

# Measuring Neural Net Robustness with Constraints

Osbert Bastani, Yani Ioannou, Leonidas Lampropoulos,  
Dimitrios Vytiniotis, Aditya Nori, Antonio Criminisi

# Verification of Learning-Based Systems



robustness/stability



fairness

# Neural Net Robustness

$f$  :



$\mapsto$  “school bus”

# Adversarial Examples

(Szegedy et al. 2014)

# Adversarial Examples



“school bus”

(Szegedy et al. 2014)

# Adversarial Examples



“school bus”

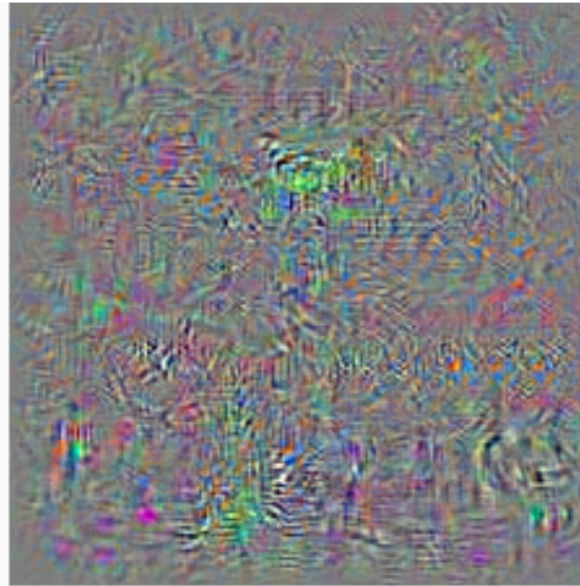
“ostrich”

(Szegedy et al. 2014)

# Adversarial Examples



“school bus”



perturbation (10×)

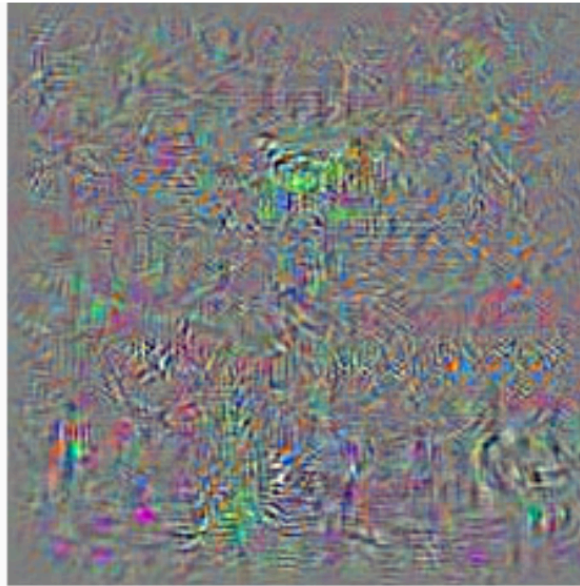
“ostrich”

(Szegedy et al. 2014)

# Adversarial Examples



“school bus”



perturbation (10×)



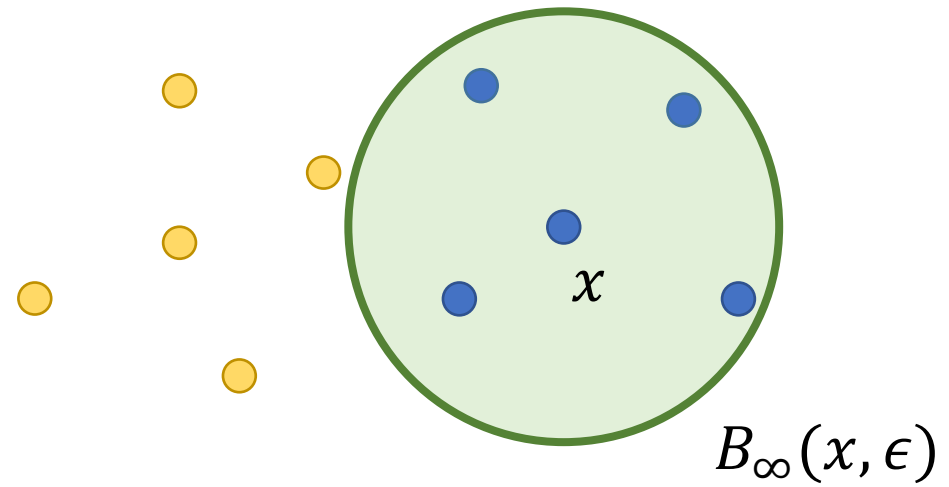
“ostrich”

(Szegedy et al. 2014)

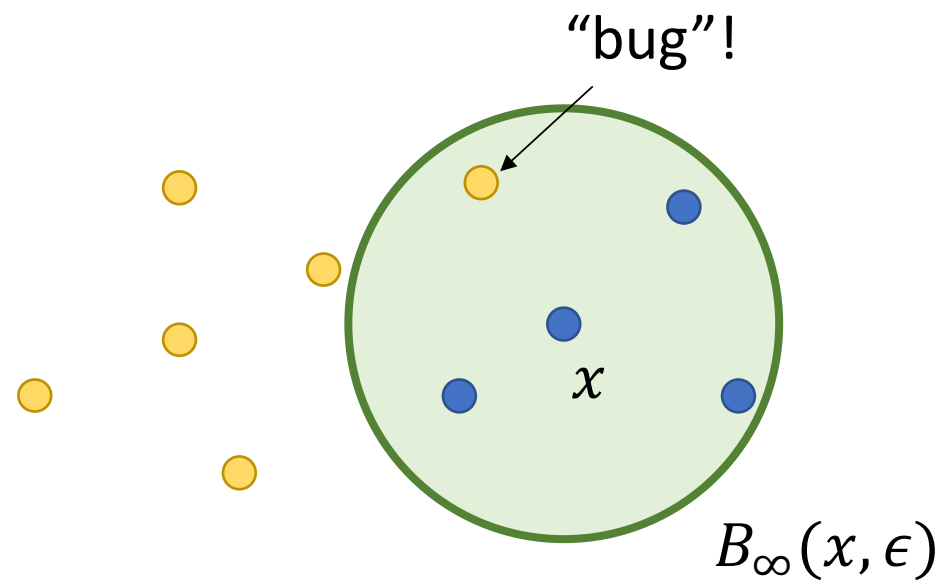


**robustness:** similar images  $\Rightarrow$  same label

**robustness:**  $\|x - x'\|_{\infty} \leq \epsilon \Rightarrow$  same label

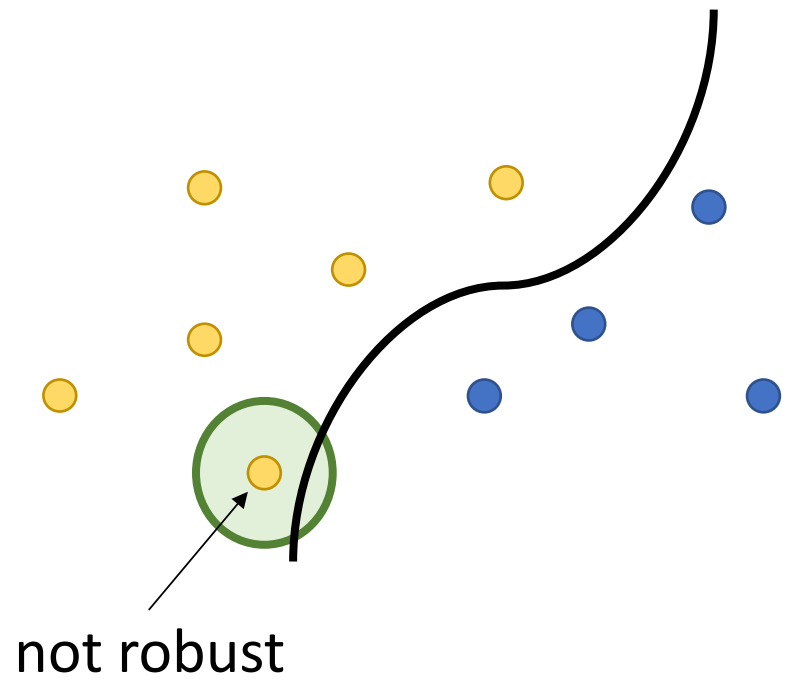


$\epsilon$ -robust at  $x$



**not**  $\epsilon$ -robust at  $x$

Can we **verify** robustness?

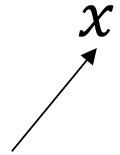


Can we **quantify** robustness?

$$\psi(f, \epsilon) = \Pr_x[f \text{ not } \epsilon\text{-robust at } x]$$

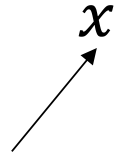


$$\psi(f, \epsilon) = \Pr[f \text{ not } \epsilon\text{-robust at } x]$$



input distribution (e.g., held-out test set)

$$\psi(f, \epsilon) = \Pr[f \text{ not } \epsilon\text{-robust at } x] \in [0, 1]$$



input distribution (e.g., held-out test set)

# Disclaimer

# Disclaimer

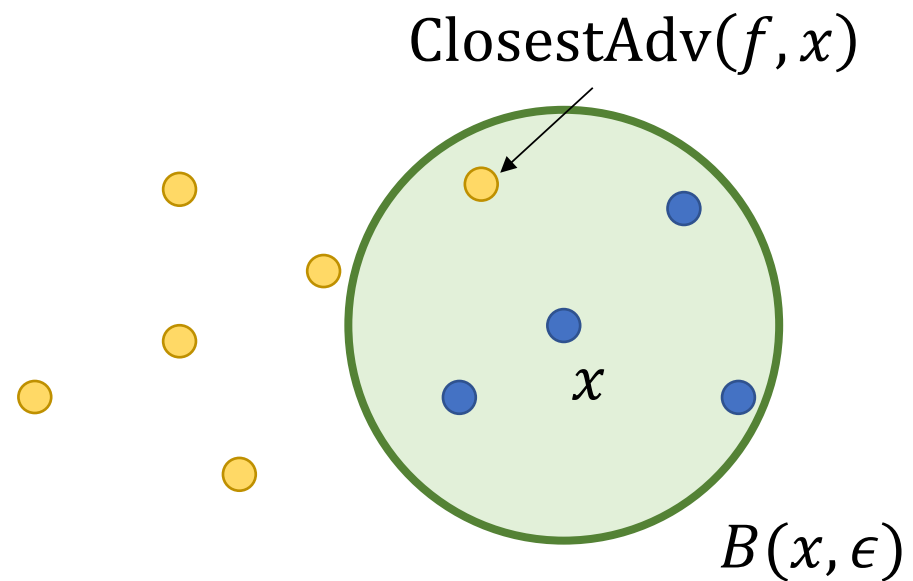
- Exactly quantifying neural net robustness did not scale

# Disclaimer

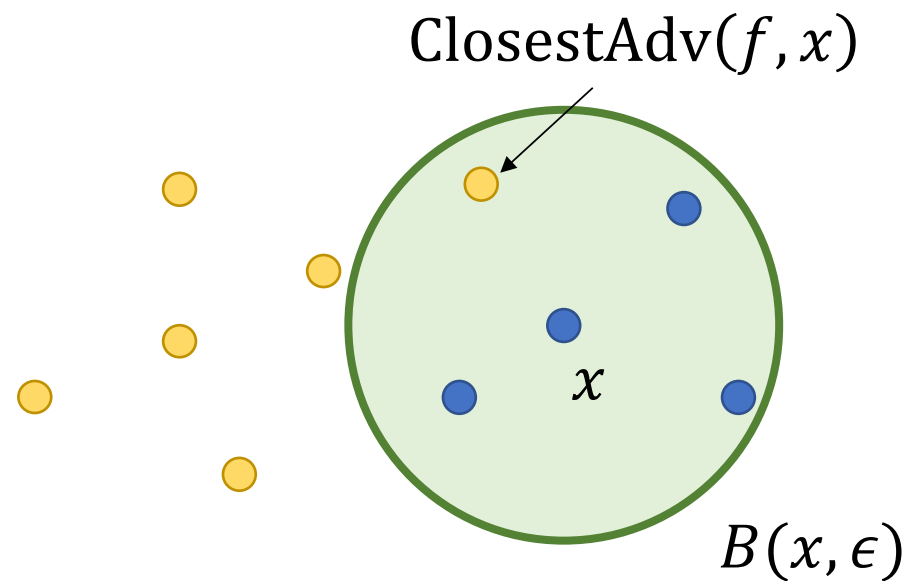
- Exactly quantifying neural net robustness did not scale
- We produce **useful approximations**

# Disclaimer

- Exactly quantifying neural net robustness did not scale
- We produce **useful approximations**
- At the end: possible paths forward?

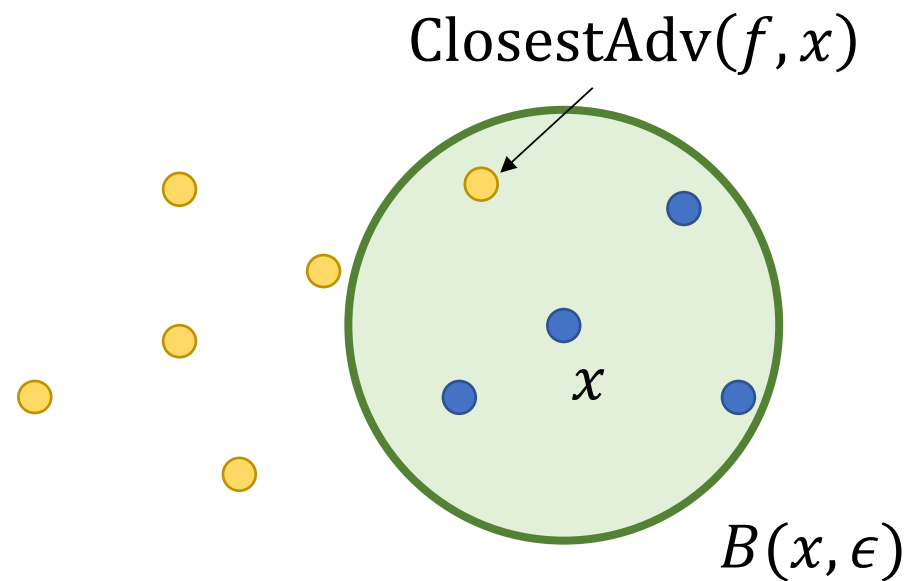


$$\psi(f, \epsilon) = \Pr_x [f \text{ not } \epsilon\text{-robust at } x]$$

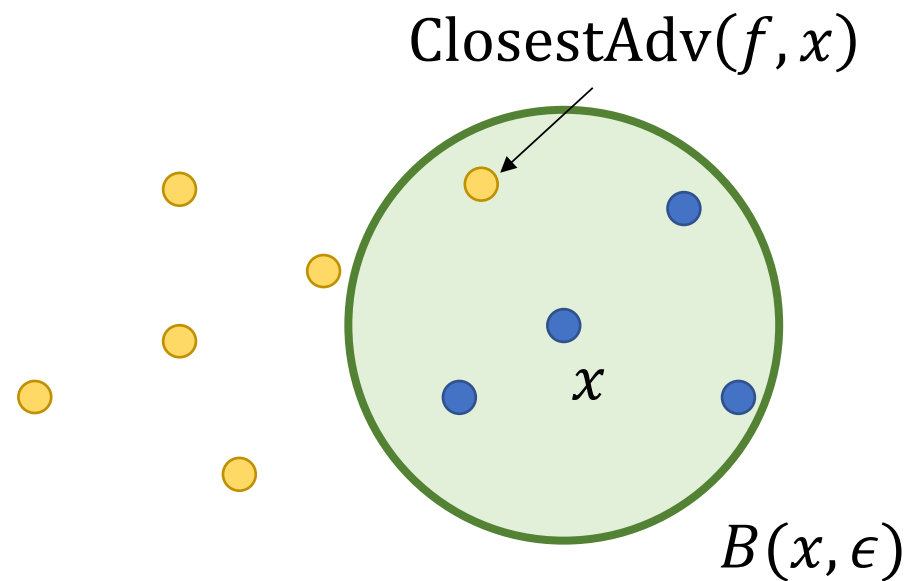


$$\psi(f, \epsilon) = \Pr_x [ \|\text{ClosestAdv}(f, x) - x\| \leq \epsilon ]$$





$$\psi(f, \epsilon) = \frac{1}{|X|} \sum_{x \in X} \mathbb{I}[\|\text{ClosestAdv}(f, x) - x\| \leq \epsilon]$$



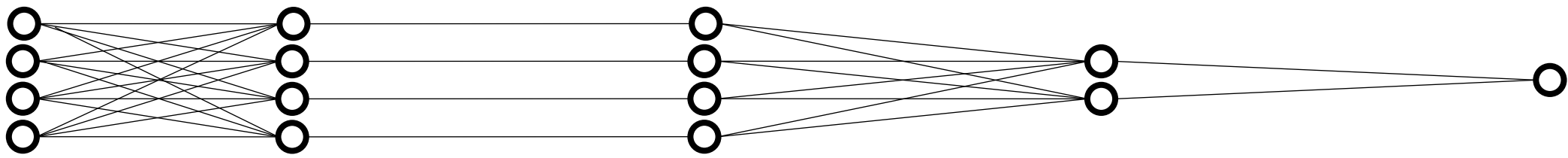
$$\psi(f, \epsilon) = \frac{1}{|X|} \sum_{x \in X} \mathbb{I}[\|\text{ClosestAdv}(f, x) - x\| \leq \epsilon]$$

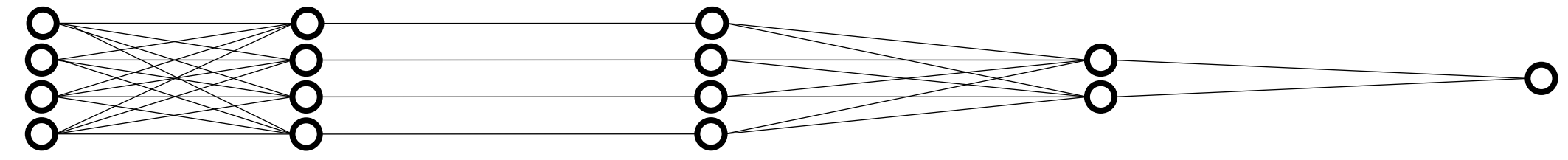
$$\begin{array}{ll} \arg \min_{x'} & \|x' - x\|_{\infty} \\ \text{subj. to} & f(x') \neq f(x) \end{array}$$

$$\begin{aligned} & \arg \min_{x'} \quad \|x' - x\|_{\infty} \\ & \text{subj. to} \quad \bigvee_{\ell \neq f(x)} f(x') = \ell \end{aligned}$$

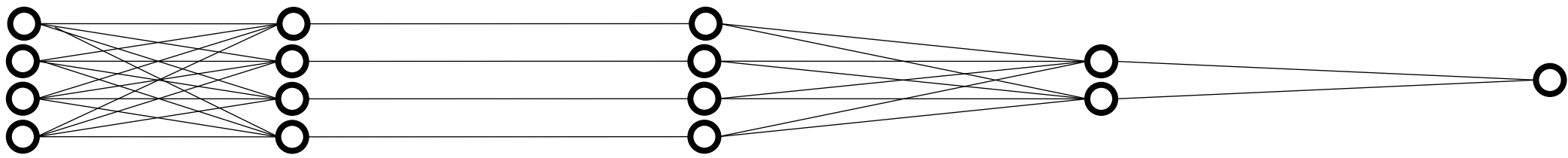
$$\begin{aligned} & \arg \min_{x'} \quad \|x' - x\|_{\infty} \\ & \text{subj. to} \quad \bigvee_{\ell \neq f(x)} f(x') = \ell \end{aligned}$$

$$\phi_f(x, \ell) = \mathbb{I}[f(x) = \ell]$$



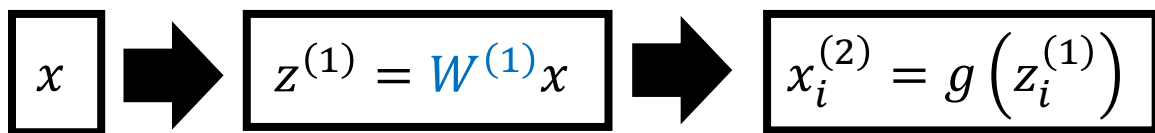
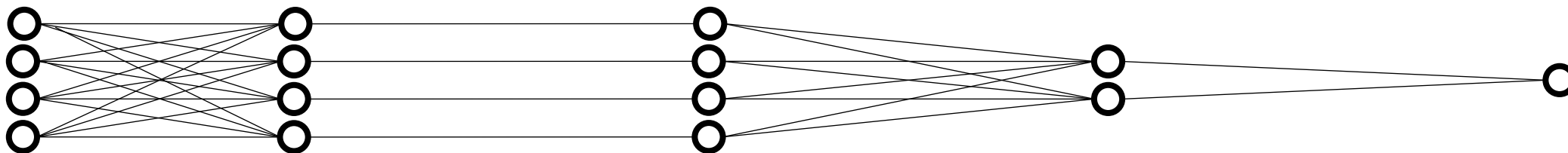


$x$

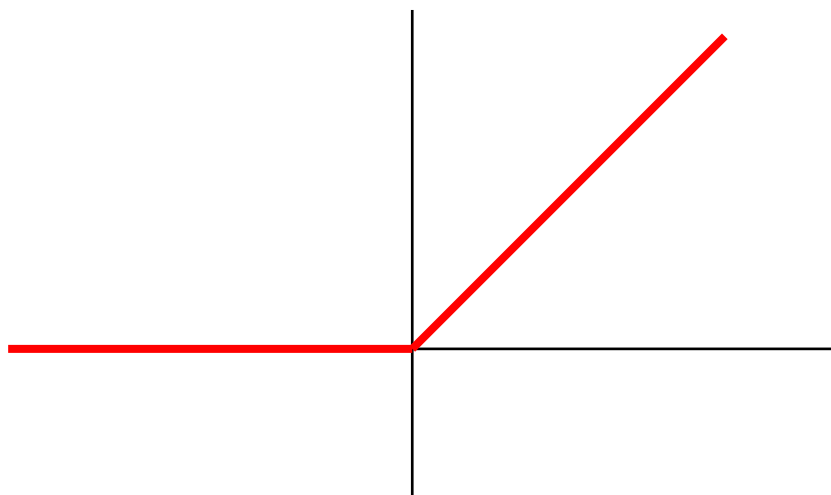


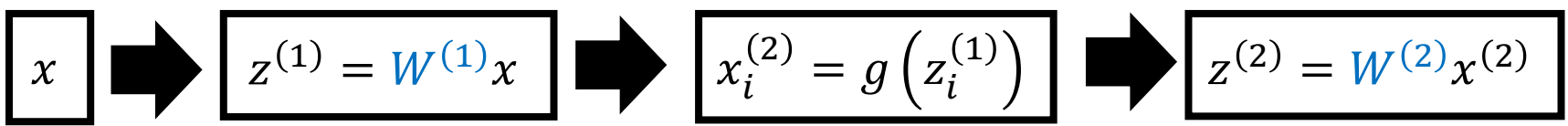
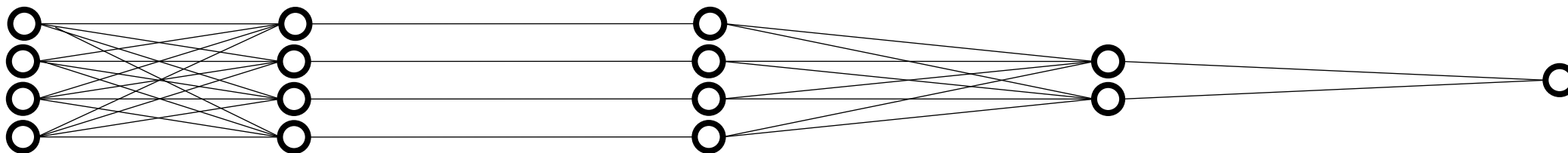
$x$   $\rightarrow$   $z^{(1)} = W^{(1)}x$



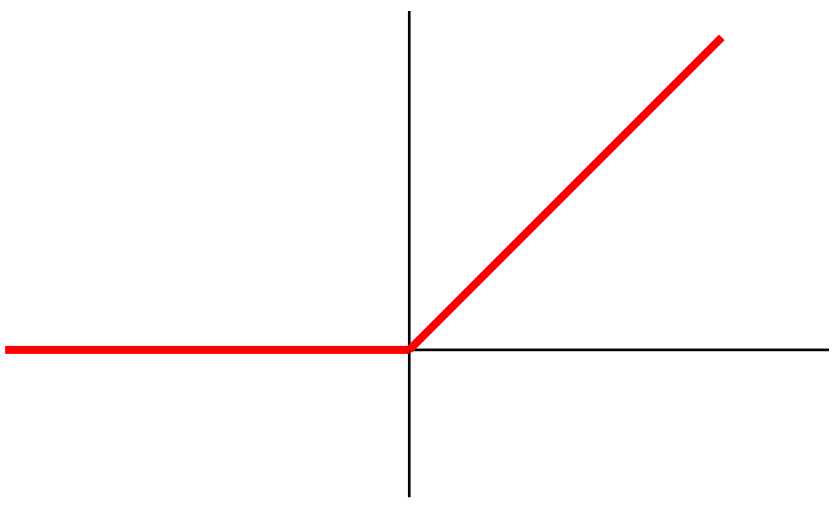


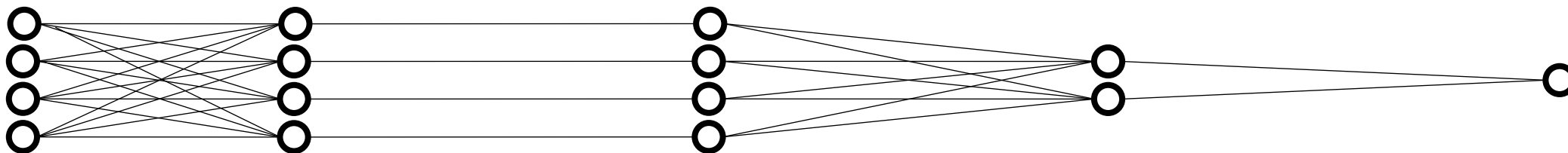
$$g(z) = \begin{cases} z & (z \geq 0) \\ 0 & (z \leq 0) \end{cases}$$



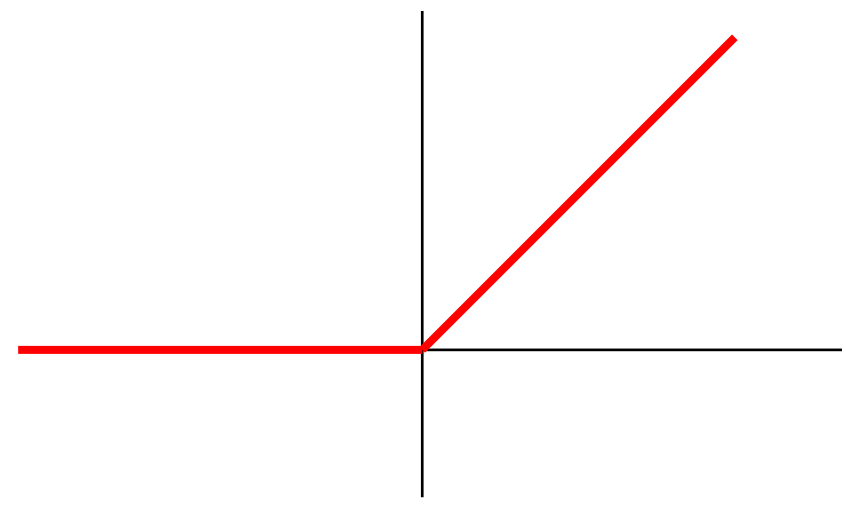


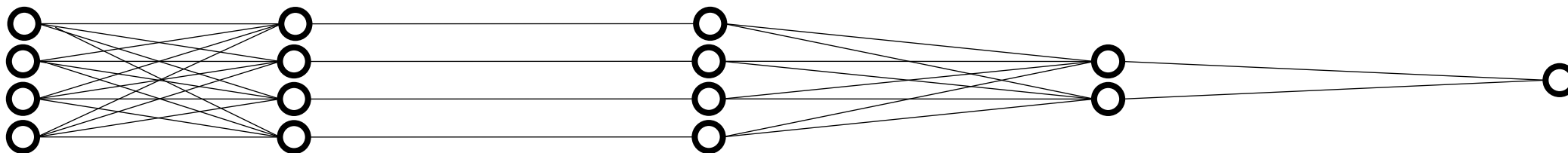
$$g(z) = \begin{cases} z & (z \geq 0) \\ 0 & (z \leq 0) \end{cases}$$





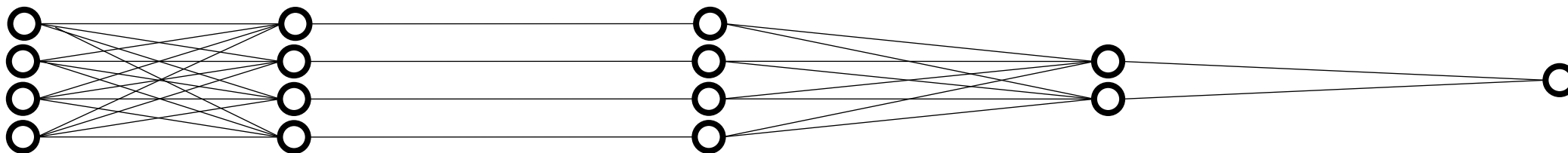
$$g(z) = \begin{cases} z & (z \geq 0) \\ 0 & (z \leq 0) \end{cases}$$





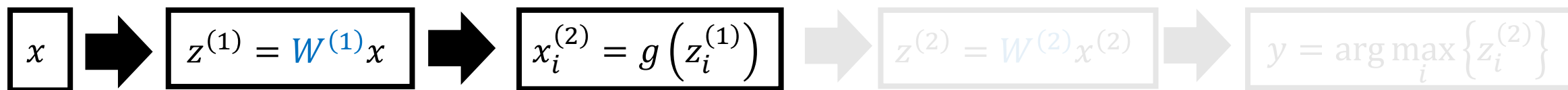
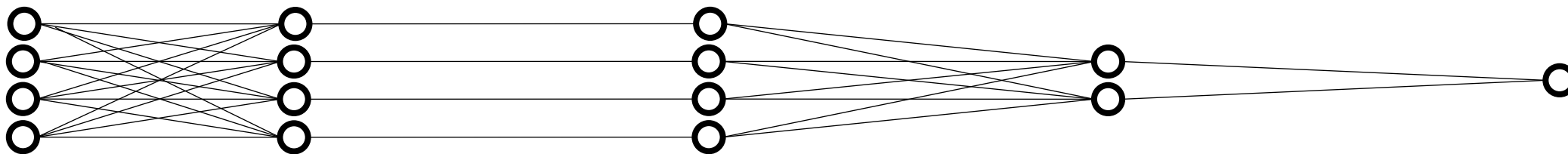
$$g(z) = \begin{cases} z & (z \geq 0) \\ 0 & (z \leq 0) \end{cases}$$

$$\phi_f(x, \ell) =$$



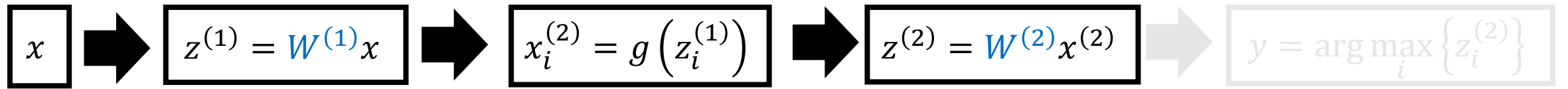
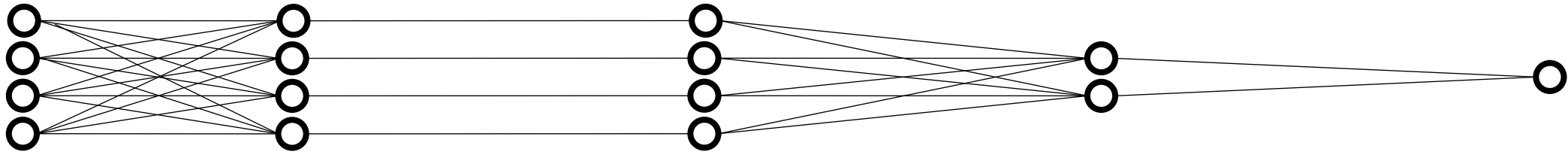
$$g(z) = \begin{cases} z & (z \geq 0) \\ 0 & (z \leq 0) \end{cases}$$

$$\phi_f(x, \ell) = (z^{(1)} = W^{(1)}x)$$



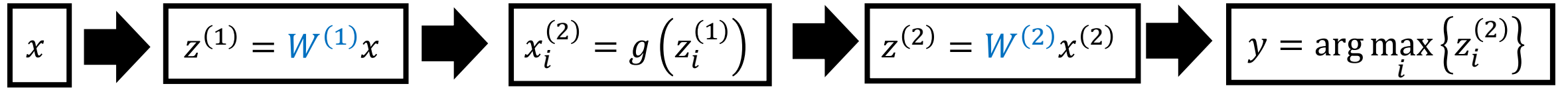
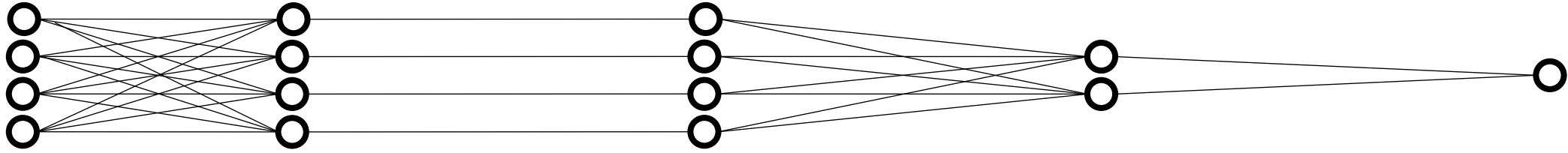
$$g(z) = \begin{cases} z & (z \geq 0) \\ 0 & (z \leq 0) \end{cases}$$

$$\phi_f(x, \ell) = (z^{(1)} = W^{(1)}x) \wedge \forall i. \left[ (z_i^{(1)} \leq 0 \wedge x_i^{(2)} = 0) \vee (z_i^{(1)} \geq 0 \wedge x_i^{(2)} = z_i^{(1)}) \right]$$



$$g(z) = \begin{cases} z & (z \geq 0) \\ 0 & (z \leq 0) \end{cases}$$

$$\phi_f(x, \ell) = (z^{(1)} = W^{(1)}x) \wedge (z^{(2)} = W^{(2)}x^{(2)}) \\ \wedge \forall i. \left[ (z_i^{(1)} \leq 0 \wedge x_i^{(2)} = 0) \vee (z_i^{(1)} \geq 0 \wedge x_i^{(2)} = z_i^{(1)}) \right]$$



$$g(z) = \begin{cases} z & (z \geq 0) \\ 0 & (z \leq 0) \end{cases}$$

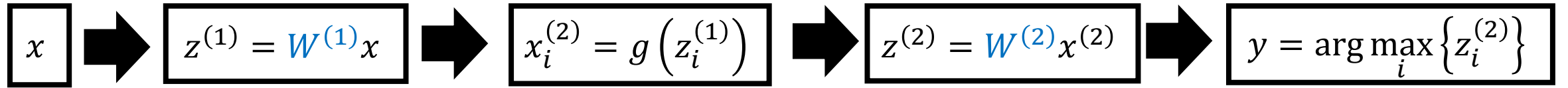
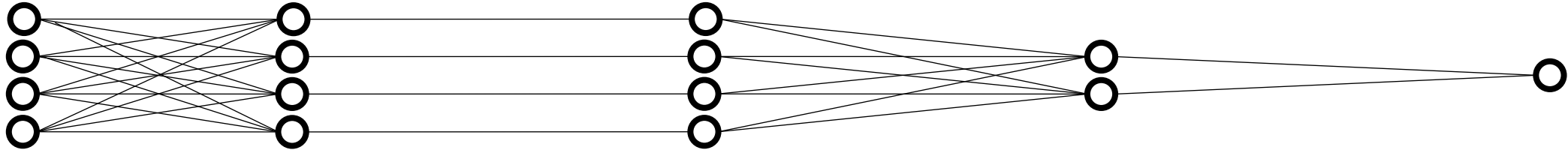
$$\begin{aligned} \phi_f(x, \ell) = & (z^{(1)} = W^{(1)}x) \wedge (z^{(2)} = W^{(2)}x^{(2)}) \\ & \wedge \forall i. \left[ (z_i^{(1)} \leq 0 \wedge x_i^{(2)} = 0) \vee (z_i^{(1)} \geq 0 \wedge x_i^{(2)} = z_i^{(1)}) \right] \\ & \wedge \forall i. \left[ z_\ell^{(2)} \geq z_i^{(2)} \right] \end{aligned}$$





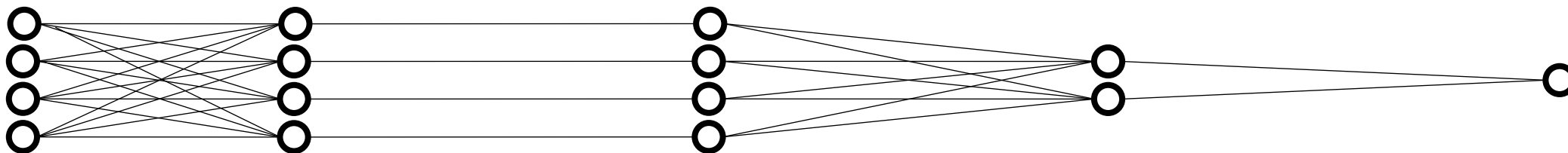
$$\phi_f(x, \ell) = \left( z^{(1)} = \mathbf{1} \wedge \forall i. \left[ \left( z_i^{(1)} \leq 0 \wedge x_i^{(2)} = z_i^{(1)} \right) \right] \right)$$

$g(\dots)$



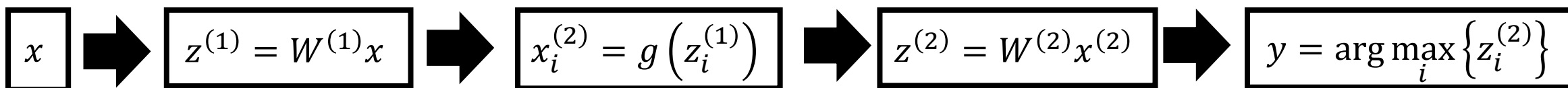
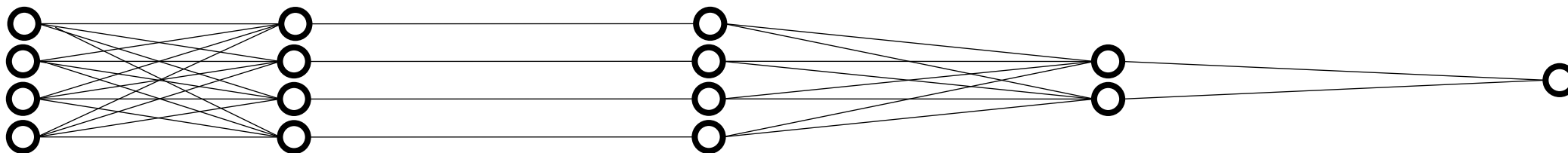
$$g(z) = \begin{cases} z & (z \geq 0) \\ 0 & (z \leq 0) \end{cases}$$

$$\begin{aligned} \phi_f(x, \ell) = & (z^{(1)} = W^{(1)}x) \wedge (z^{(2)} = W^{(2)}x^{(2)}) \\ & \wedge \forall i. \left[ (z_i^{(1)} \leq 0 \wedge x_i^{(2)} = 0) \vee (z_i^{(1)} \geq 0 \wedge x_i^{(2)} = z_i^{(1)}) \right] \\ & \wedge \forall i. \left[ z_\ell^{(2)} \geq z_i^{(2)} \right] \end{aligned}$$



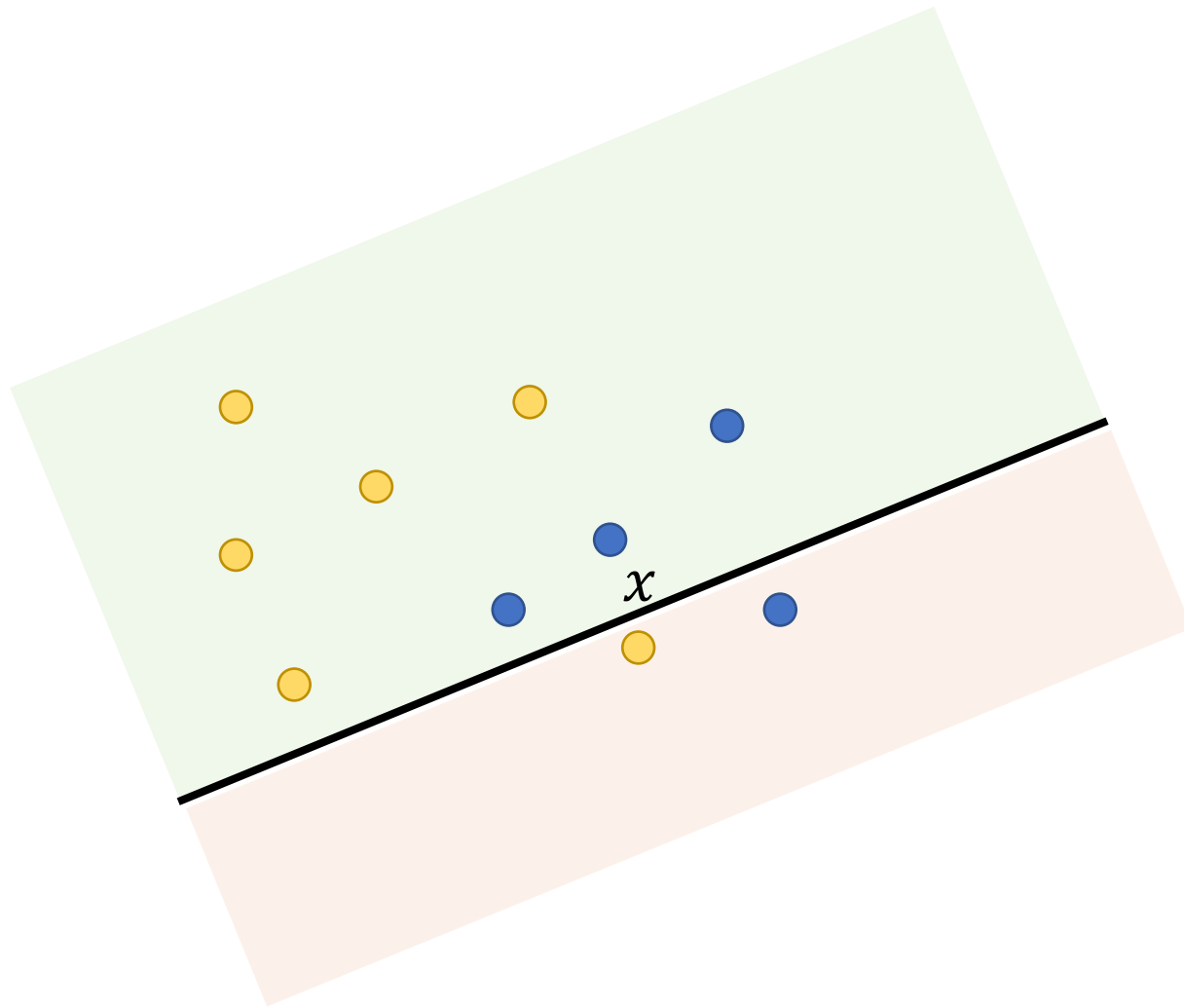
$$g(z) = \begin{cases} z & (z \geq 0) \\ 0 & (z \leq 0) \end{cases}$$

$$\begin{aligned} \phi_f(x, \ell) = & (z^{(1)} = W^{(1)}x) \wedge (z^{(2)} = W^{(2)}x^{(2)}) \\ & \wedge \forall i. \left[ \left( \cancel{z_i^{(1)} \leq 0 \wedge x_i^{(2)} = 0} \vee \left( z_i^{(1)} \geq 0 \wedge x_i^{(2)} = z_i^{(1)} \right) \right) \right] \\ & \wedge \forall i. \left[ z_\ell^{(2)} \geq z_i^{(2)} \right] \end{aligned}$$

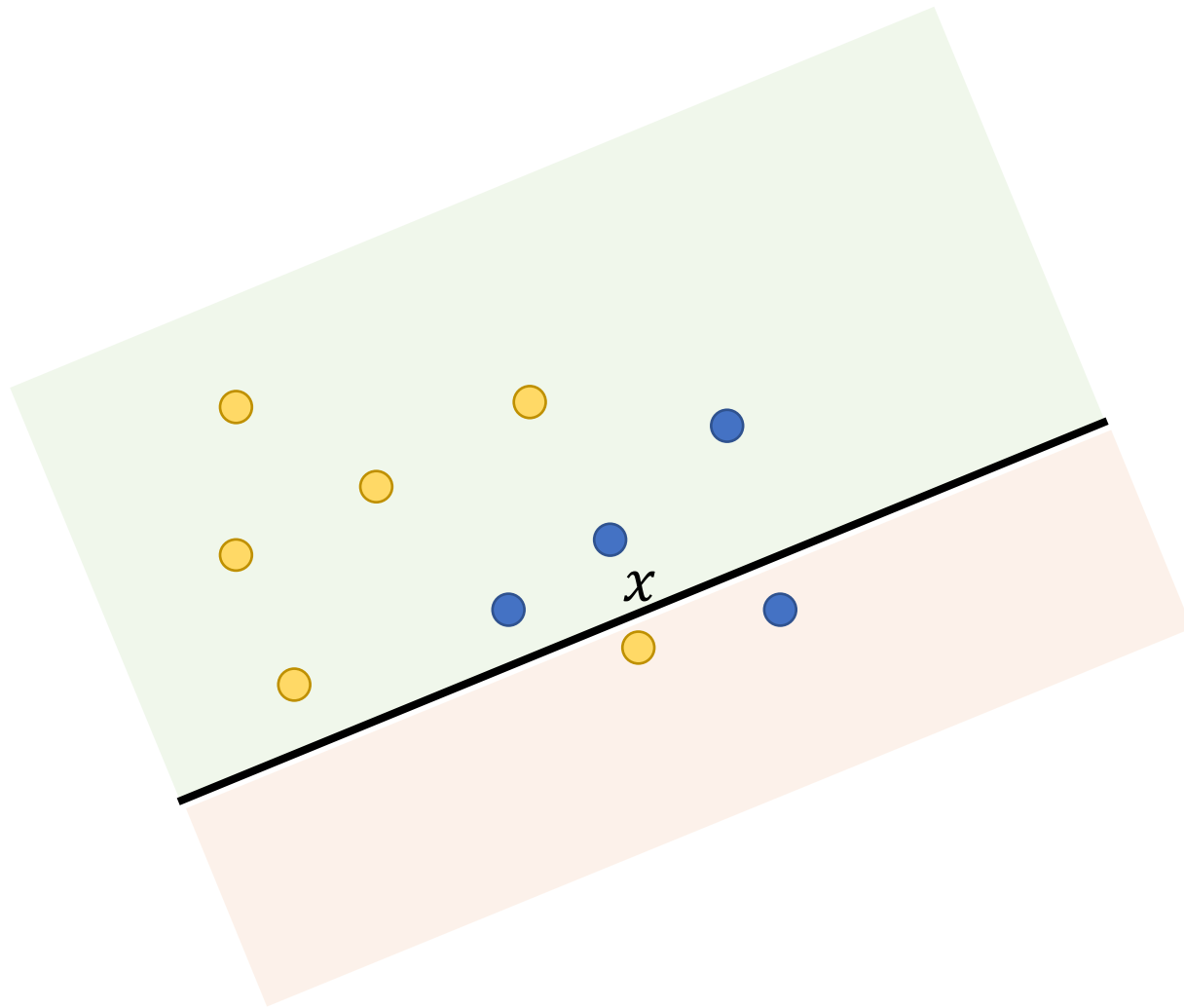


$$g(z) = \begin{cases} z & (z \geq 0) \\ 0 & (z \leq 0) \end{cases}$$

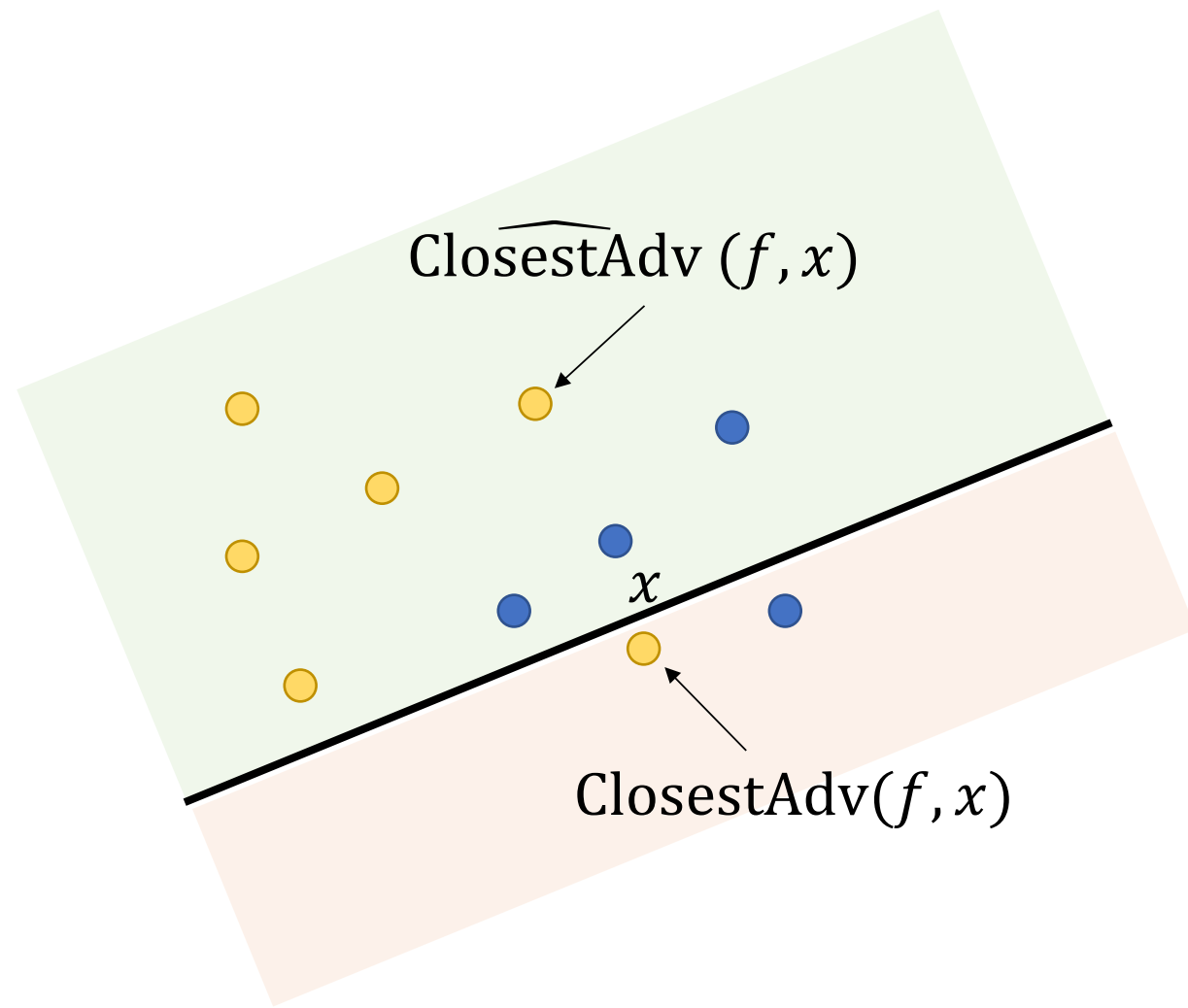
$$\begin{aligned} \phi_f(x, \ell) = & (z^{(1)} = W^{(1)}x) \wedge (z^{(2)} = W^{(2)}x^{(2)}) \\ & \wedge \forall i. [z_i^{(1)} \geq 0 \wedge x_i^{(2)} = z_i^{(1)}] \\ & \wedge \forall i. [z_\ell^{(2)} \geq z_i^{(2)}] \end{aligned}$$



$$\left( z_i^{(1)} \leq 0 \wedge x_i^{(2)} = 0 \right) \vee \left( z_i^{(1)} \geq 0 \wedge x_i^{(2)} = z_i^{(1)} \right)$$



$$\left( z_i^{(1)} \leq 0 \wedge x_i^{(2)} = 0 \right) \vee \left( z_i^{(1)} \geq 0 \wedge x_i^{(2)} = z_i^{(1)} \right)$$



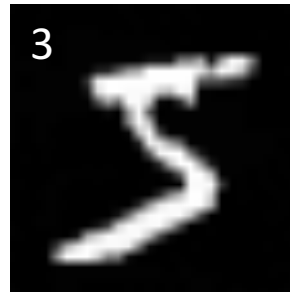
$$\left( z_i^{(1)} \leq 0 \wedge x_i^{(2)} = 0 \right) \vee \left( z_i^{(1)} \geq 0 \wedge x_i^{(2)} = z_i^{(1)} \right)$$

# MNIST

original ( $x$ )



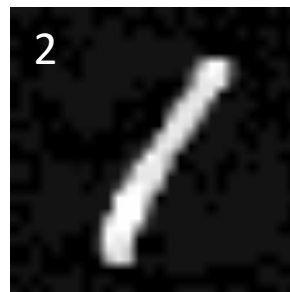
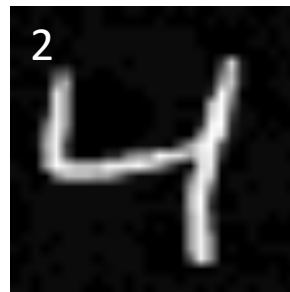
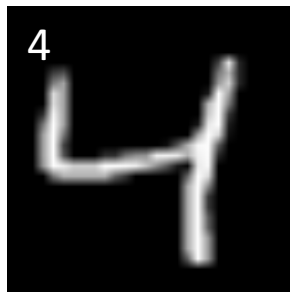
mislabeled ( $x'$ )



$10|x - x'|$



$100|x - x'|$





# MNIST

Neural Net	Accuracy (%)	Adversarial Frequency (%)	
		Baseline	Our Algo.
LeNet (Original)	99.08	1.32	7.15
Baseline ( $T = 1$ )	99.14	1.02	6.89
Baseline ( $T = 2$ )	99.15	0.99	6.97
Our Algo. ( $T = 1$ )	99.17	1.18	5.40
Our Algo. ( $T = 2$ )	99.23	1.12	5.03

$\epsilon = 20$  pixels

# Fine Tuning (Goodfellow 2015)

**input:**  $X_{\text{train}}$

**while true:**

$f$  = train neural network on  $X_{\text{train}}$

$X_{\text{adv}}$  = find adversarial examples

$X_{\text{train}} = X_{\text{train}} \cup X_{\text{adv}}$

# MNIST

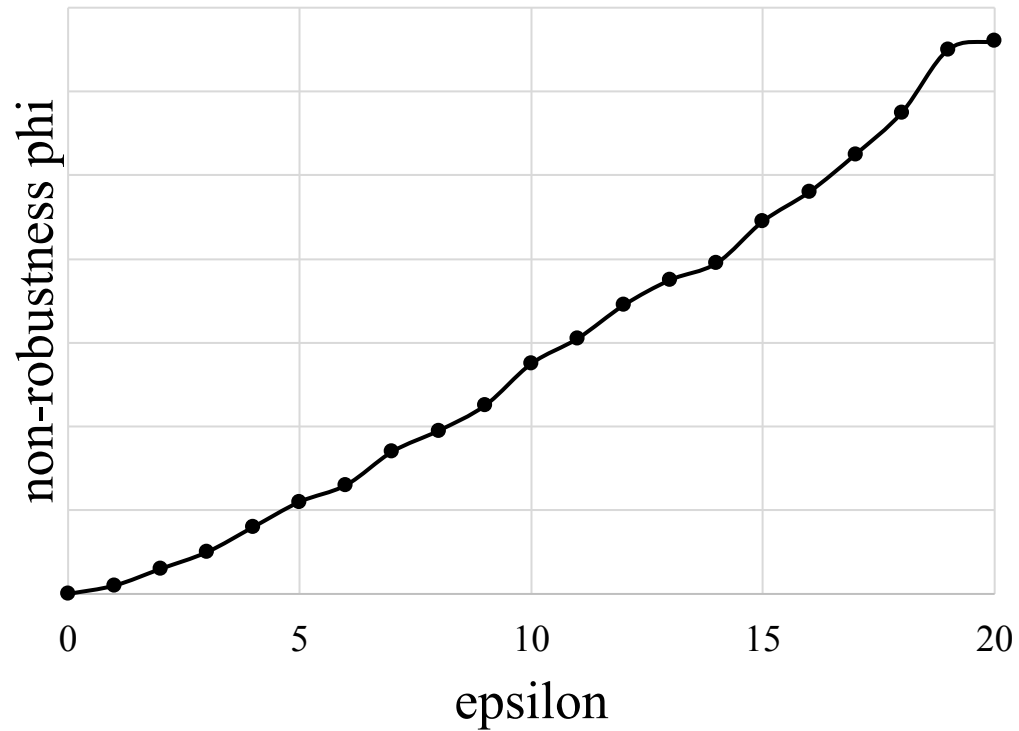
Neural Net	Accuracy (%)	Adversarial Frequency (%)	
		Baseline	Our Algo.
LeNet (Original)	99.08	1.32	7.15
Baseline ( $T = 1$ )	99.14	1.02	6.89
Baseline ( $T = 2$ )	99.15	0.99	6.97
Our Algo. ( $T = 1$ )	99.17	1.18	5.40
Our Algo. ( $T = 2$ )	99.23	1.12	5.03

$\epsilon = 20$  pixels

# Improving Robustness?

# Improving Robustness?

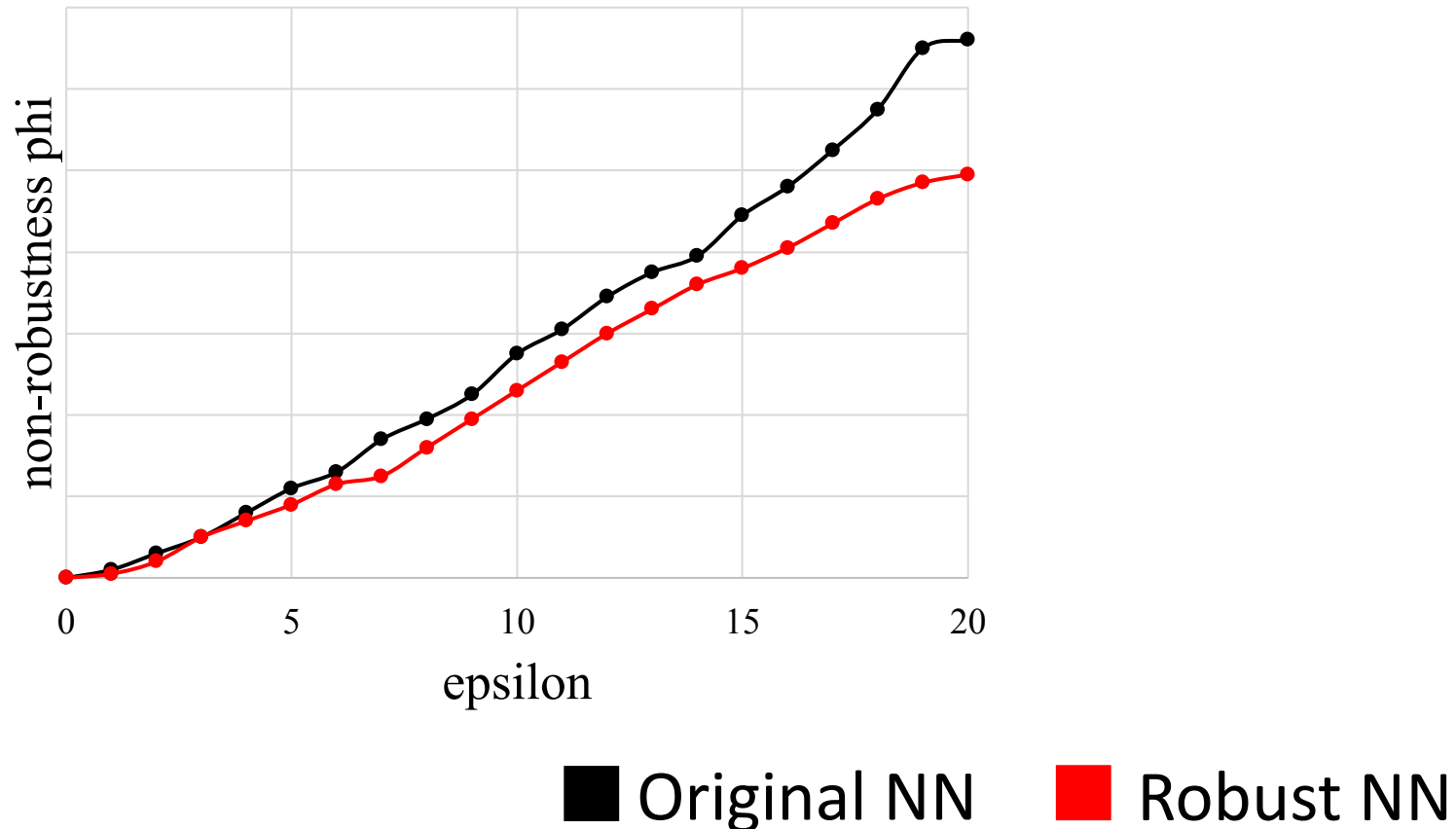
Algorithm's Own Metric



■ Original NN

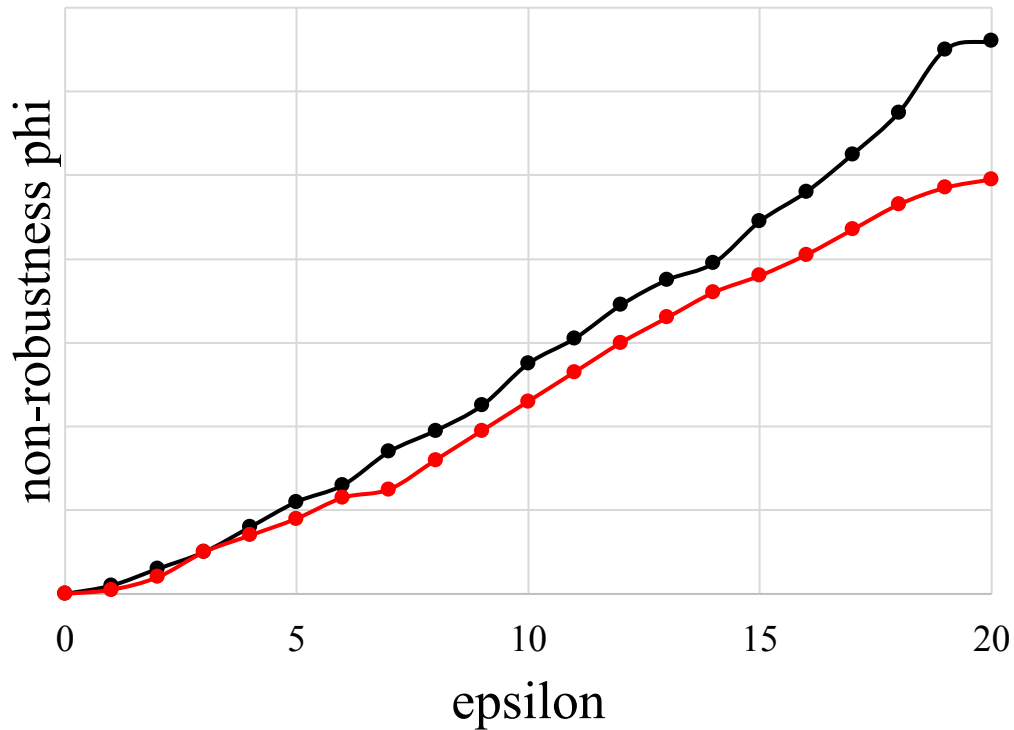
# Improving Robustness?

Algorithm's Own Metric

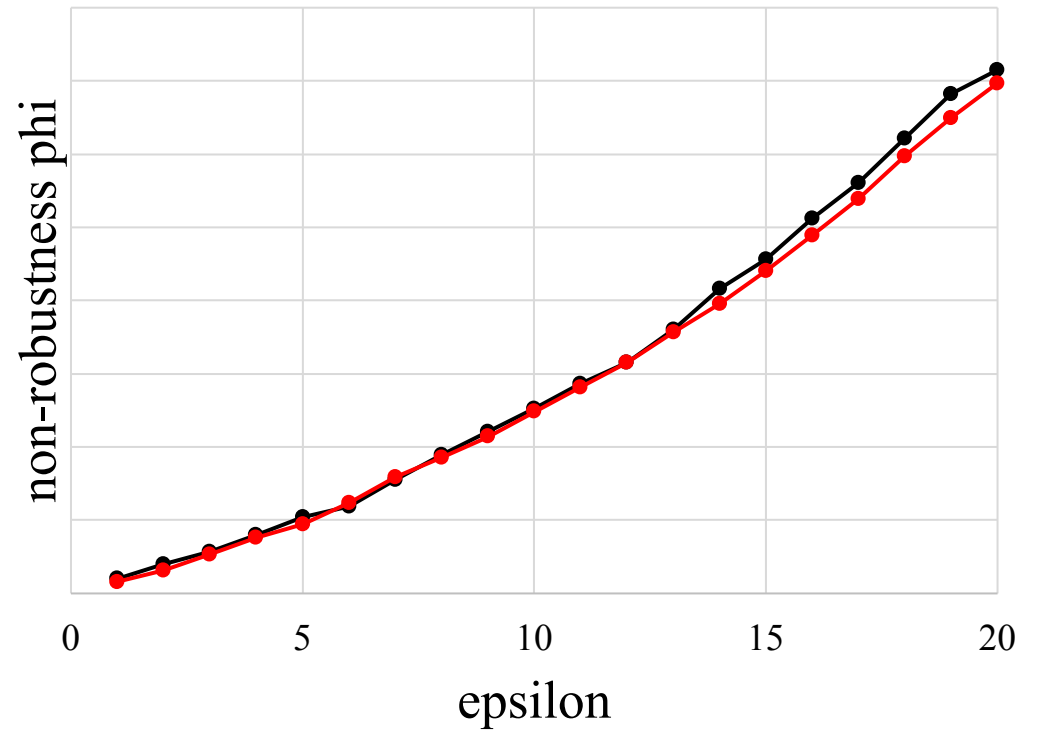


# Improving Robustness?

## Algorithm's Own Metric

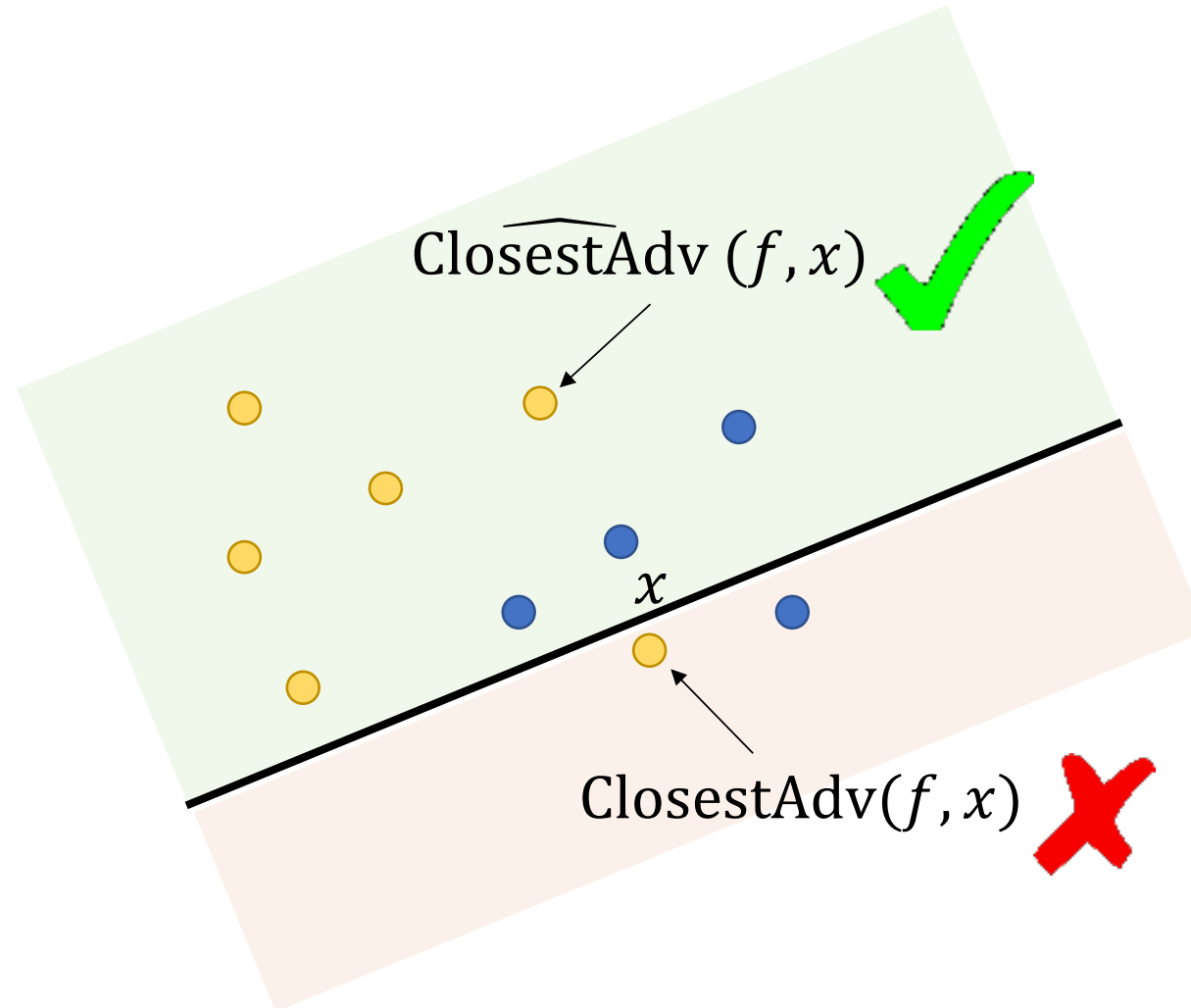


## Our Metric



■ Original NN    ■ Robust NN

# Improving Robustness?



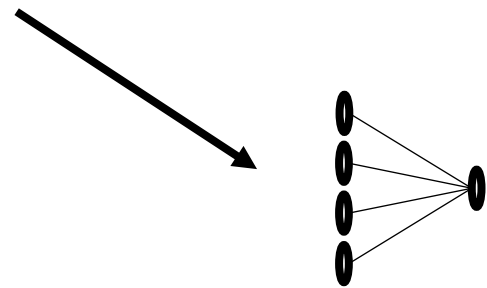
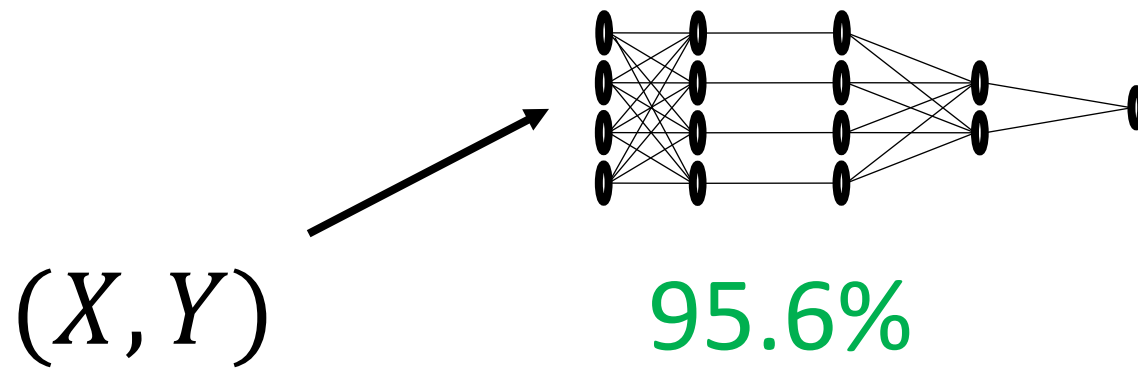


# Related Work

- Finding adversarial examples
  - Szegedy et al. 2014, Goodfellow et al. 2015, ...
- Verifying robustness
  - Reluplex: SMT solver for theory of ReLU (Katz et al. 2017)
  - DLV: Discretize search space (Huang et al. 2017)
  - Mixed integer programs (Tjeng et al. 2017)
  - AI<sup>2</sup>: Abstract interpretation (Gehr et al. 2018)

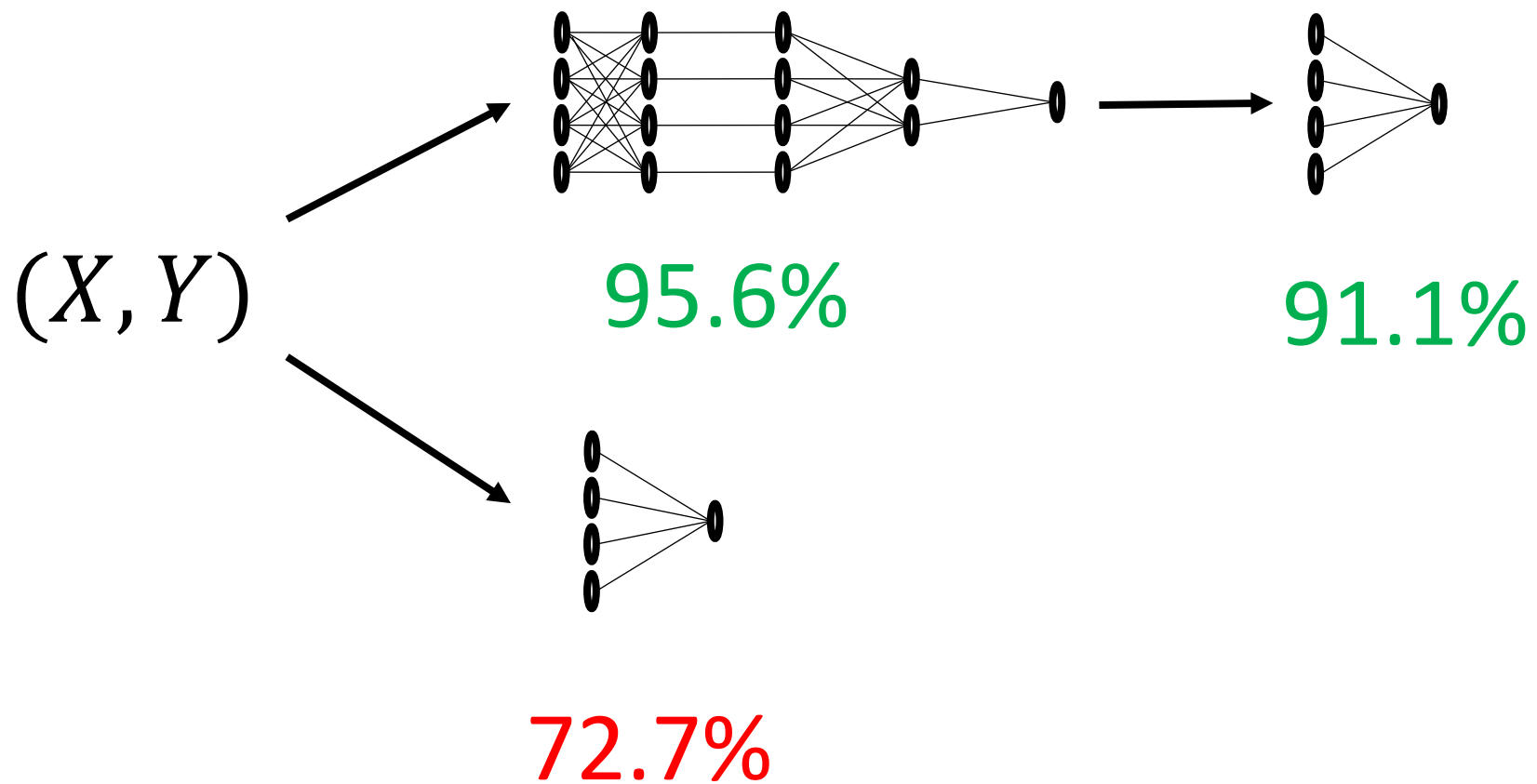
# Verifiable Models via Model Compression (Ongoing)

Osbert Bastani, Evan Pu, Armando Solar-Lezama

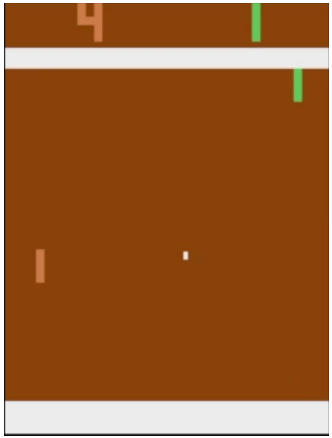


72.7%

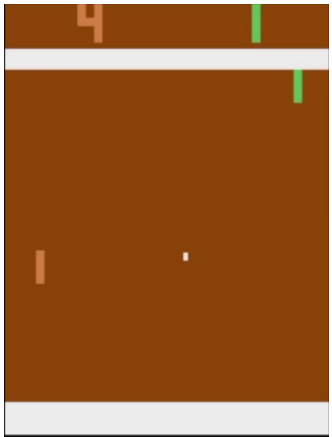
(Ba et al. 2014)



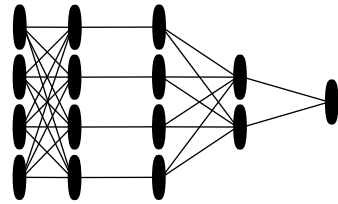
(Ba et al. 2014)



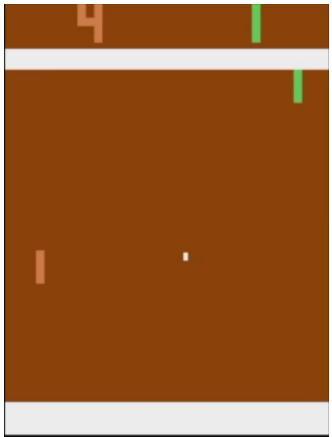
**control  
problem**



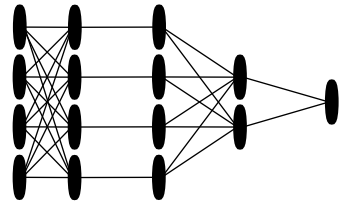
**control  
problem**



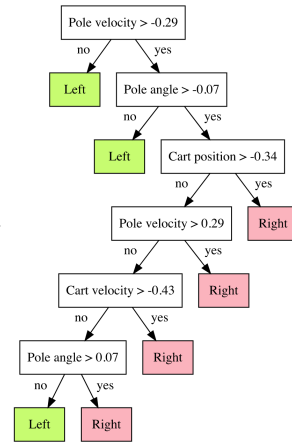
**deep RL**



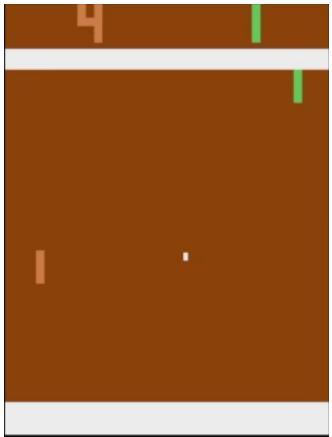
**control  
problem**



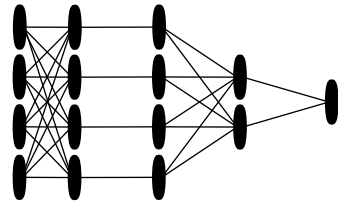
**deep RL**



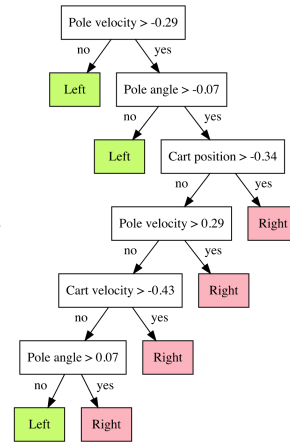
**verifiable  
controller**



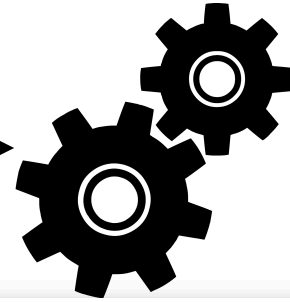
**control  
problem**



**deep RL**

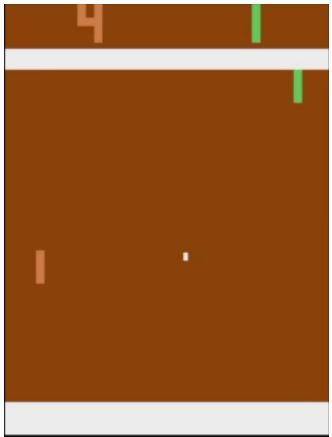


**verifiable  
controller**

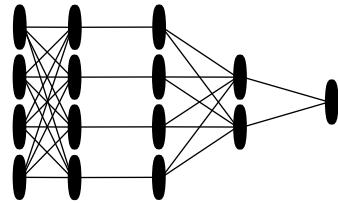


**verification**

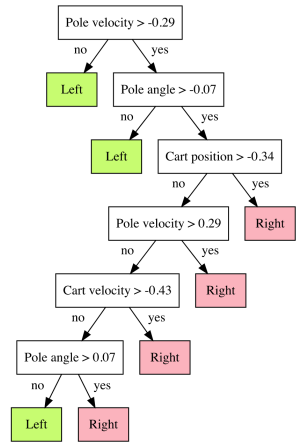




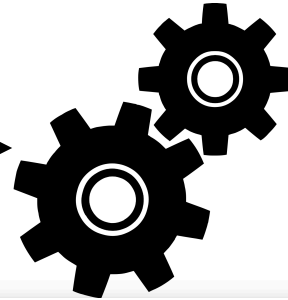
**control problem**



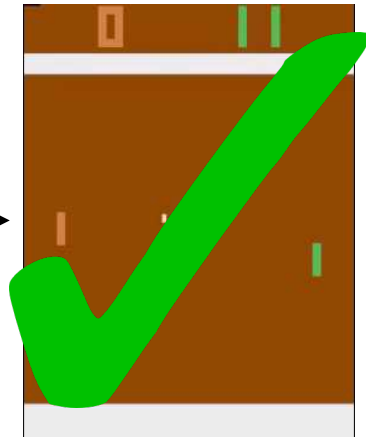
**deep RL**



**verifiable controller**



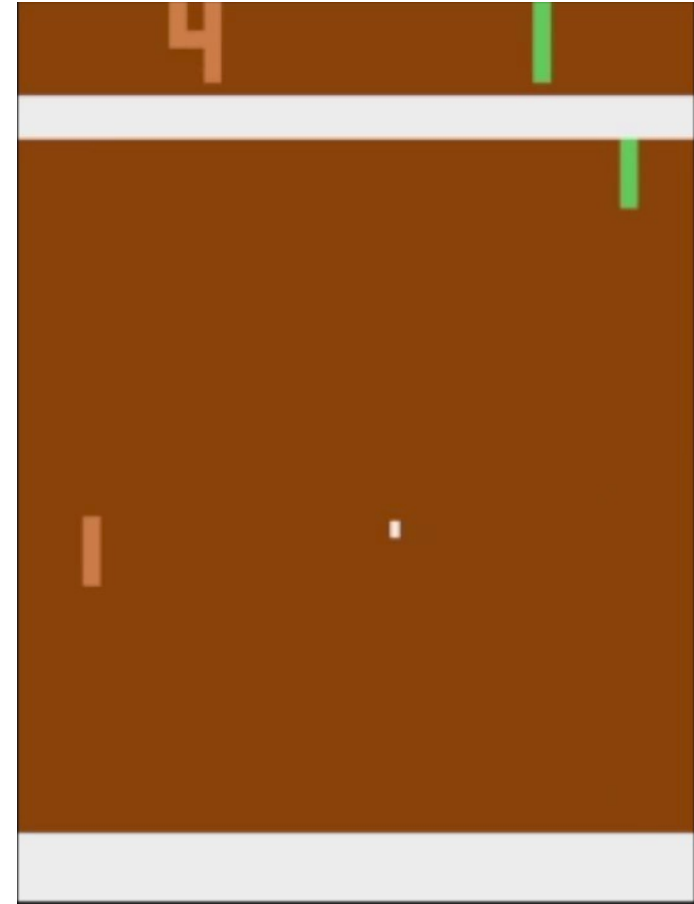
**verification**



**provably robust controller**

- **Problem setup**

- State space: 12 dimensional
- Action space: Move paddle up/down

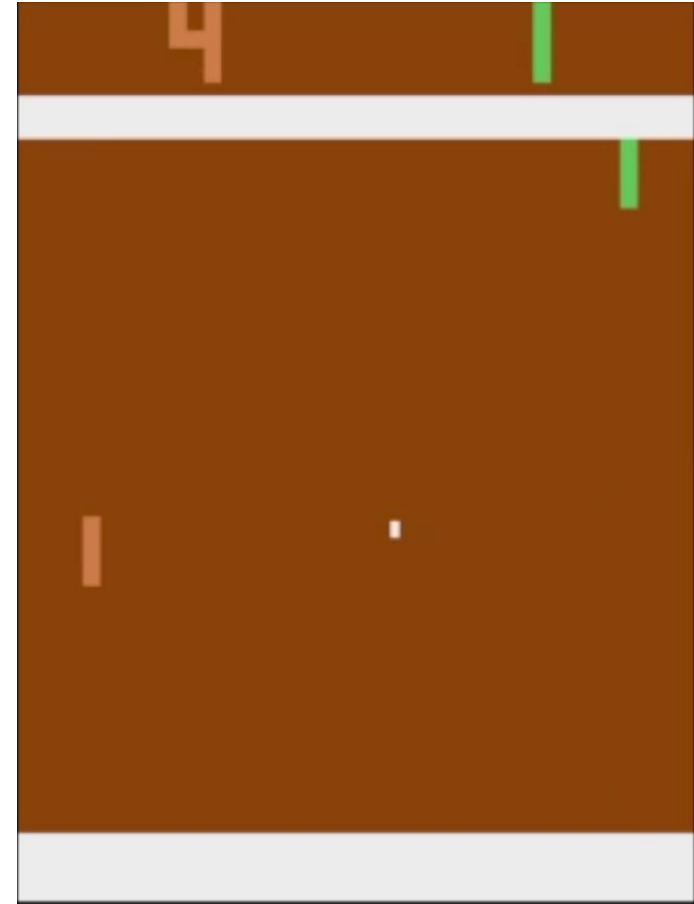


- **Problem setup**

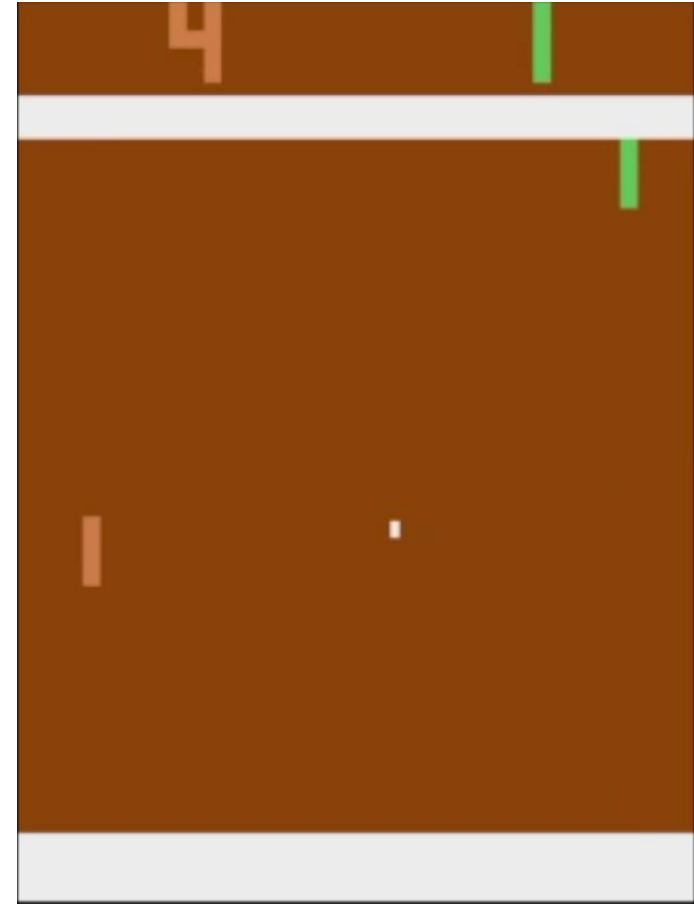
- State space: 12 dimensional
- Action space: Move paddle up/down

- **Policy**

- Decision tree
- 4354 leaf nodes



- **Problem setup**
  - State space: 12 dimensional
  - Action space: Move paddle up/down
- **Policy**
  - Decision tree
  - 4354 leaf nodes
- **Verification (Robust at one point)**
  - 4354 linear program calls
  - 52.8 seconds



# Verification for learning-based systems

- Robustness, stability, etc.