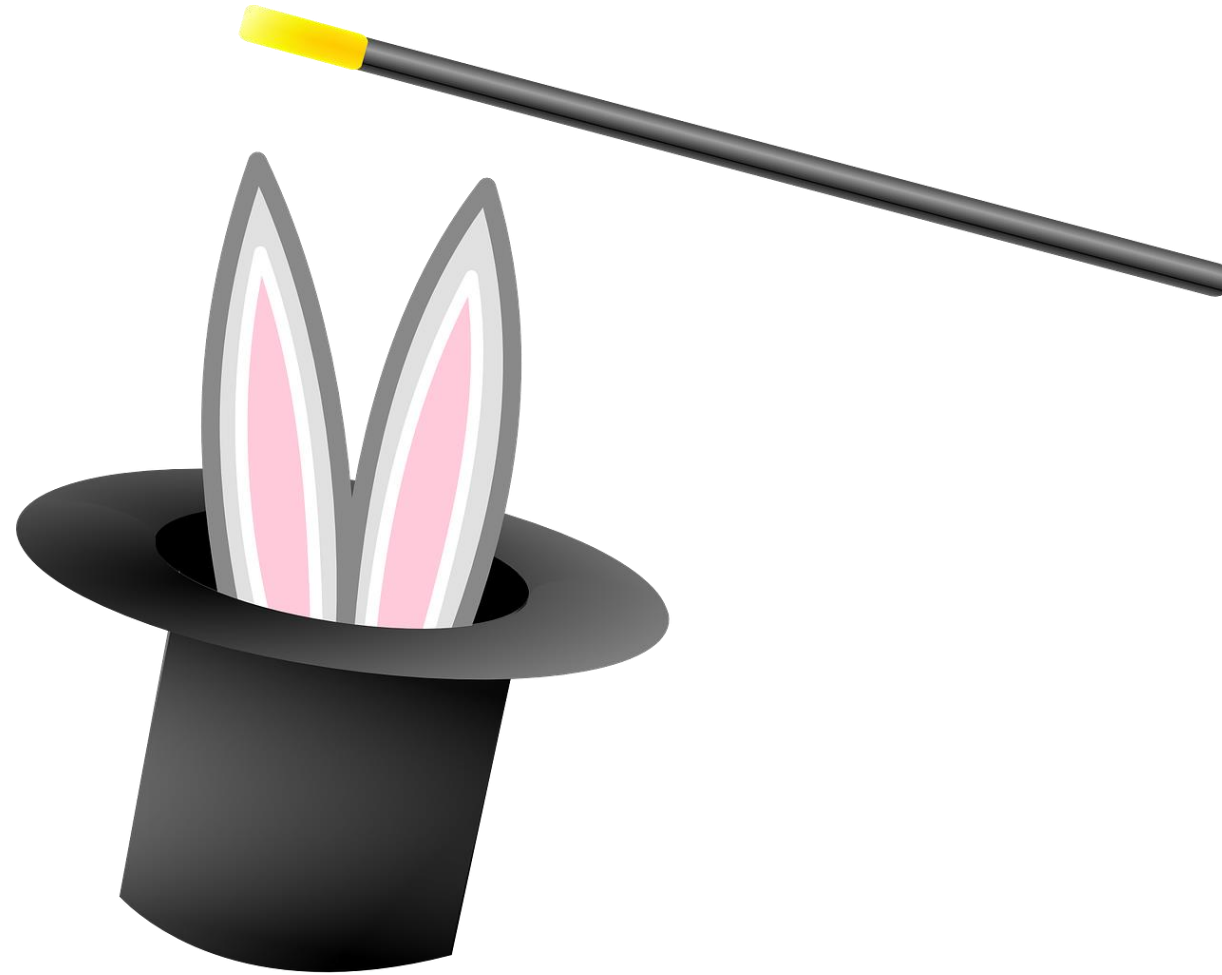Attacker → Internet → Bank Website → SQL Database

# DEMONSTRATION

*CROSSES FINGERS*

# No Tricks Up Our Sleeve

**FOOL ME ONCE**

# The Model

Code → Binary → Action

**Regular** Program

Code / Data → Model → Action

**AI** Program

# AI Programming

# 1. Object Oriented

# 2. Functional

# 3. Machine Learning

# Machine Learning 101

**Maze Solving Robot**

# Machine Learning 101

**Maze Solving Robot**

| State | Action | Reward |
|---|---|---|
| 0,0 | Up Right Down Left | |
| 0,1 | Up Right Down Left | |
| 0,2 | Up Right Down Left | |

# Machine Learning 101

**Maze Solving Robot**

| State | Action | Reward |
|-------|--------|--------|
| 0,0 | Up<br>Right<br>Down<br>Left | **-50** |
| 0,1 | Up<br>Right<br>Down<br>Left | |
| 0,2 | Up<br>Right<br>Down<br>Left | |

# Machine Learning 101

**Maze Solving Robot**

| State | Action | Reward |
|---|---|---|
| 0,0 | Up<br>Right<br>Down<br>Left | -1<br><br><br>-50 |
| 0,1 | Up<br>Right<br>Down<br>Left | |
| 0,2 | Up<br>Right<br>Down<br>Left | |

# Machine Learning 101

**Maze Solving Robot**

| State | Action | Reward |
|---|---|---|
| 0,0 | Up<br>Right<br>Down<br>Left | -1<br><br><br>-50 |
| 0,1 | Up<br>Right<br>Down<br>Left | <br>-50<br><br> |
| 0,2 | Up<br>Right<br>Down<br>Left | |

# Machine Learning 101

**Maze Solving Robot**

| State | Action | Reward |
|-------|--------|--------|
| 0,0 | Up<br>Right<br>Down<br>Left | -1<br>-50<br>-50<br>-50 |
| 0,1 | Up<br>Right<br>Down<br>Left | -1<br>-50<br>-1<br>-50 |
| 0,2 | Up<br>Right<br>Down<br>Left | -50<br>-50<br>-1<br>+50 |

# Machine Learning 101

**Chess Playing Robot**

# Machine Learning 101

**Chess Playing Robot**

| State | Action | Reward |
|-------|--------|--------|
| ??? | ??? ??? ??? ??? | -?? -?? -?? -?? |

$\sim \mathbf{10^{47}}$ States in Chess

Can't store it all

# Function Approximation

IS NOT MAGIC

| State | Action | Reward |
|-------|--------|--------|
| 0,0 | Up<br>Right<br>Down<br>Left | -1<br>-50<br>-50<br>-50 |
| 0,1 | Up<br>Right<br>Down<br>Left | -1<br>-50<br>-1<br>-50 |
| 0,2 | Up<br>Right<br>Down<br>Left | -50<br>-50<br>-1<br>+50 |

# Function Approximation

**IS NOT MAGIC**

| State | Action | Reward |
|-------|--------|--------|
| 0,0 | Up<br>Right<br>Down<br>Left | **-1**<br>**-50**<br>**-50**<br>**-50** |
| 0,1 | Up<br>Right<br>Down<br>Left | **-1**<br>**-50**<br>**-1**<br>**-50** |
| 0,2 | Up<br>Right<br>Down<br>Left | **-50**<br>**-50**<br>**-1**<br>**+50** |

## Is *just* one of these:



(A function)

# Math
NOT EVEN ONCE

STOCHASTIC GRADIENT DESCENT

BILINEAR INTERPOLATION

LINEAR REGRESSION

CONVOLUTIONAL TRANSPOSITION

# Function Approximation
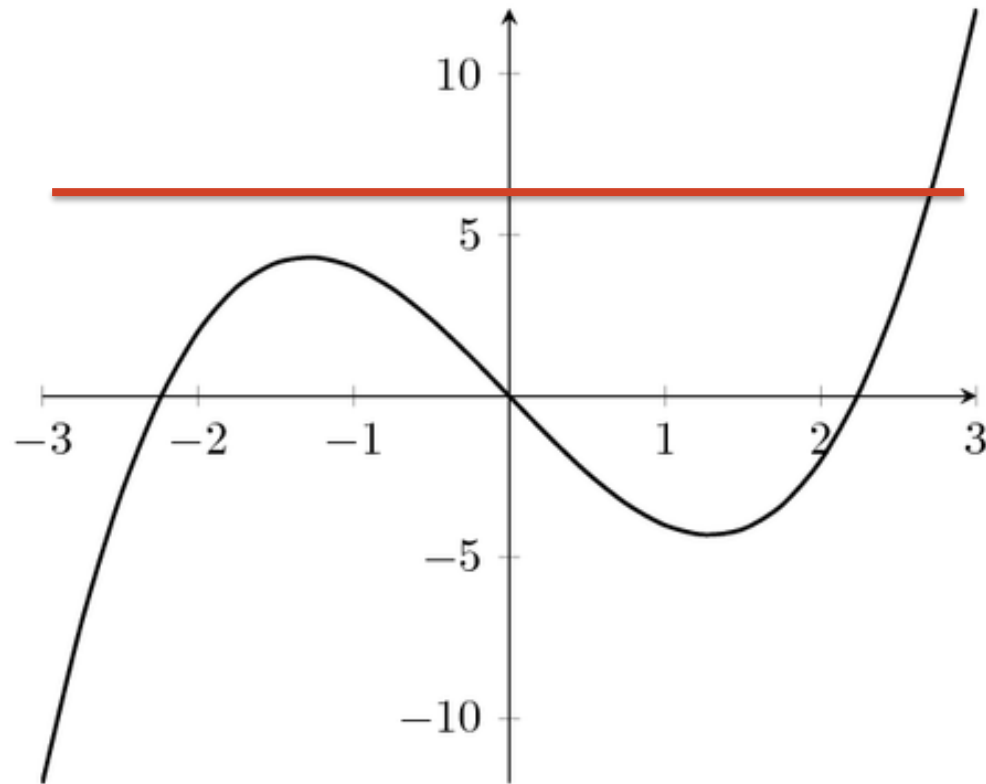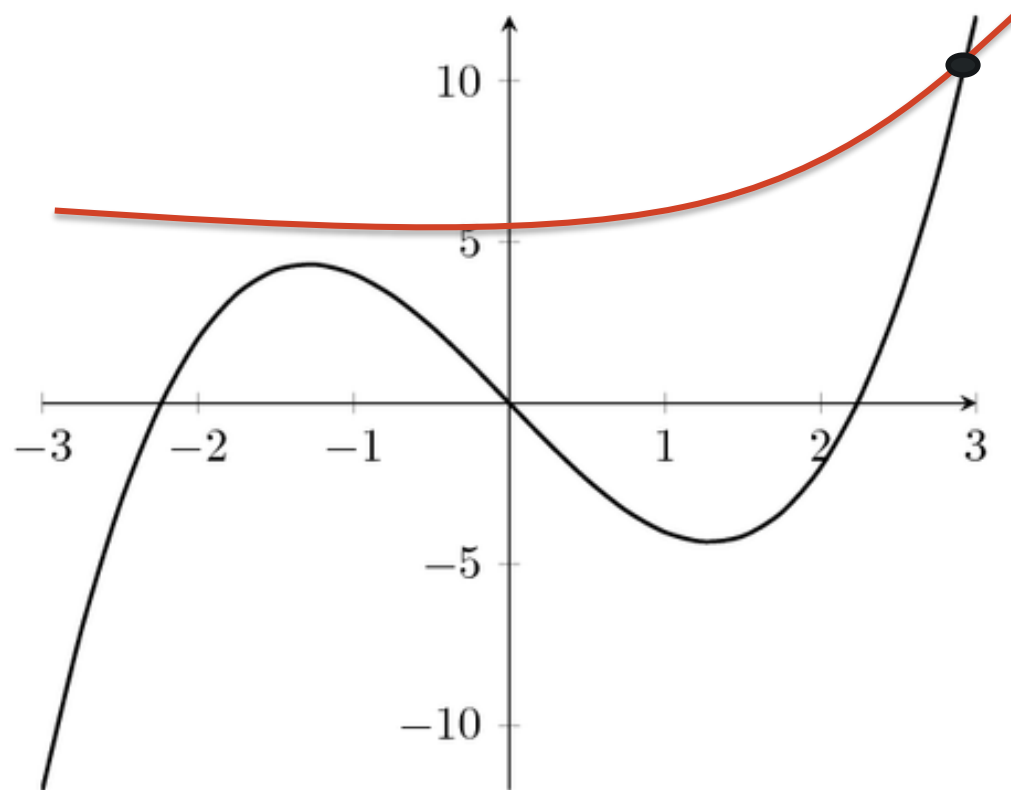
(A function)

# Function Approximation
IS NOT MAGIC



(A function)

# Function Approximation
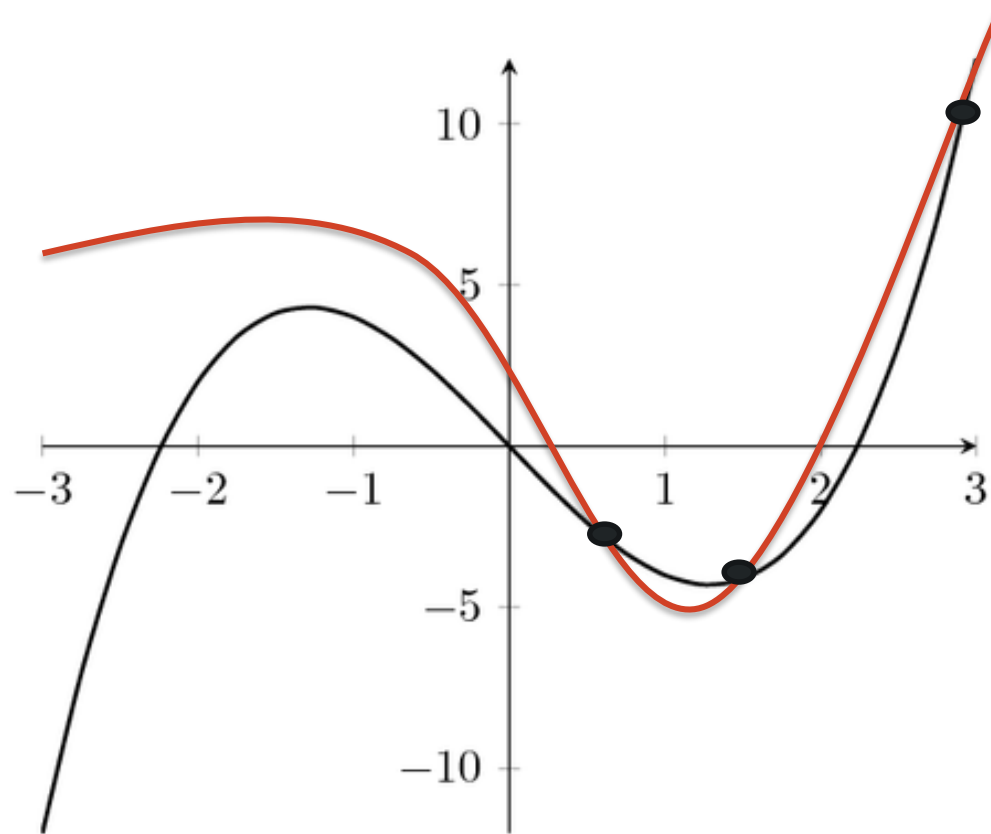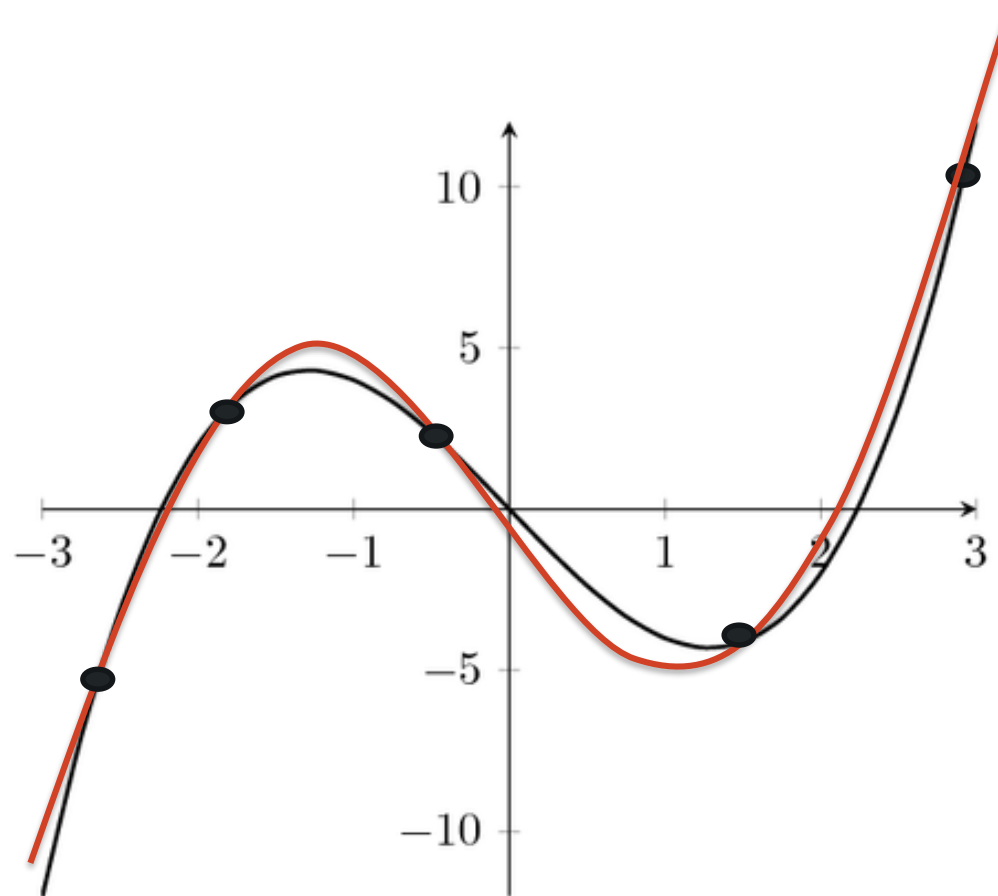
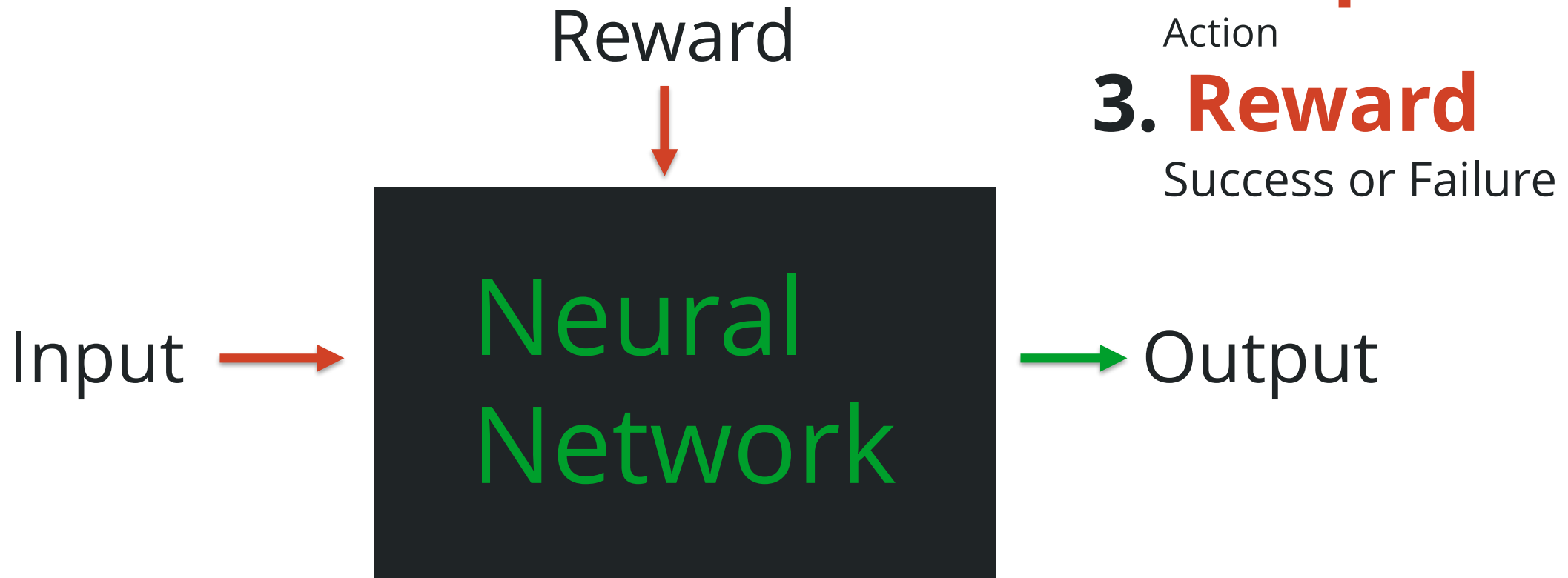**IS NOT MAGIC**



(A function)

# Function Approximation

(A function)

# Function Approximation

**IS NOT MAGIC**



(A function)

# The Neural Network
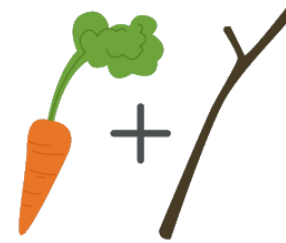
1. **Input**
Environment

2. **Output**
Action

3. **Reward**
Success or Failure
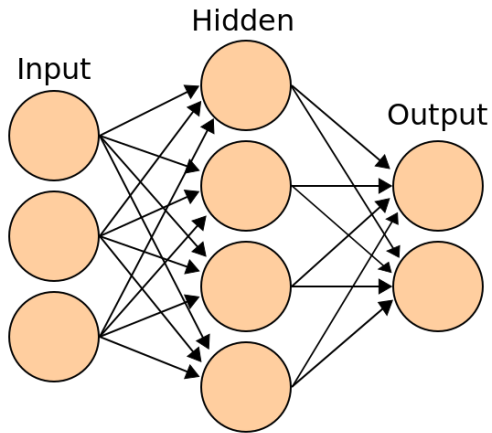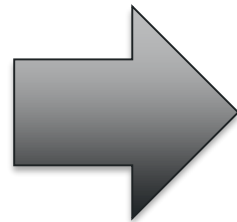
Reward

↓

Input →

## Neural Network

→ Output

# Example: Chess

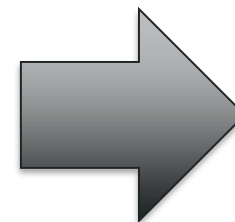**Reward**: Material gained / lost

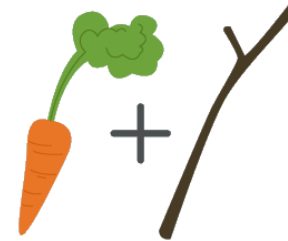**Input**: Current Piece Positions

Neural Network

Input

Hidden

Output

**Output**: Move one piece

# Example: DeepHack

**Reward**: HTTP Status (200/500)

**SELECT * FRO**

**Input**: Query String

Input

Hidden

Output

Neural Network

**M**

**Output**: Next Character

# Autocomplete Game

A lot of foo

d? t? s?

# Training

- Harvest good labeled data

- Bootstrap your model with experiences

- ... Or get your users to do it for you

Select all images below
that match this one:

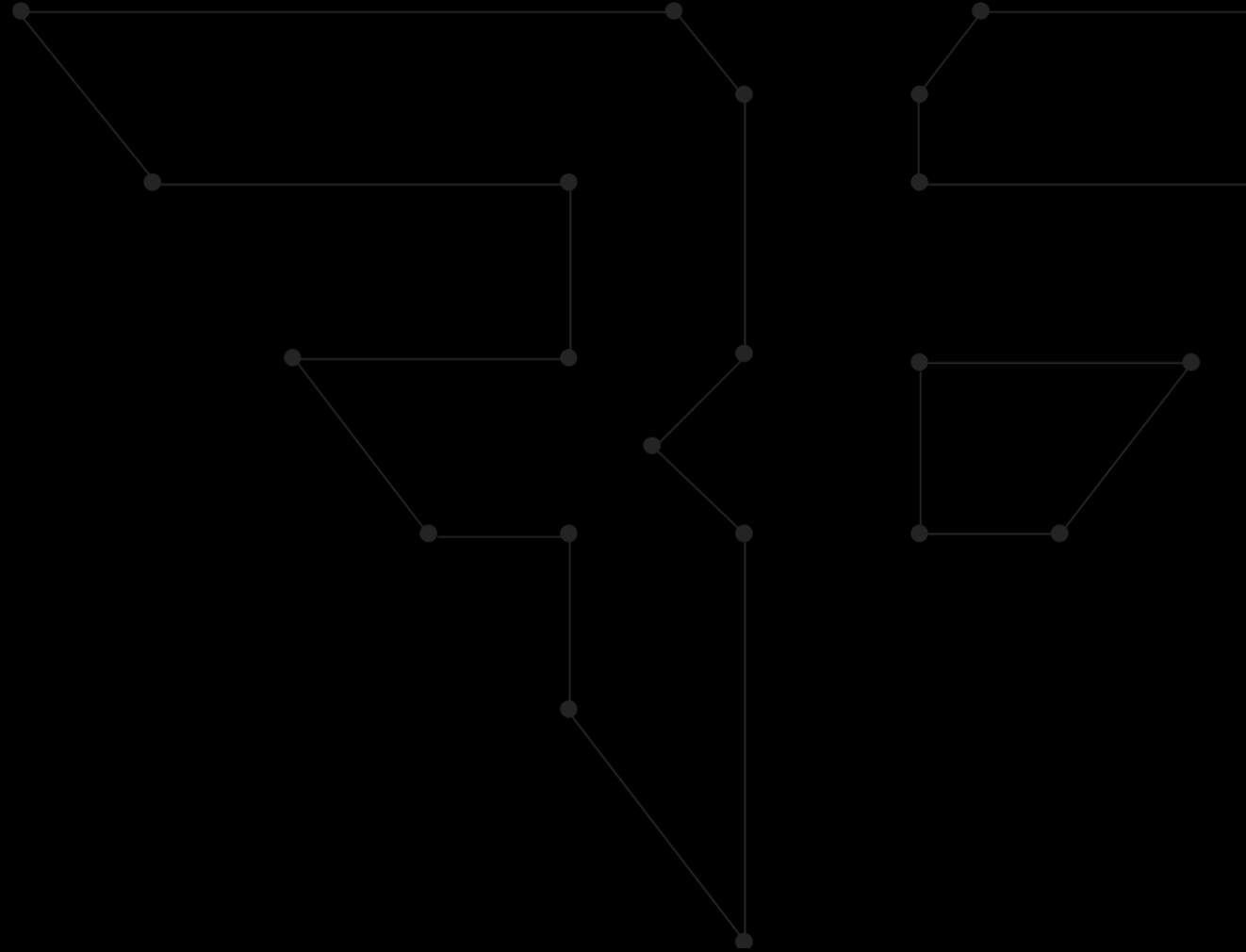# So What?

I DON'T GET IT

# DEMONSTRATION
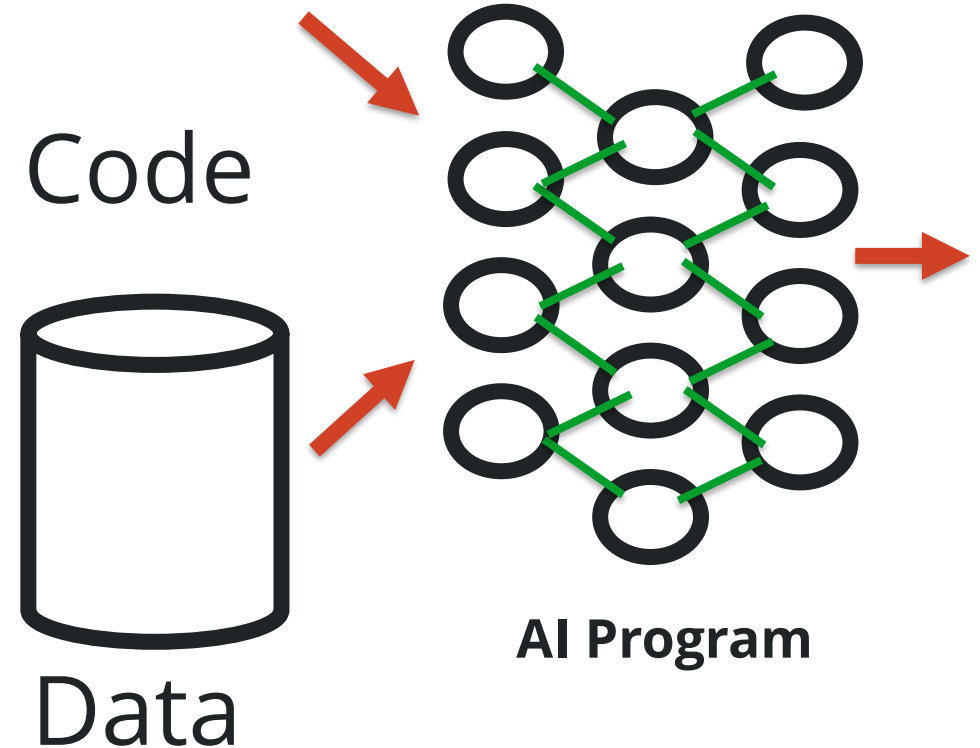
*CROSSES FINGERS*

# Lessons Learned

- **Quality training data is important**
  - And hard to get
  - Garbage in, Garbage out

- **Be careful about what you reward**
  - You will get more of it

- **Get a GPU**
  - Or better yet, a lot of them

# Other Considerations

- **Inherent Proprietary-ness**
- **Unreliable factor**
  - Stochastic
  - Undebuggable
    - o Trained too much?
    - o Not enough?
- **Power Imbalance**
  - Computational / Data

Code

Data

**AI Program**

# Future Work

- **Instrumented Webapp Fuzzer**

  - Available data?

- **Password Bruteforcing**

  - Context aware

- **Service Identification**

  - What's behind an open port?

- **Bad at**:

  - Finding new *classes* of vulnerabilities

# Questions!

**www.bishopfox.com**

**github.com/bishopfox**

**careers@bishopfox.com**

# Attributions (Images in Slides)

[Artificial intelligence brain image](#)

[Hacker image](#)

[Magic hat image](#)

[Chessboard image](#)

[Function plot graph image](#)

[Artificial neural network image](#)

[Chess pieces image](#)

[Carrot+Stick<Love image](#)