Machine Learning Systems

Lecture 12: Security of ML Systems (Adversarial ML)

Pooyan Jamshidi

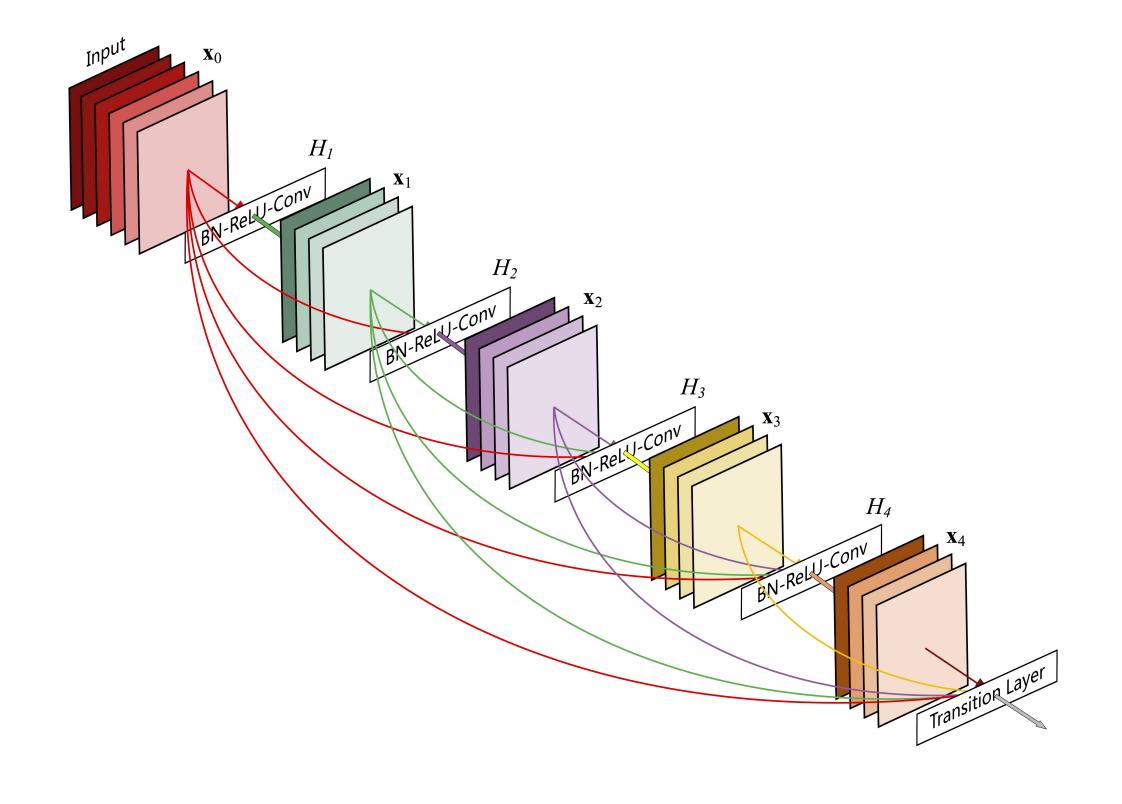


So what this talk is about?

The Security of

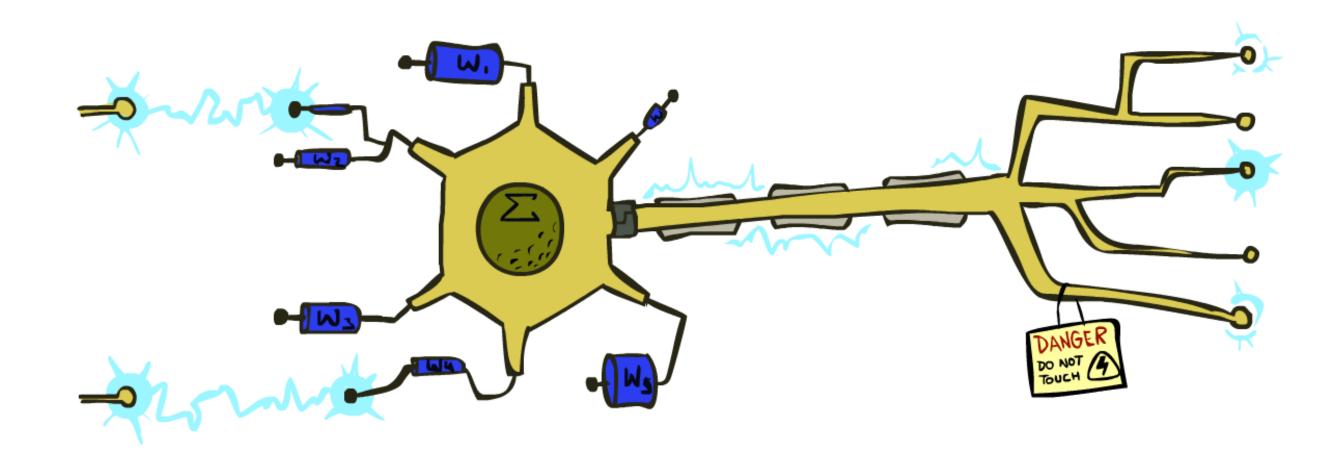
Machine Learning

Deep

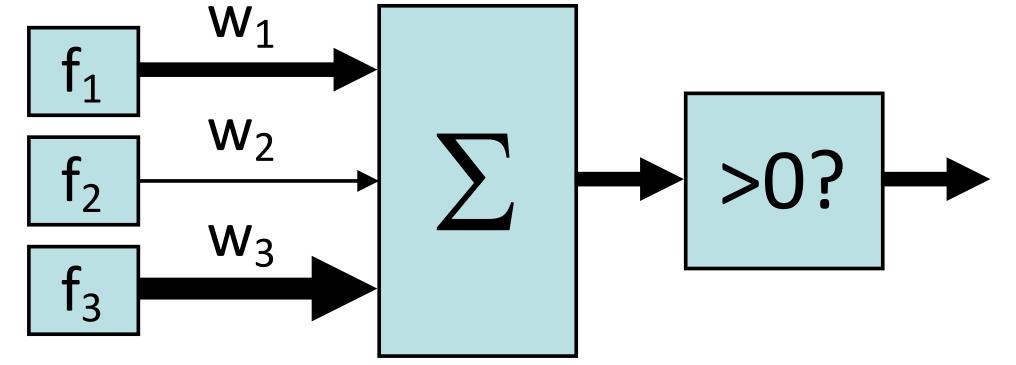


A quick recap about Supervised Machine Learning

Linear Classifiers

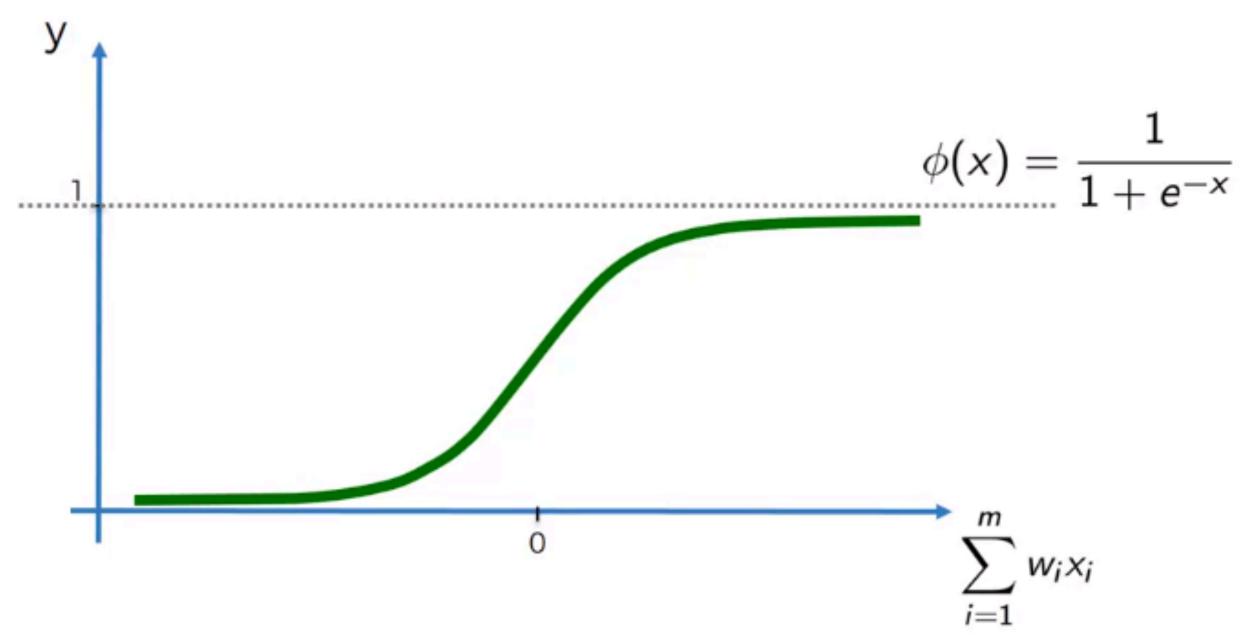


$$activation_{w}(x) = \sum_{i} w_{i} \cdot f_{i}(x) = w \cdot f(x)$$



How to get probabilistic decisions?

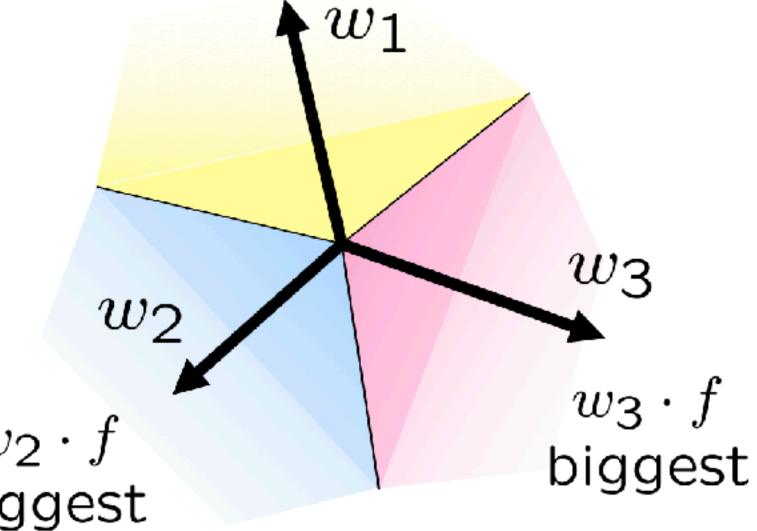
- Activation: $z = w \cdot f(x)$
- If z very positive -> want probability going to 1
- If z very negative -> want probability going to 0

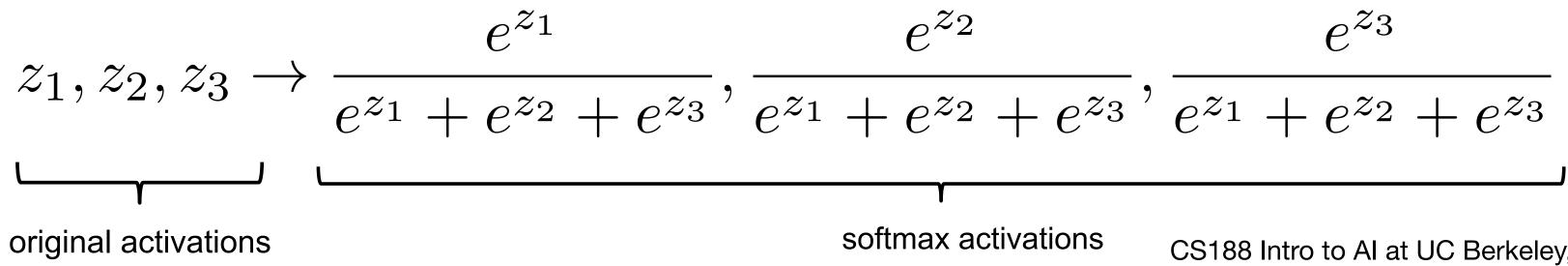


Multiclass Logistic Regression

- Multi-class linear classification
 - A weight vector for each class: w_v
 - Score (activation) of a class y: $w_v \cdot f(x)$
 - Prediction w/highest score wins: $y = argmax_v w_v \cdot f(x)$ biggest
- How to make the scores into probabilities?

 $w_1 \cdot f$ biggest





Best w?

Maximum likelihood estimation:

$$\max_{w} \ ll(w) = \max_{w} \ \sum_{i} \log P(y^{(i)}|x^{(i)};w)$$

$$P(y^{(i)}|x^{(i)};w) = \frac{e^{w_{y^{(i)}} \cdot f(x^{(i)})}}{\sum_{y} e^{w_{y} \cdot f(x^{(i)})}}$$

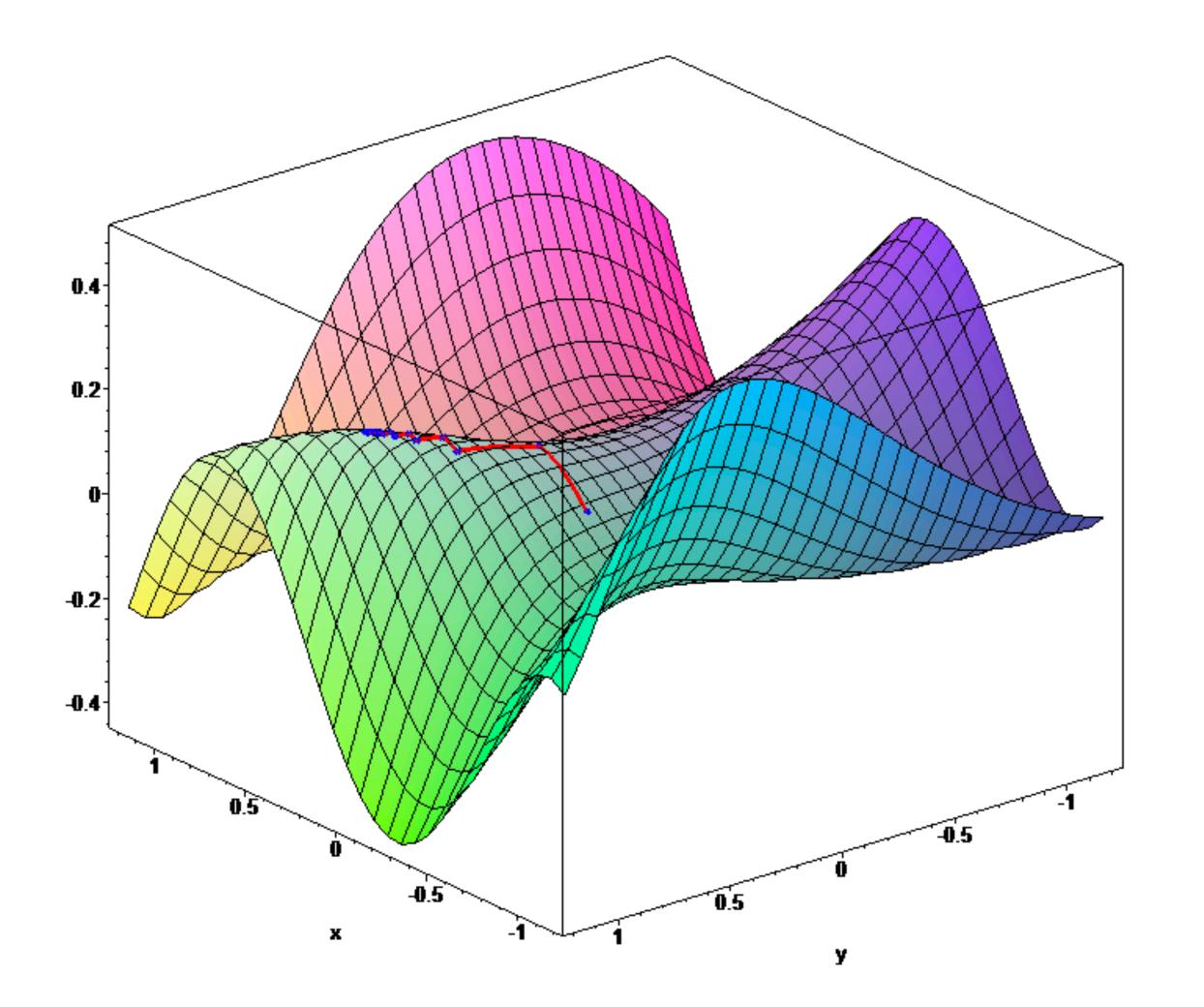
How do we solve the optimization problem?

$$\max_{w} ll(w) = \max_{w} \sum_{i} \log P(y^{(i)}|x^{(i)};w)$$

$$g(w)$$

Gradient Ascent

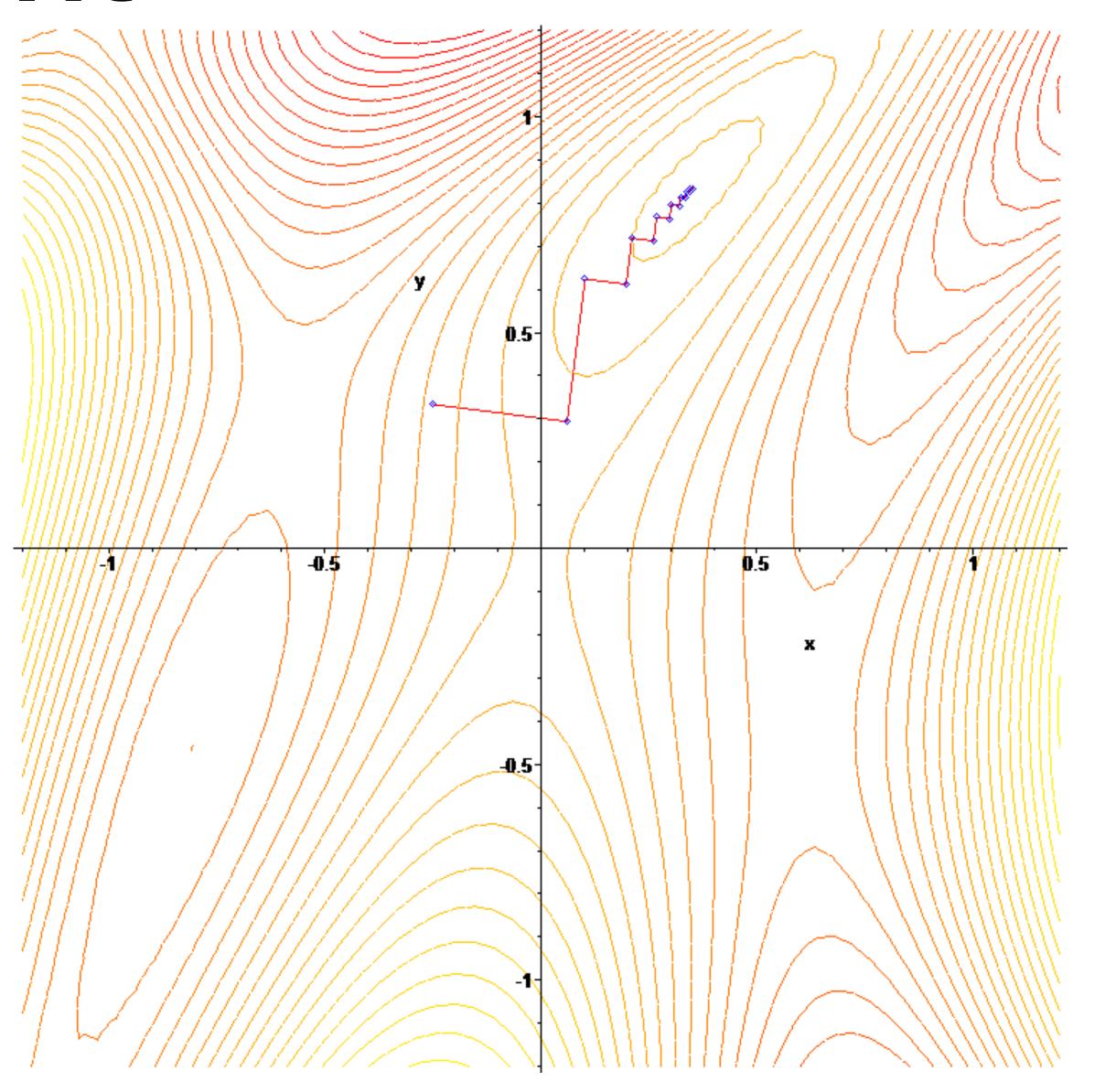
$$w \leftarrow w + \alpha * \nabla g(w)$$



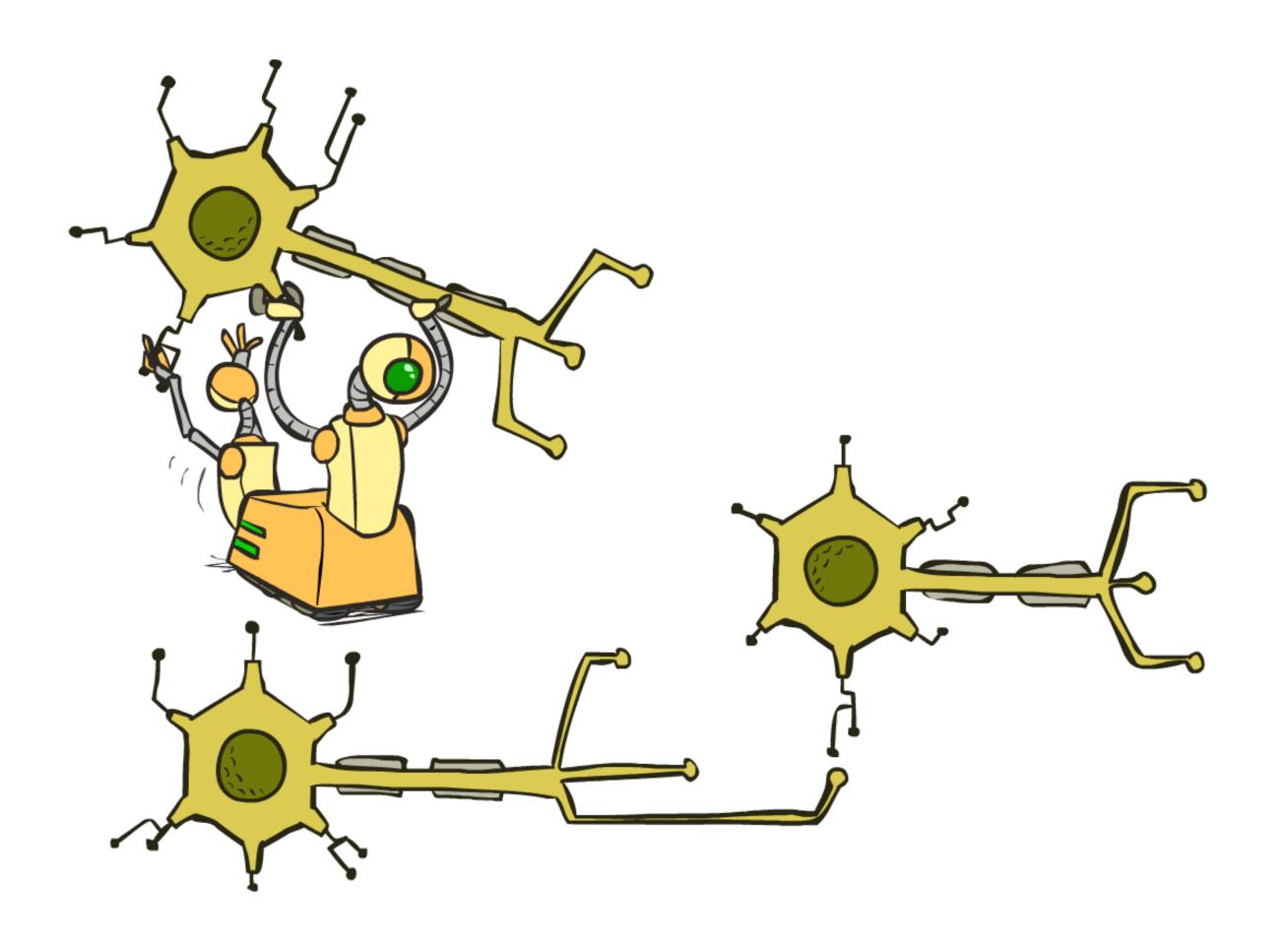
Gradient Ascent

$$w \leftarrow w + \alpha * \nabla g(w)$$

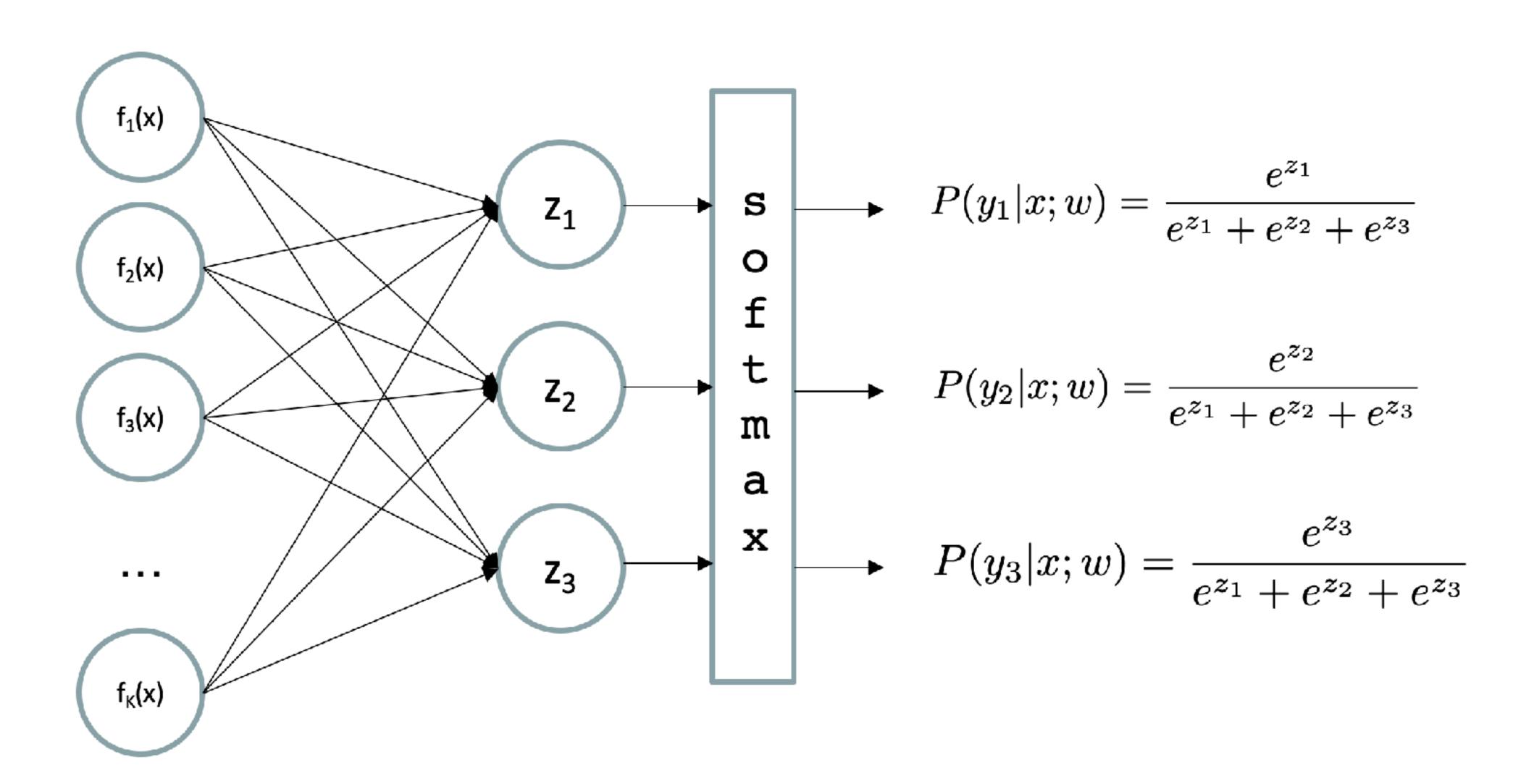
$$\nabla g = \begin{bmatrix} \frac{\partial g}{\partial w_1} \\ \frac{\partial g}{\partial w_2} \\ \vdots \\ \frac{\partial g}{\partial w_n} \end{bmatrix}$$



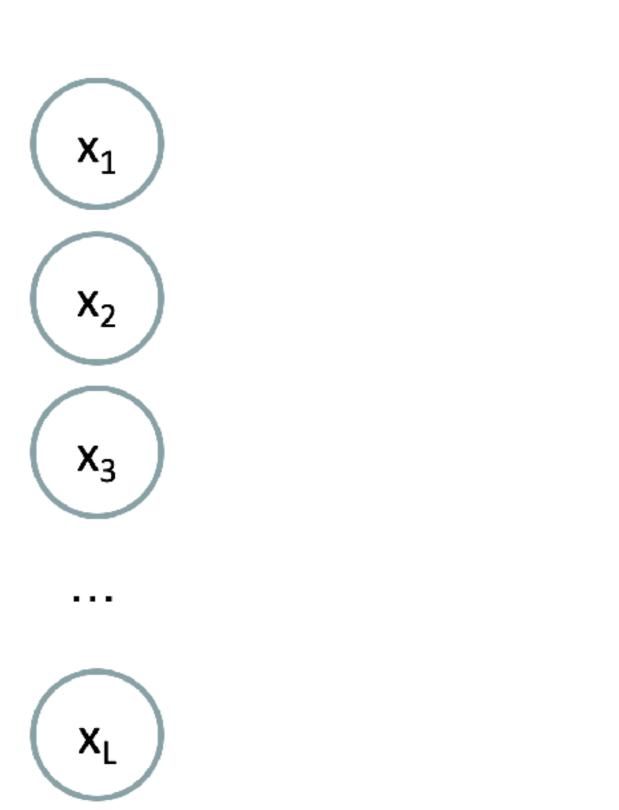
Neural Networks

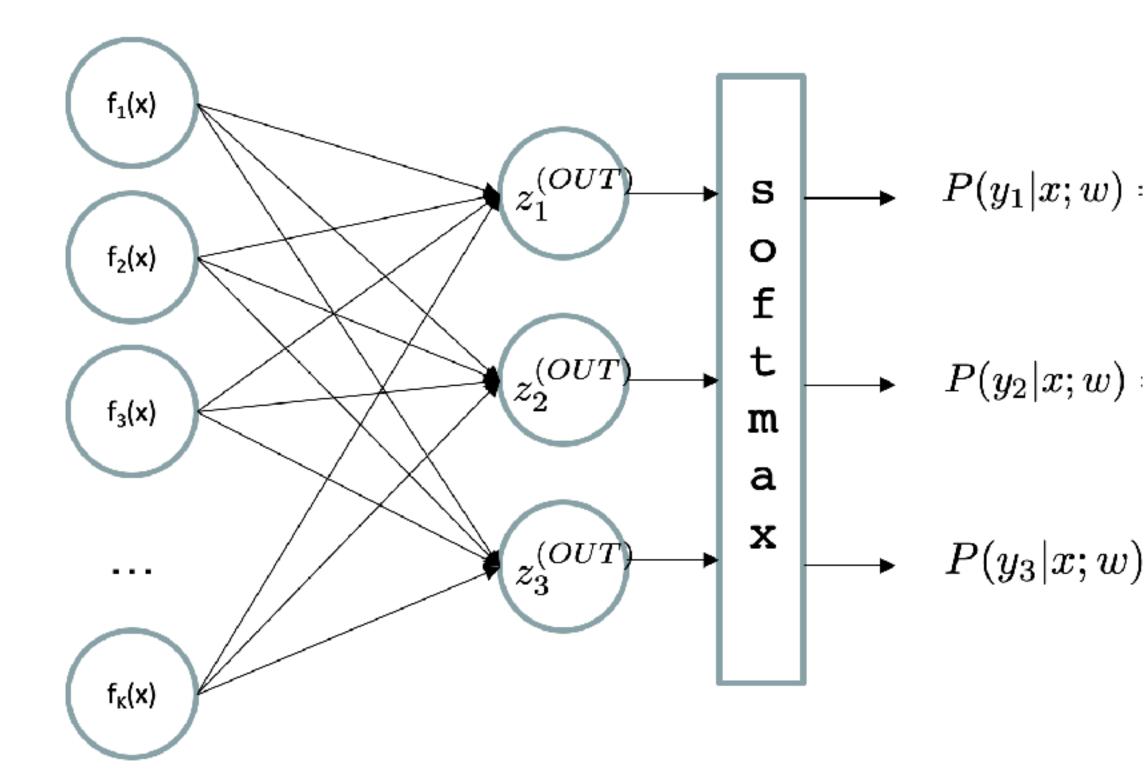


Multi-class Logistic Regression is a special case of neural network

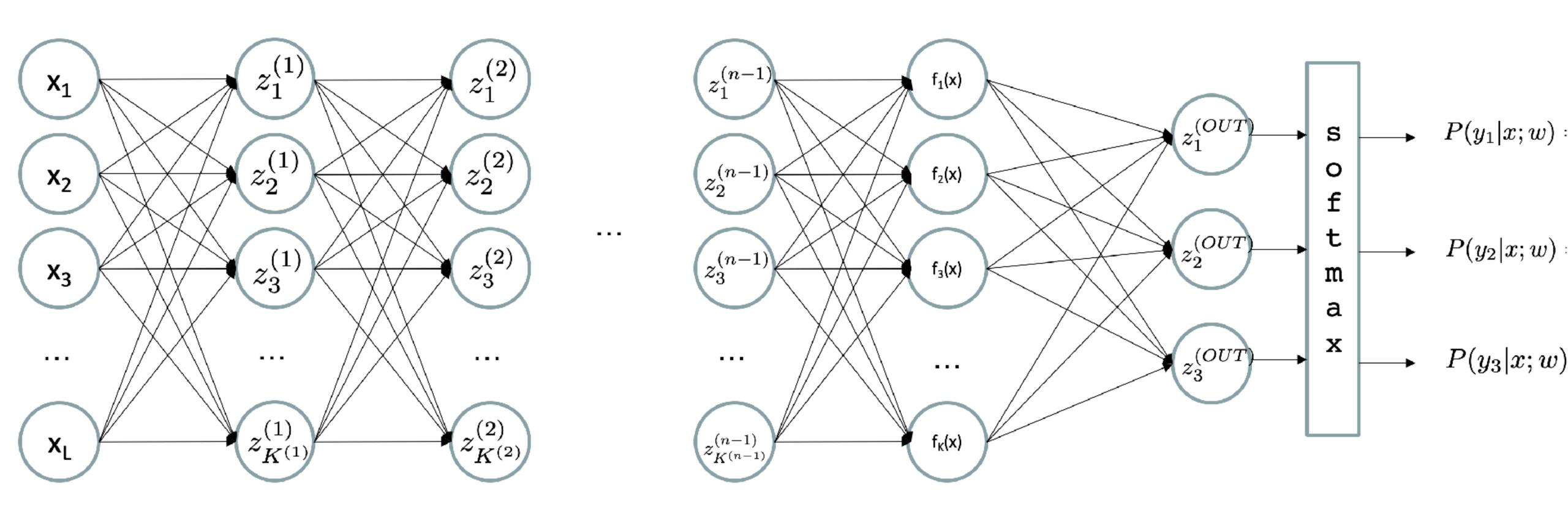


Deep Neural Network = Also learn the features!

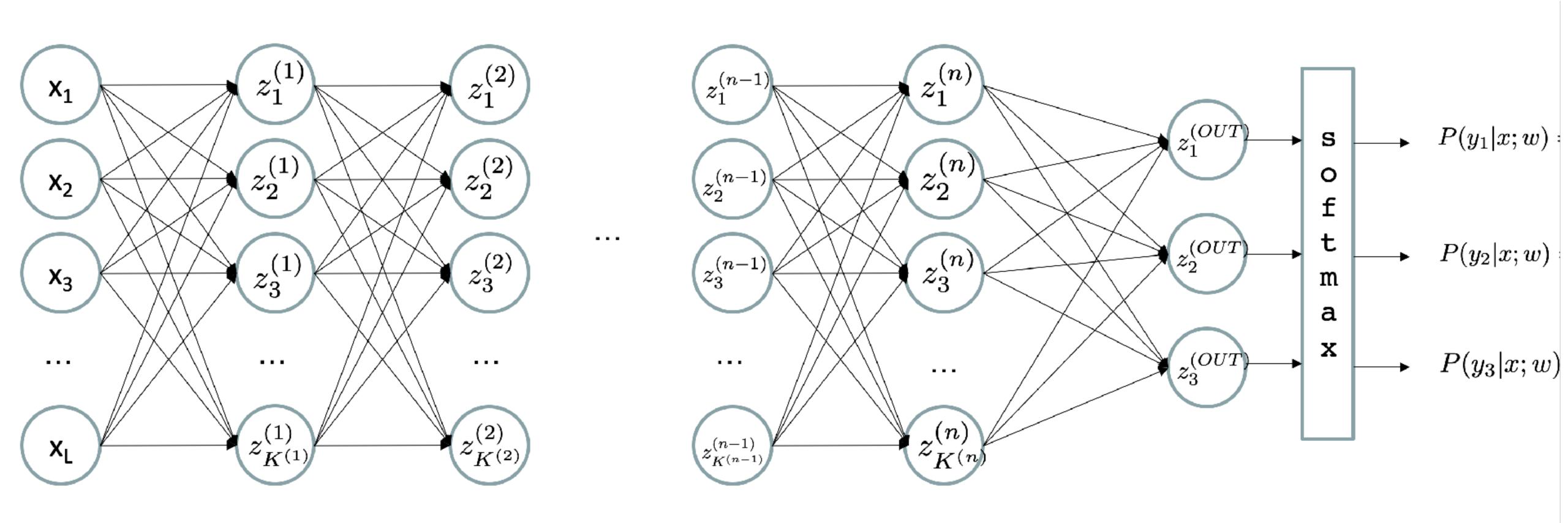




Deep Neural Network = Also learn the features!



Deep Neural Network = Also learn the features!



$$z_i^{(k)} = g(\sum_j W_{i,j}^{(k-1,k)} z_j^{(k-1)})$$

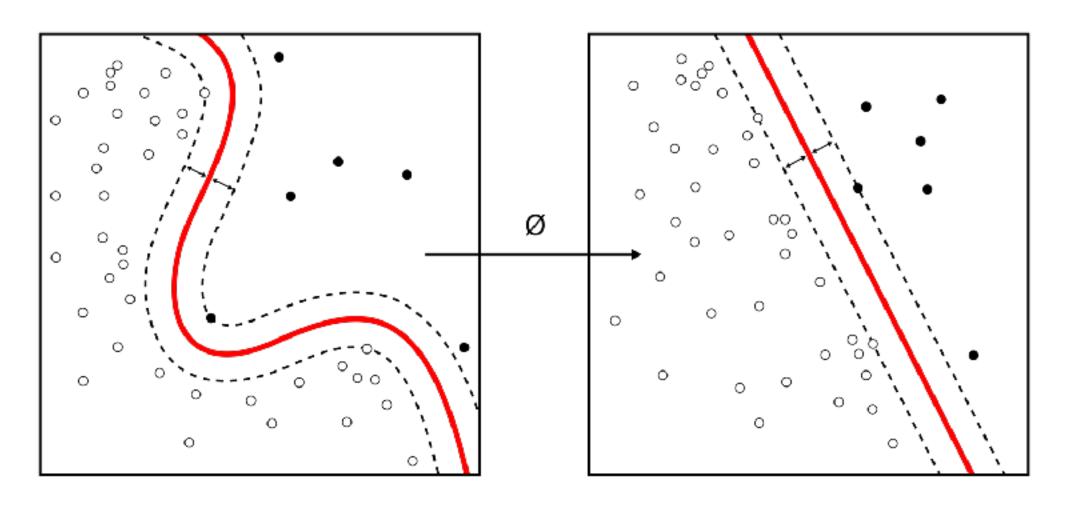
g = nonlinear activation function

Training the deep neural network is just like logistic regression

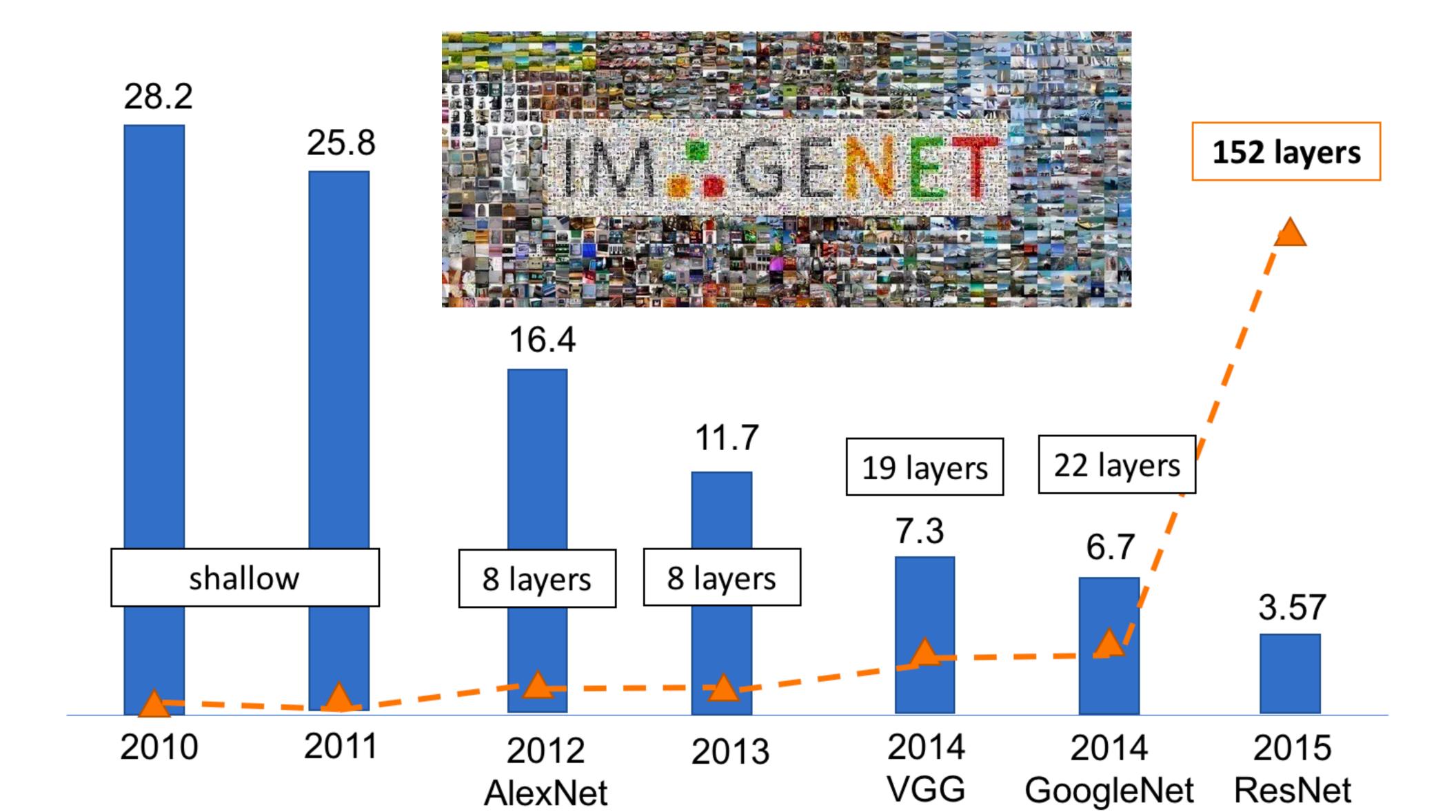
$$\max_{w} \ ll(w) = \max_{w} \ \sum_{i} \log P(y^{(i)}|x^{(i)};w)$$

- just run gradient ascent
- + stop when log likelihood of hold-out data starts to decrease

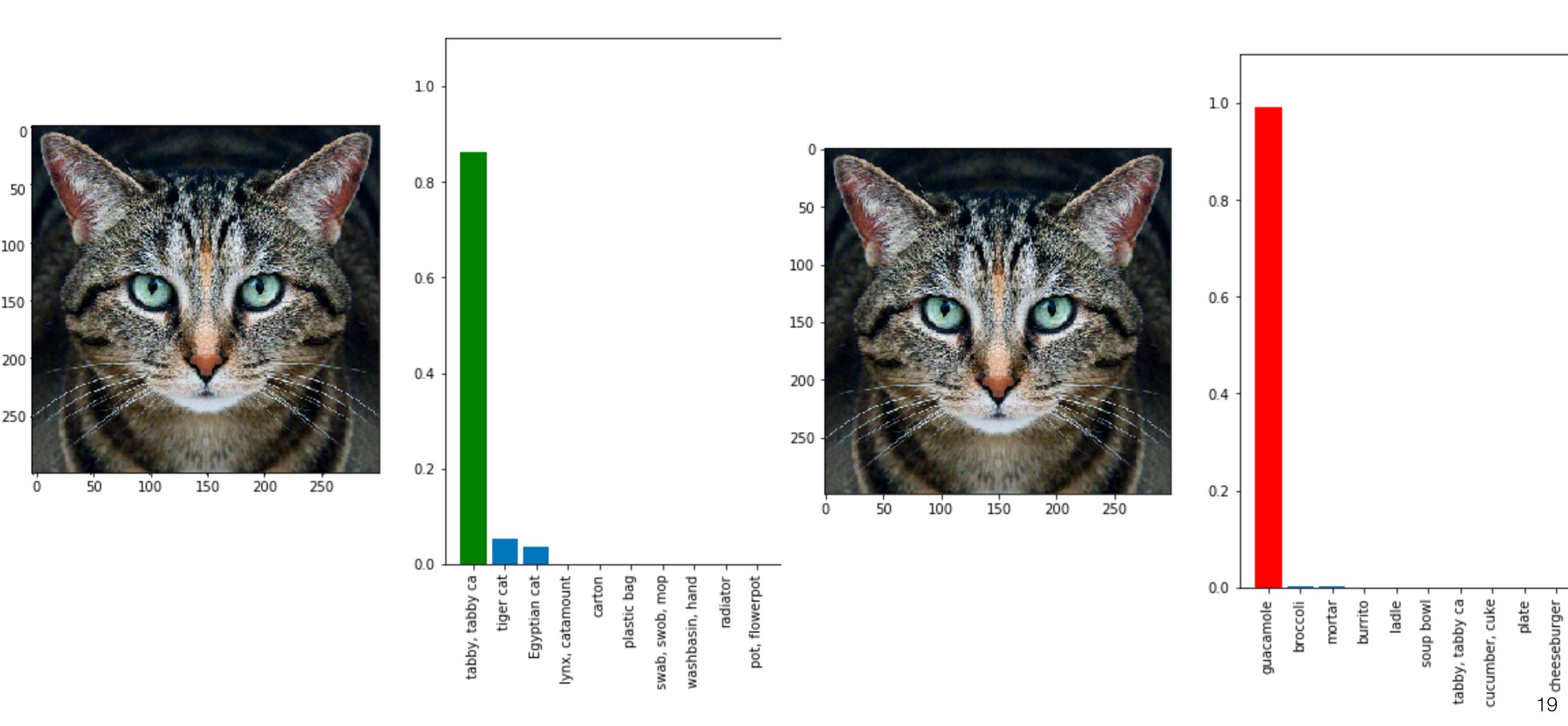
ML developed a rich theory to guide us here (and this was its only goal)



Machine Learning: The Success Story



Deep neural networks can be easily fooled



Neural networks can be tricked



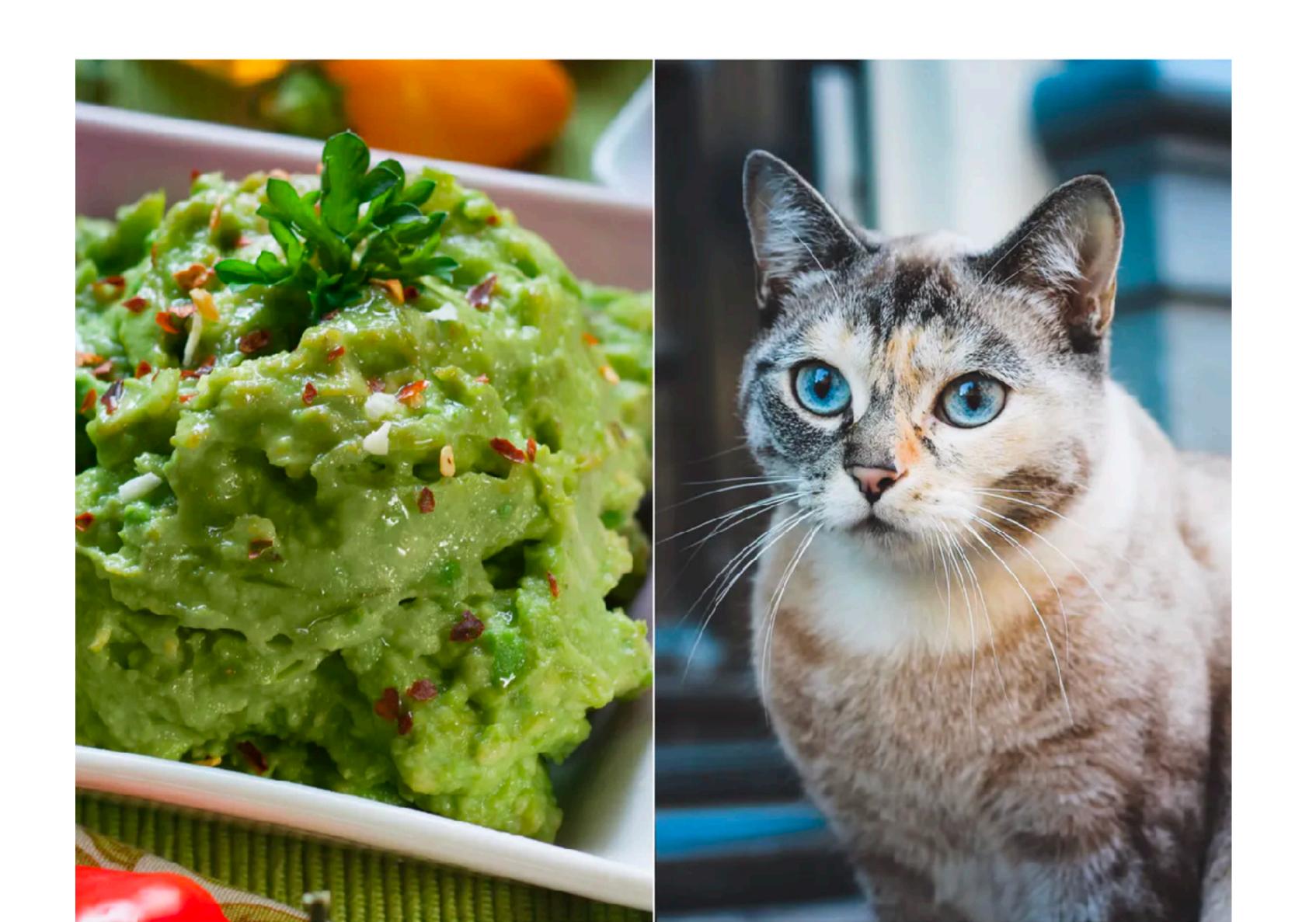
89% tabby cat

Adversarial Perturbation



98% guacamole

Yes, neural networks can be tricked that easily

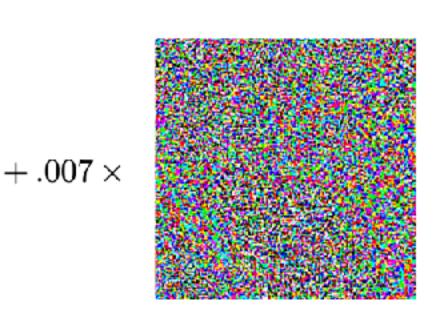


Adversarial Examples



57.7% confidence

 \boldsymbol{x} "panda"



 $sign(\nabla_{\boldsymbol{x}}J(\boldsymbol{\theta},\boldsymbol{x},y))$ "nematode" 8.2% confidence



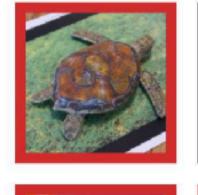
 $\epsilon \operatorname{sign}(\nabla_{\boldsymbol{x}} J(\boldsymbol{\theta}, \boldsymbol{x}, y))$ "gibbon" 99.3 % confidence



dog

[Goodfellow et al. 2014]: Imperceptible noise can fool DNN classifiers

[Engstrom, Tran, Tsipras, Schmidt, Madry 2018]: **Rotation + Translation can fool classifiers**





















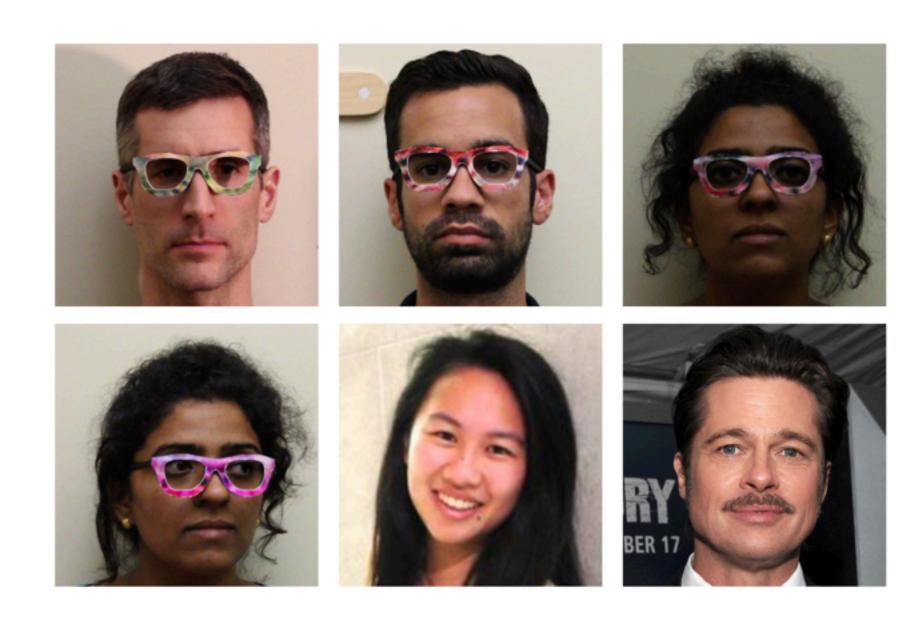




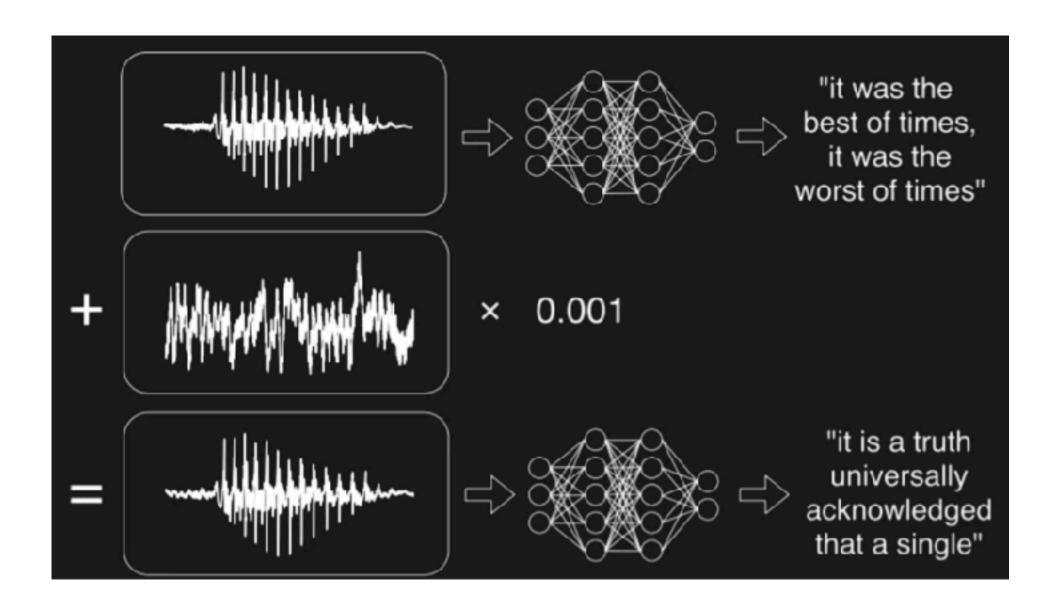
classified as other

[Athalye, Engstrom, Ilyas, Kwok 2017]: 3D-printed model classified as rifle from most viewpoints

Adversarial Examples (Security)

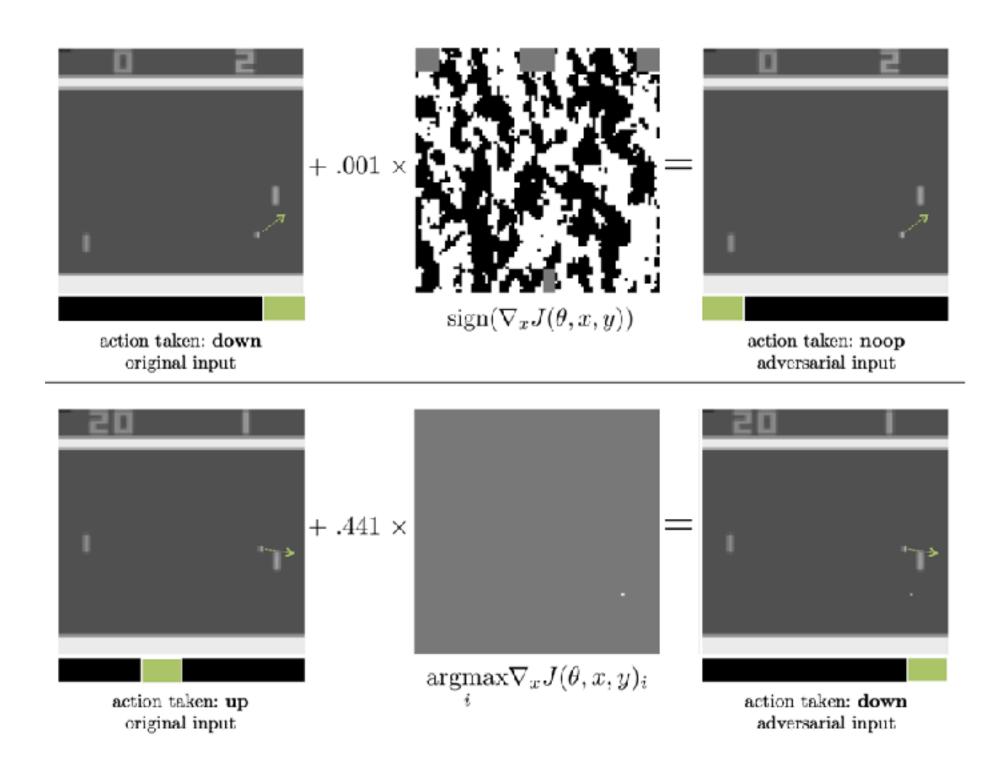


[Sharif et al. 2016]: Glasses the fool face classifiers



[Carlini et al. 2016]: Voice commands that are imperceptible by humans

Adversarial Examples (RL, NLP)



[Huang et al. 2017]: Small input changes can decrease RL performance

Article: Super Bowl 50

Paragraph: "Peyton Manning became the first quarter-back ever to lead two different teams to multiple Super Bowls. He is also the oldest quarterback ever to play in a Super Bowl at age 39. The past record was held by John Elway, who led the Broncos to victory in Super Bowl XXXIII at age 38 and is currently Denver's Executive Vice President of Football Operations and General Manager. Quarterback Jeff Dean had jersey number 37 in Champ Bowl XXXIV."

Question: "What is the name of the quarterback who was 38 in Super Bowl XXXIII?"

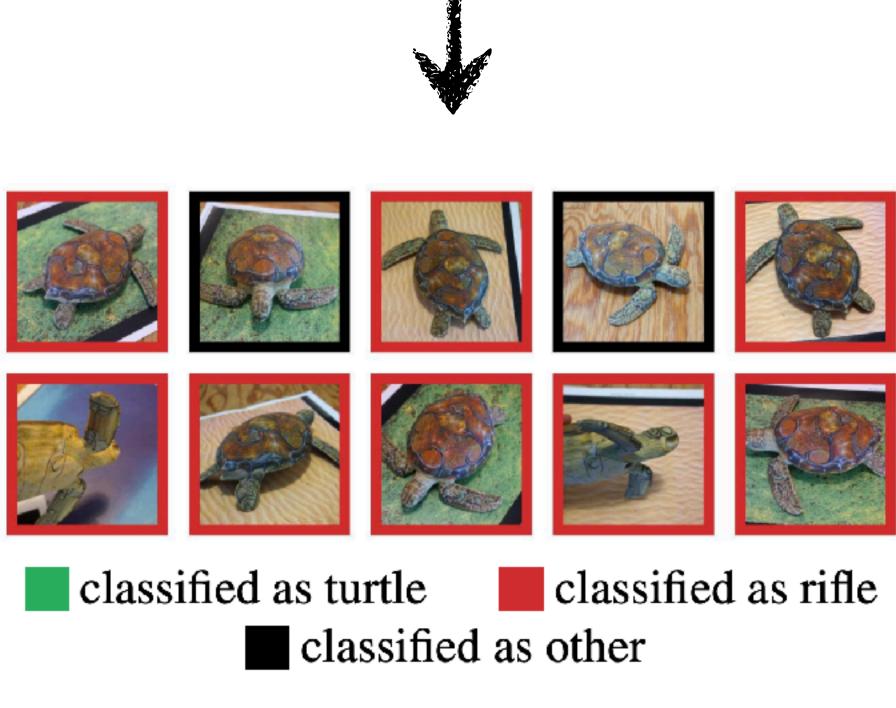
Original Prediction: John Elway

Prediction under adversary: Jeff Dean

[Jia Liang 2017]: Irrelevant sentences confused reading comprehension systems

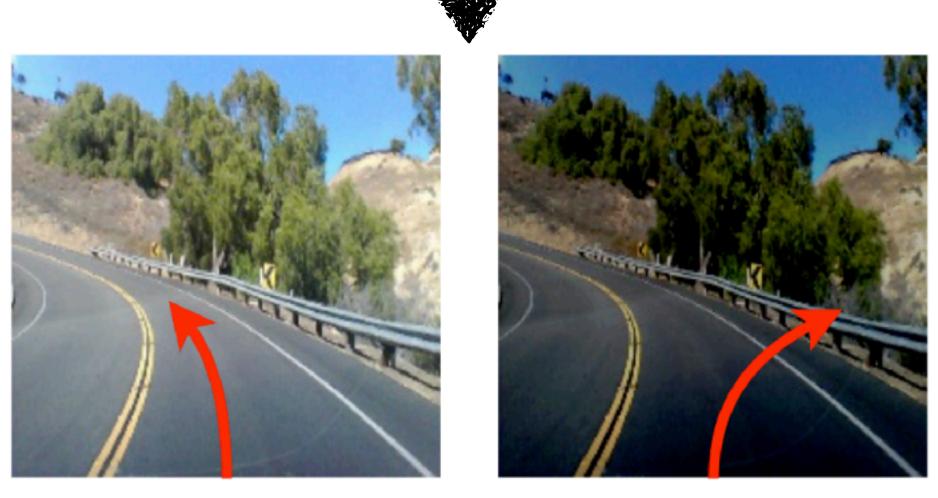
Should we be worried?

Probably not here!



[Athalye, Engstrom, Ilyas, Kwok 2017]: 3D-printed model classified as rifle from most viewpoints

But we should be worried here!



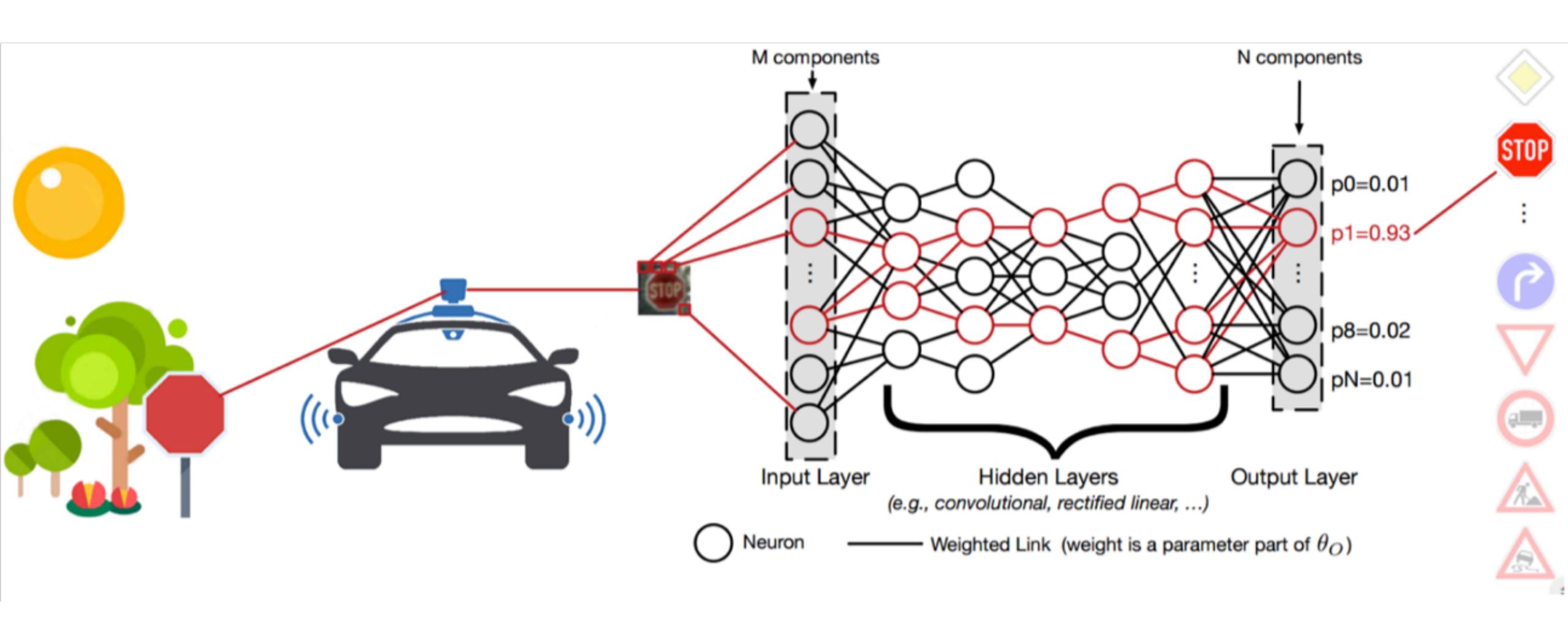
[Pei et al. 2017]: DeepXplore: Automated Whitebox Testing of Deep Learning Systems



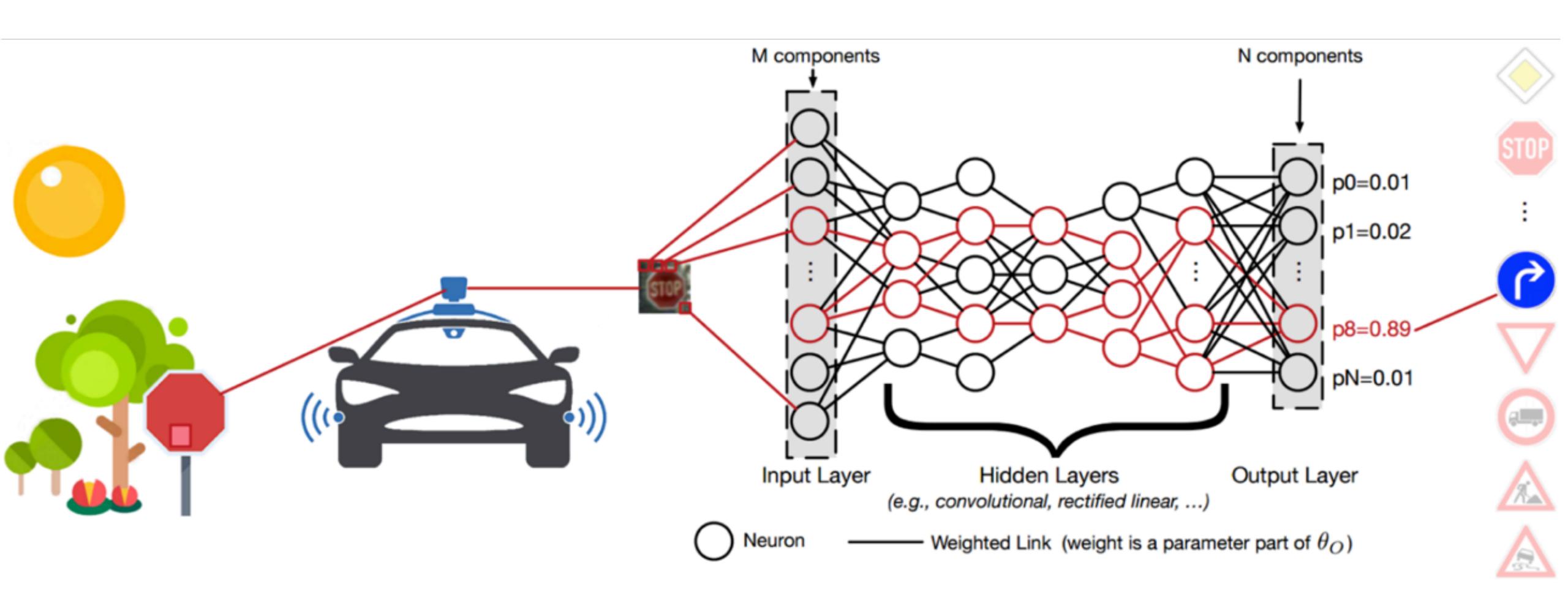


[Tian et al. 2017]: DeepTest: Automated Testing of Deep-Neural-Network-driven Autonomous Cars

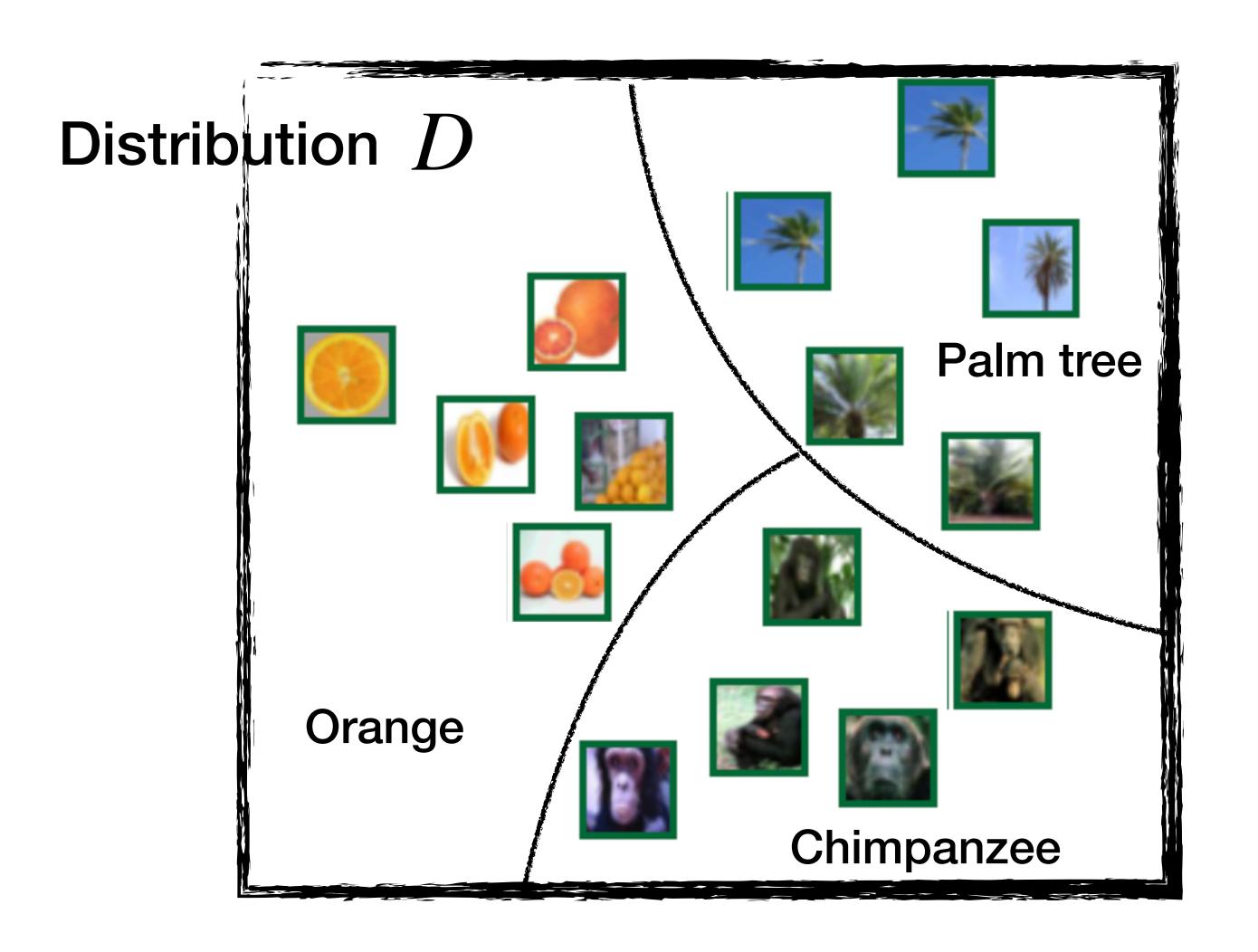
Should we be worried?

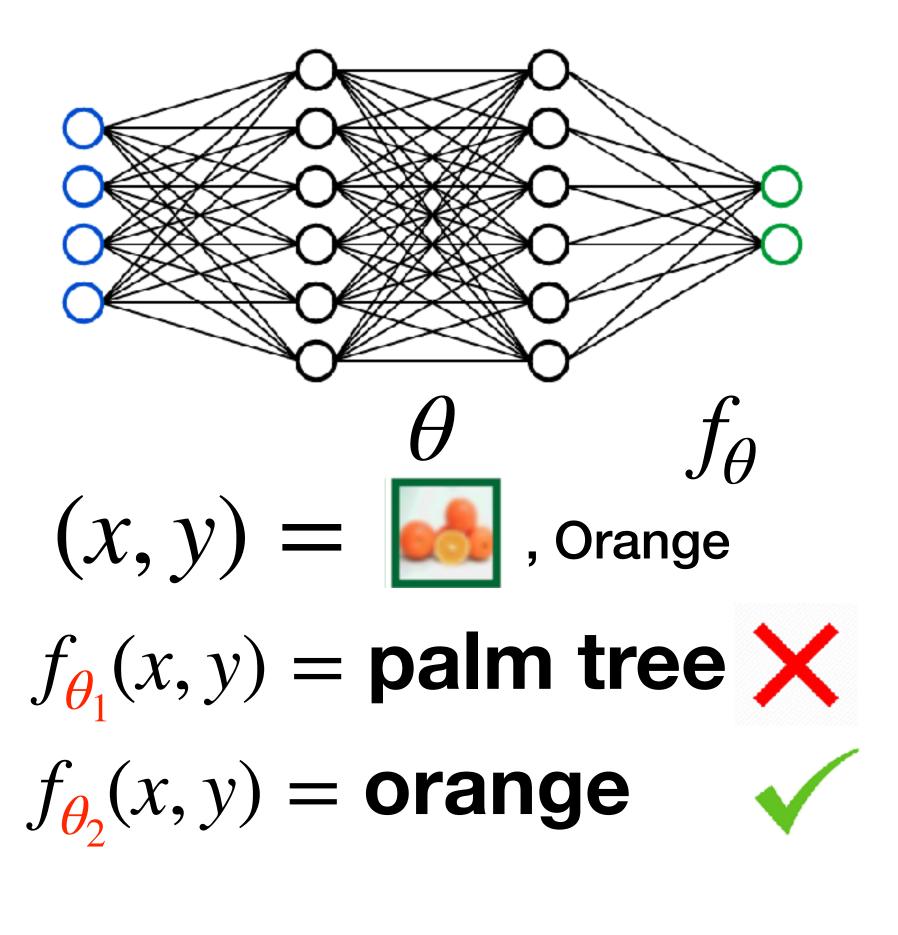


Should we be worried?



Where Do Adversarial Examples Come From?





Goal of ML:

Find θ^* such that $\mathbb{E}_{(x,y)\sim D}\mathcal{L}(\theta^*,x,y) \text{ Is small }$

Where Do Adversarial Examples Come From?

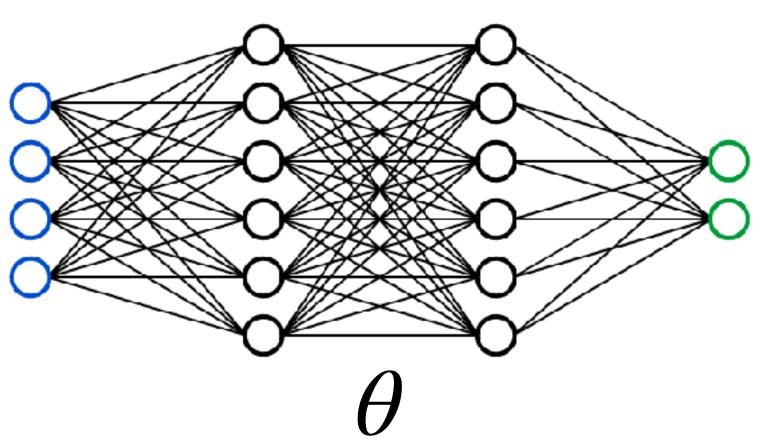
$$min_{\theta}\mathcal{L}(\theta, x, y)$$

$$max_{\delta}\mathcal{L}(\theta, x+\delta, y)$$

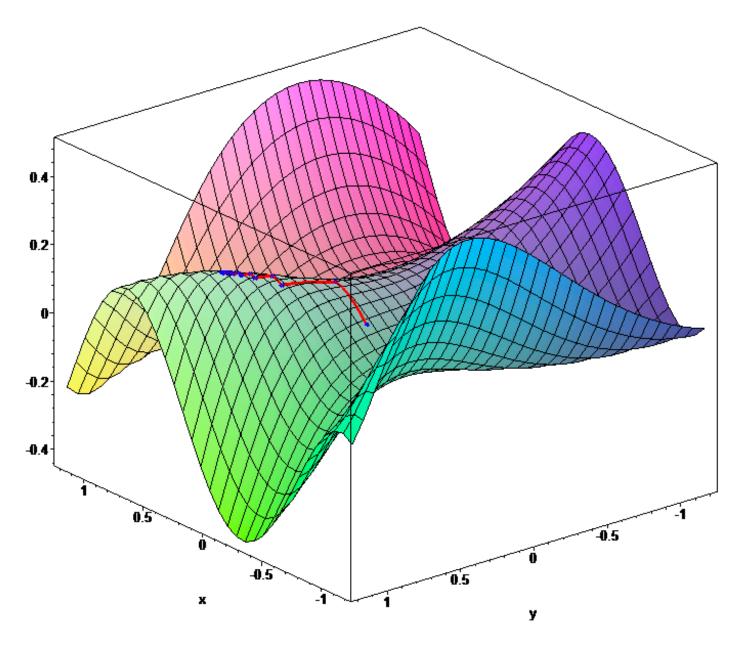
$$\|\delta\|_p \leq \epsilon$$

"Adversarial vulnerability is a direct result of our models' sensitivity to well-generalizing features in the data."

[Ilyas et al. 2019]: Adversarial Examples Are Not Bugs, They Are Features



Gradient Descent to find good parameters

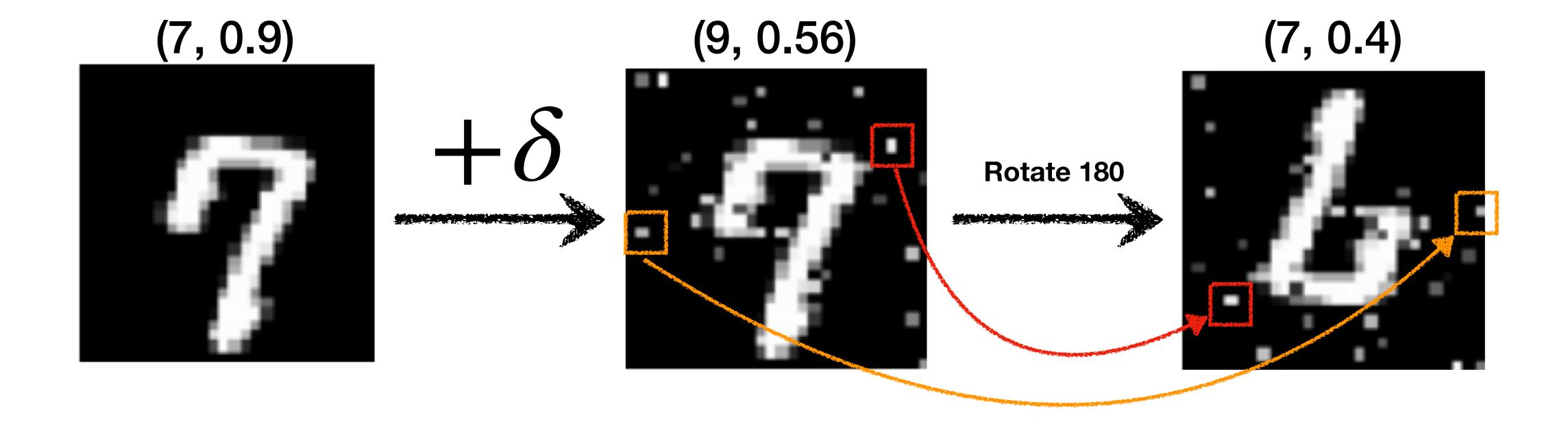


ATHENA:

A Framework for Defending Machine Learning Systems Against Adversarial Attacks



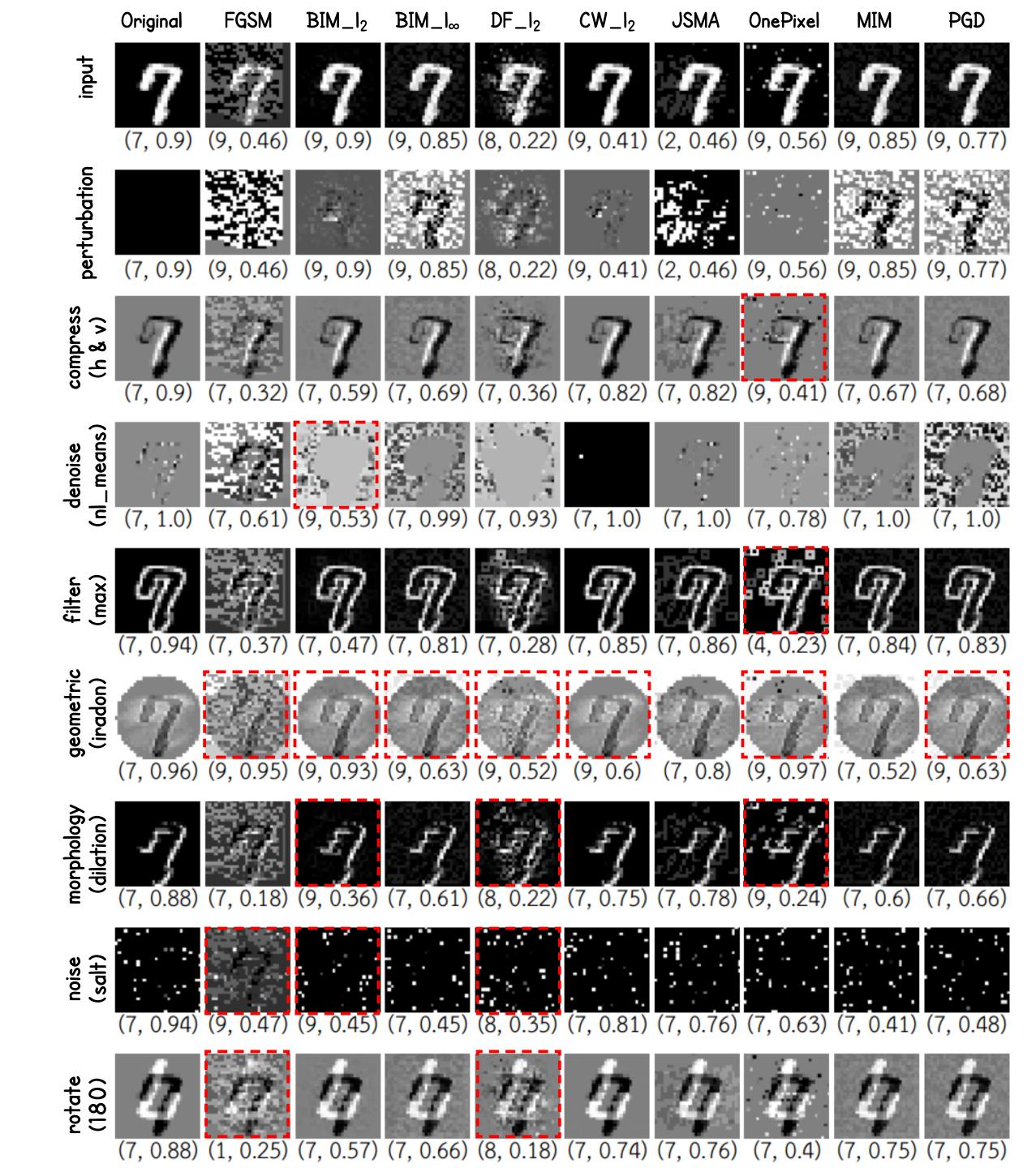
Key idea behind our approach: Input transformation



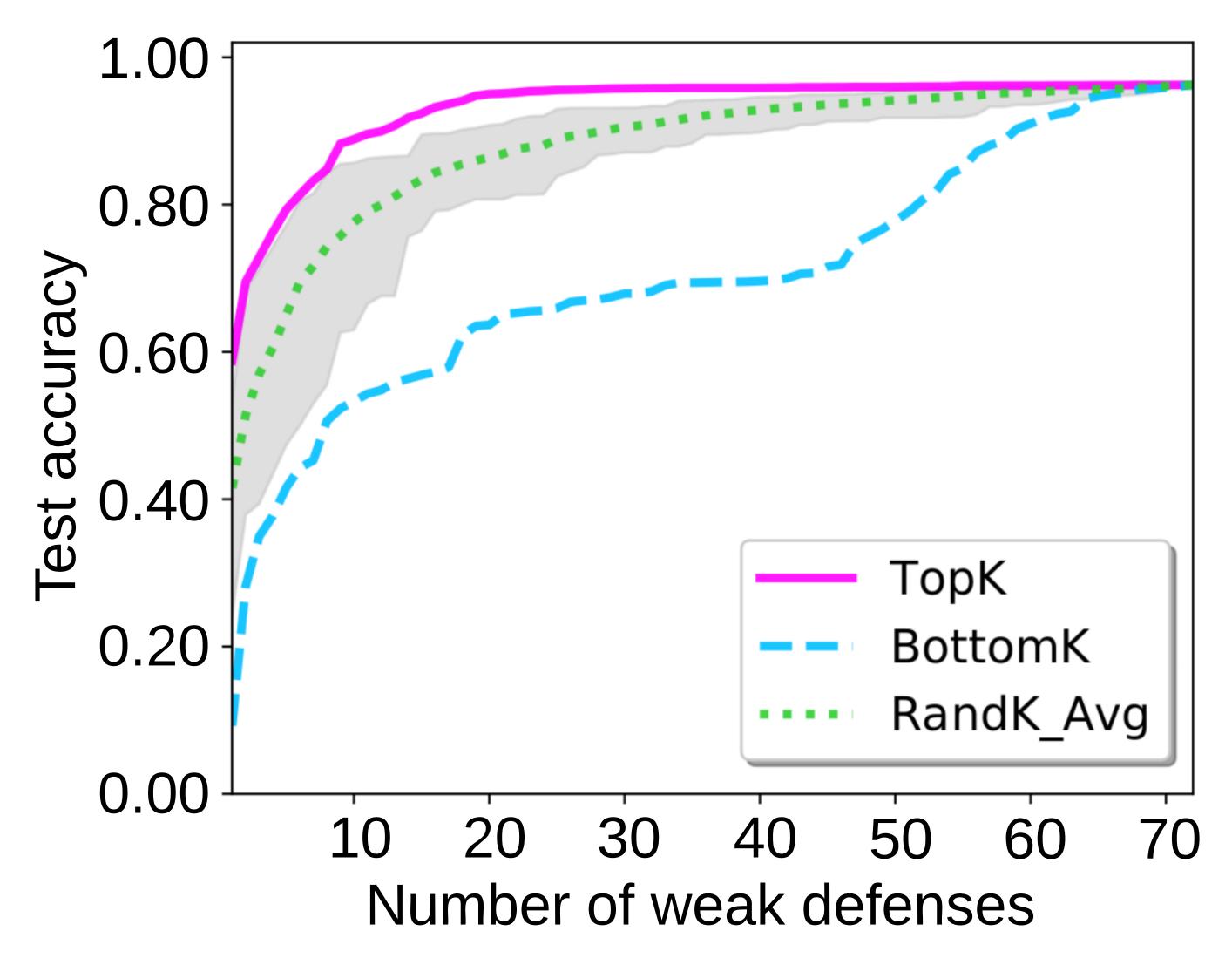
Insights

- Effectiveness of WDs varies
- WDs complement each other
- -> A defense based on *ensemble of WDs* can be independent of particular type of adversarial attack

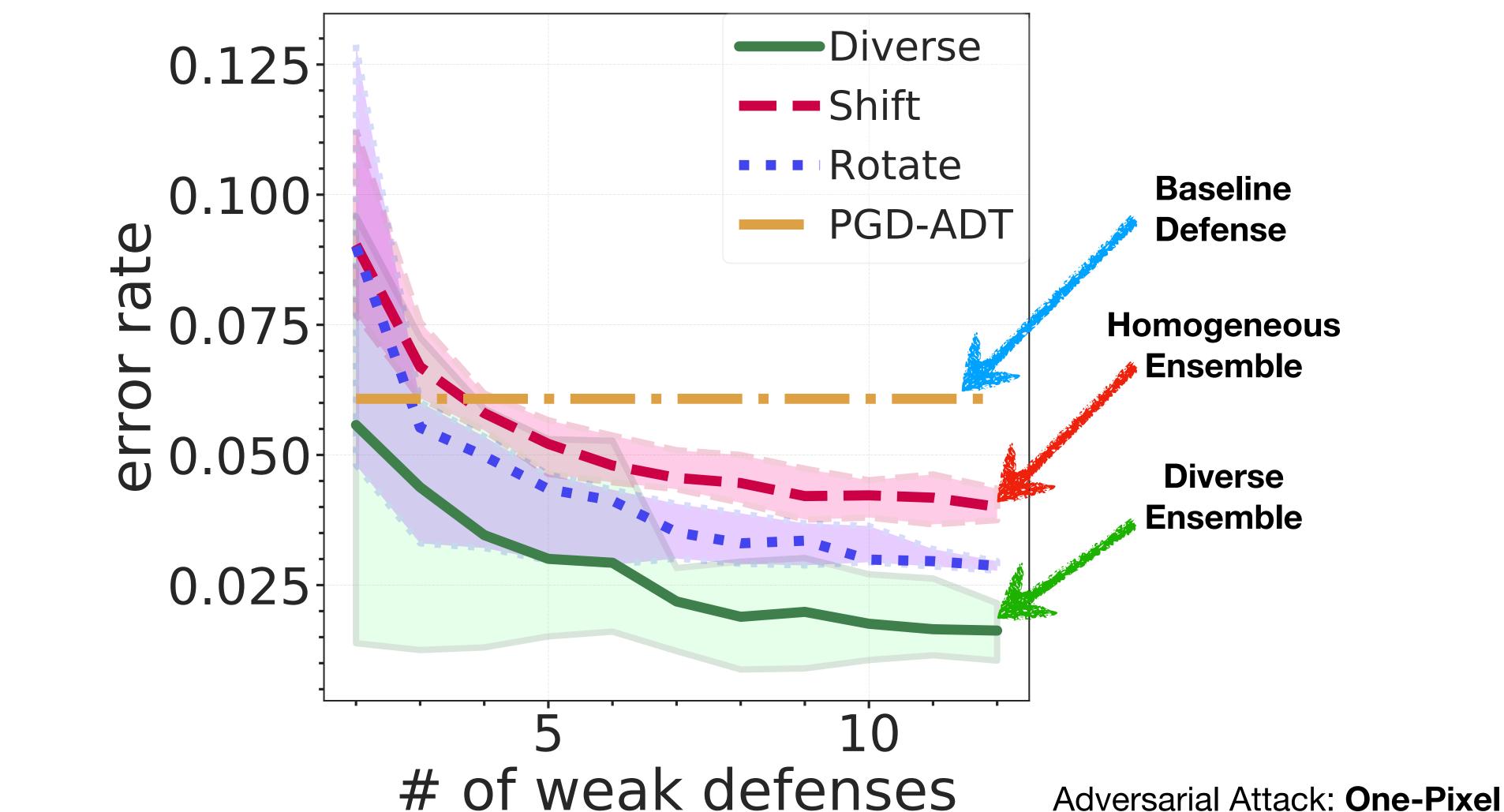
WD: Weak Defense



Quality and quantity of weak defenses matter

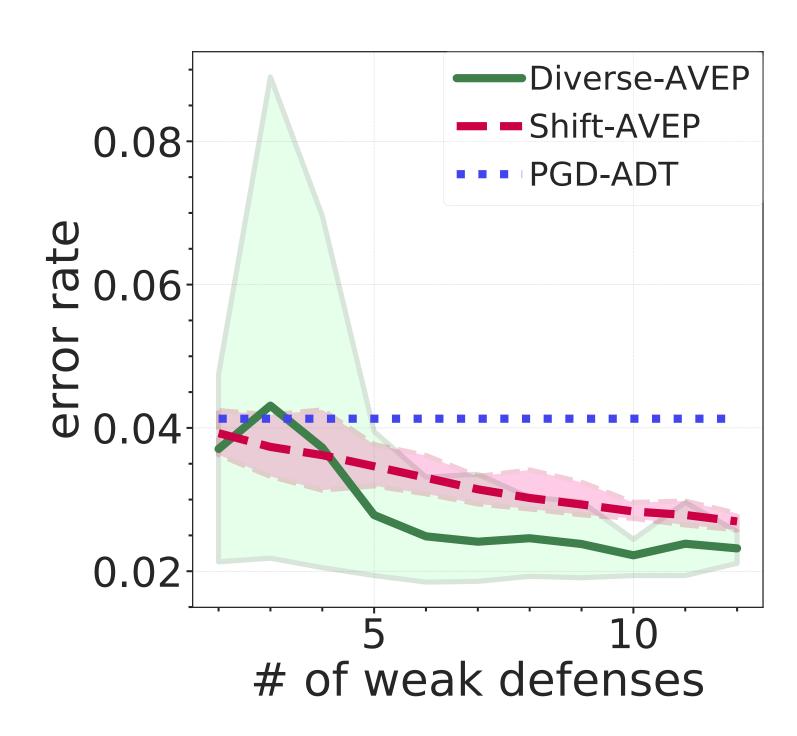


Diversity of weak defenses matters

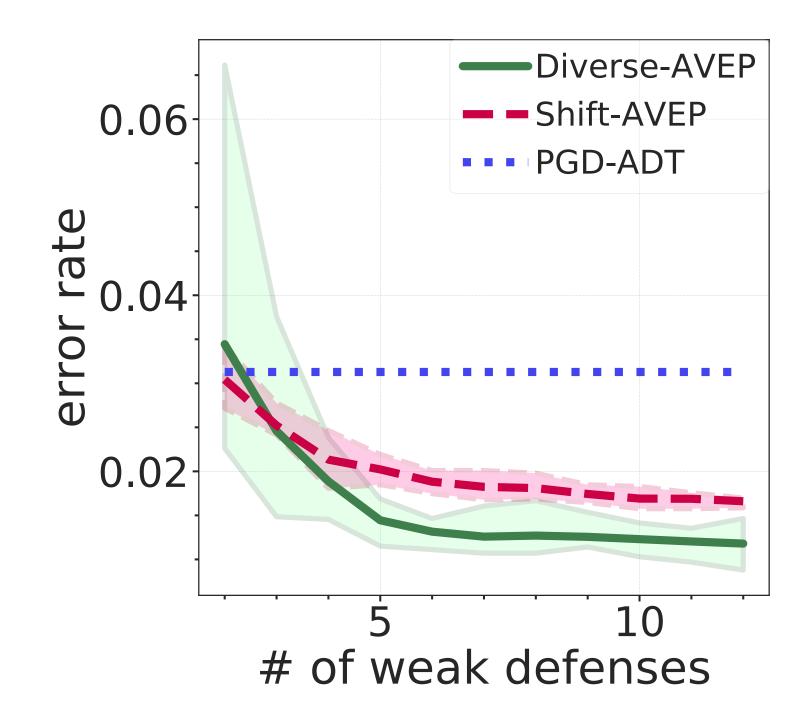


Error Rate Undefended: 0.5588
PGD-ADT: Adversarial Training

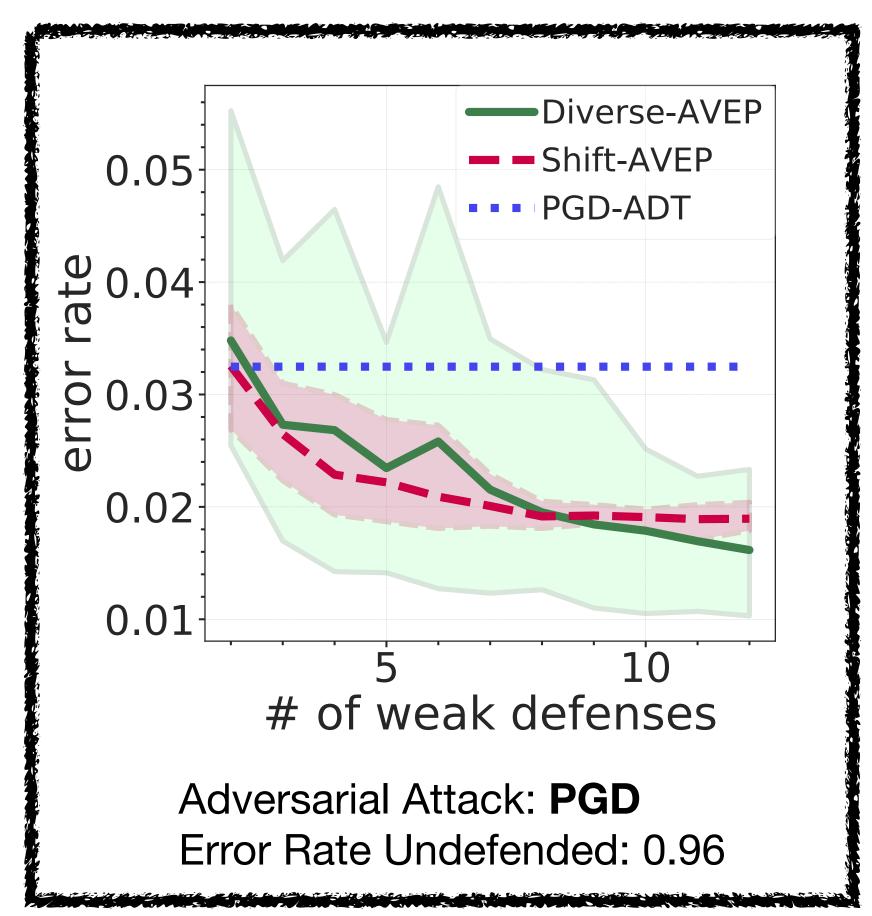
Diversity of weak defenses matters



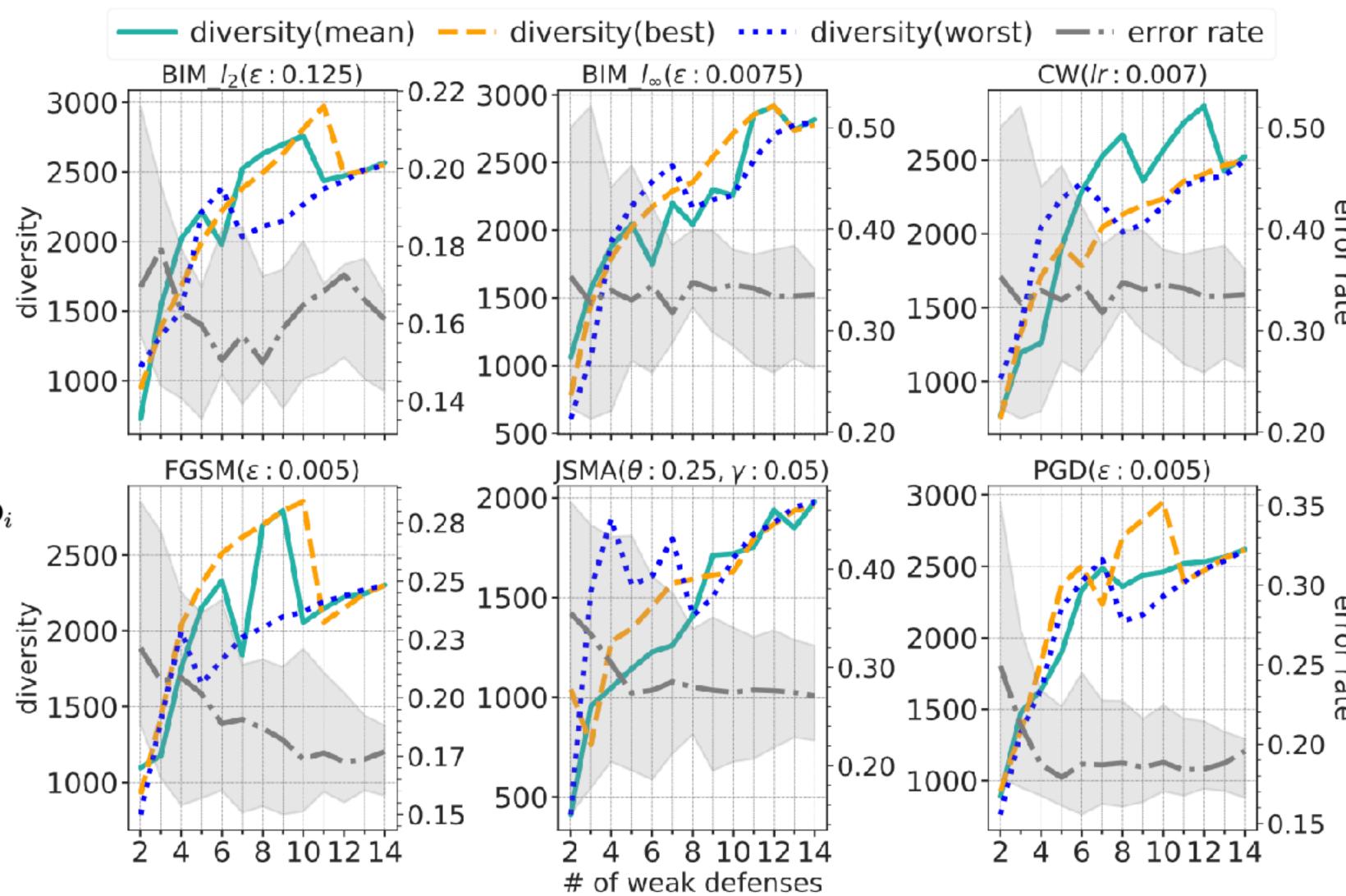
Adversarial Attack: **BIM_I2** Error Rate Undefended: 0.92



Adversarial Attack: **MIM**Error Rate Undefended: 0.94



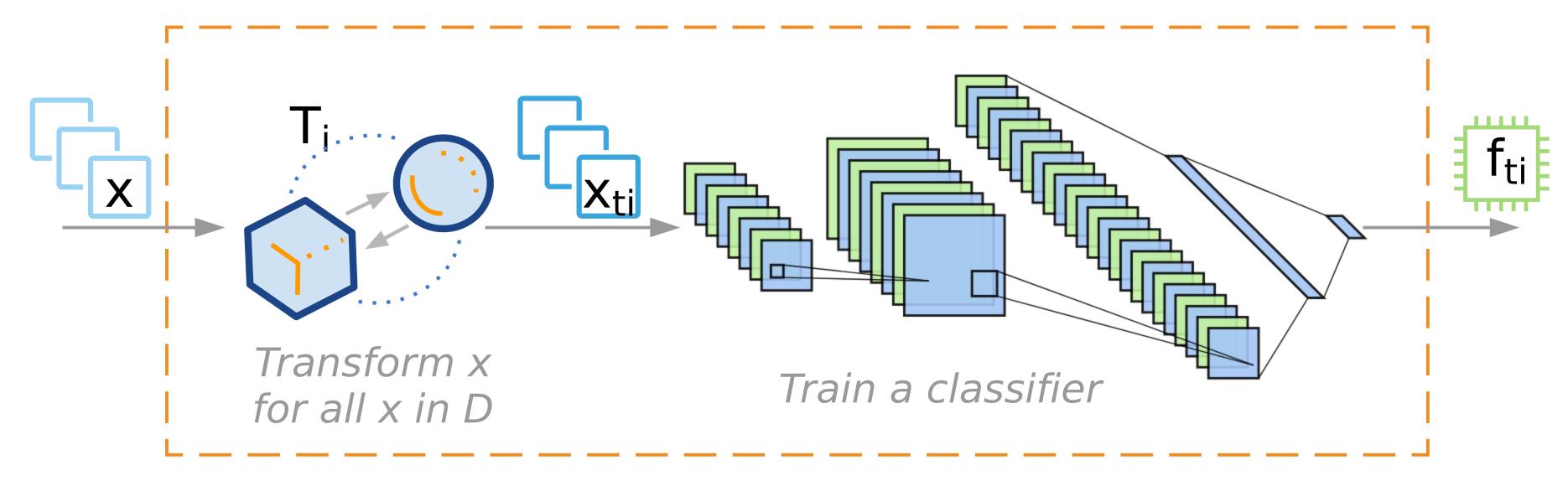
Diversity of weak defenses matters



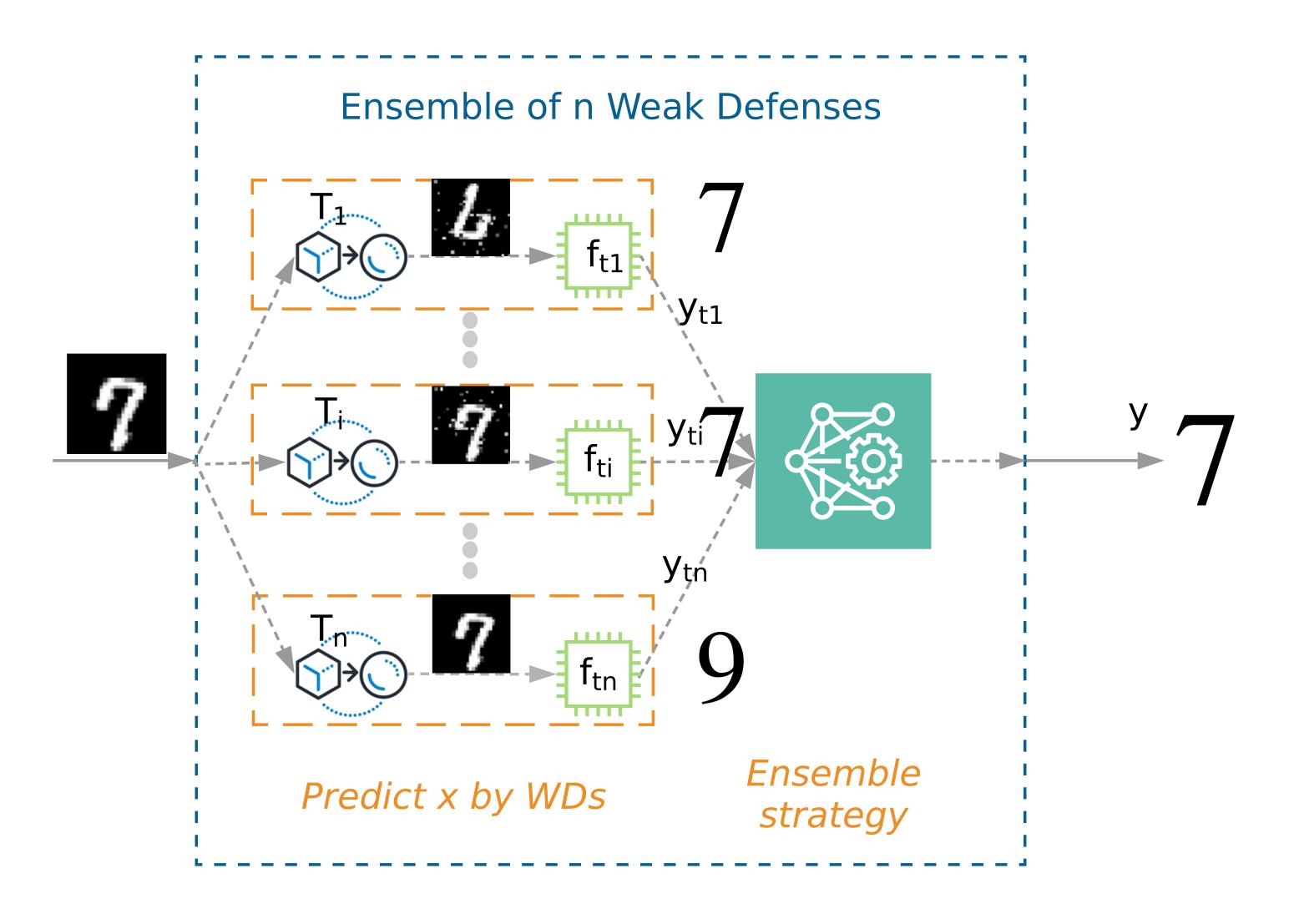
 $\psi = (\min_{i \in \{1, ..., K\}} |S_i|) - (|\bigcap_i S_i|)$

 S_i is the set of examples correctly predicted by WD_i

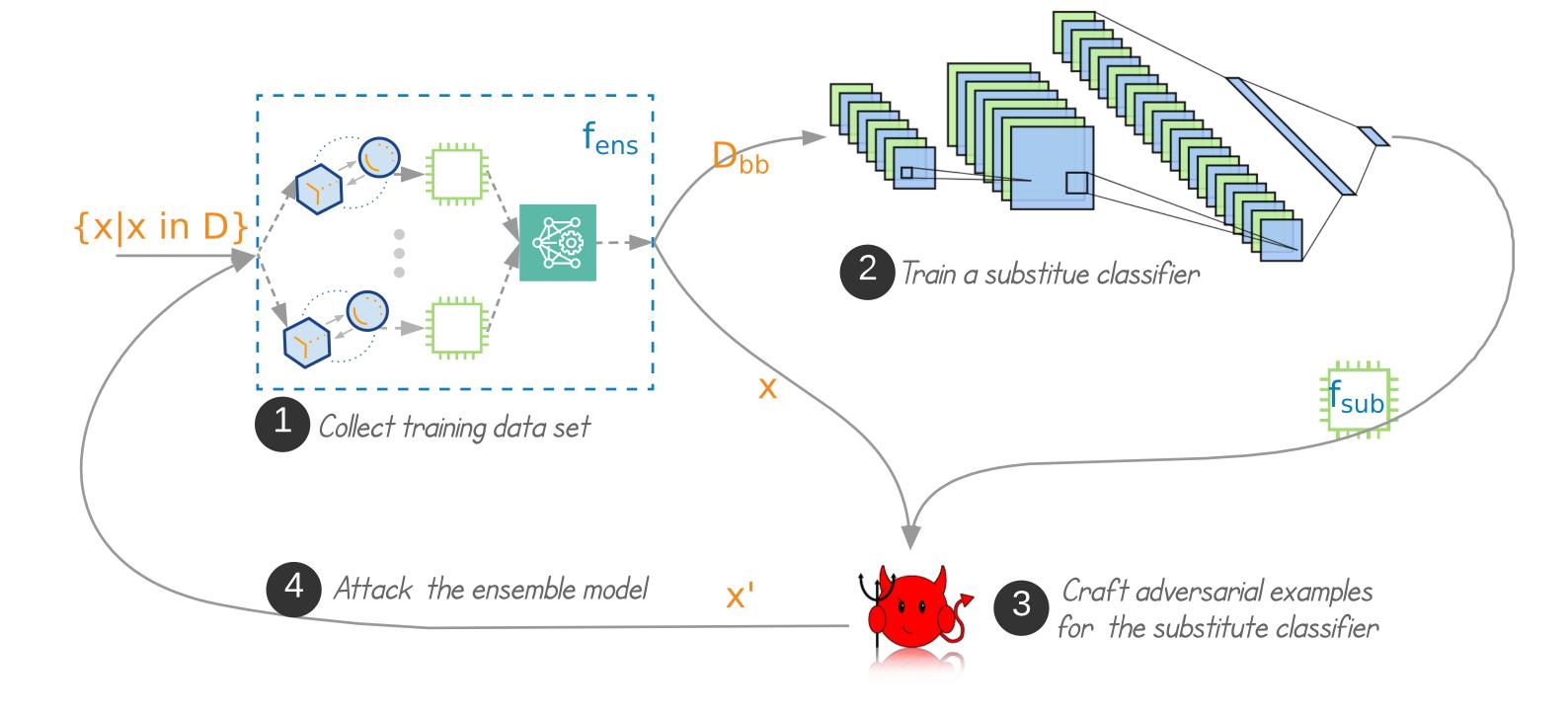
Each weak defense is essentially a model trained on a particular type of transformation



ATHENA produces the final output based on agreement between weak defenses at deployment time



Evaluation



Threat model: What we can assume about the knowledge of the adversary and its strength



Target

Knows the parameters of

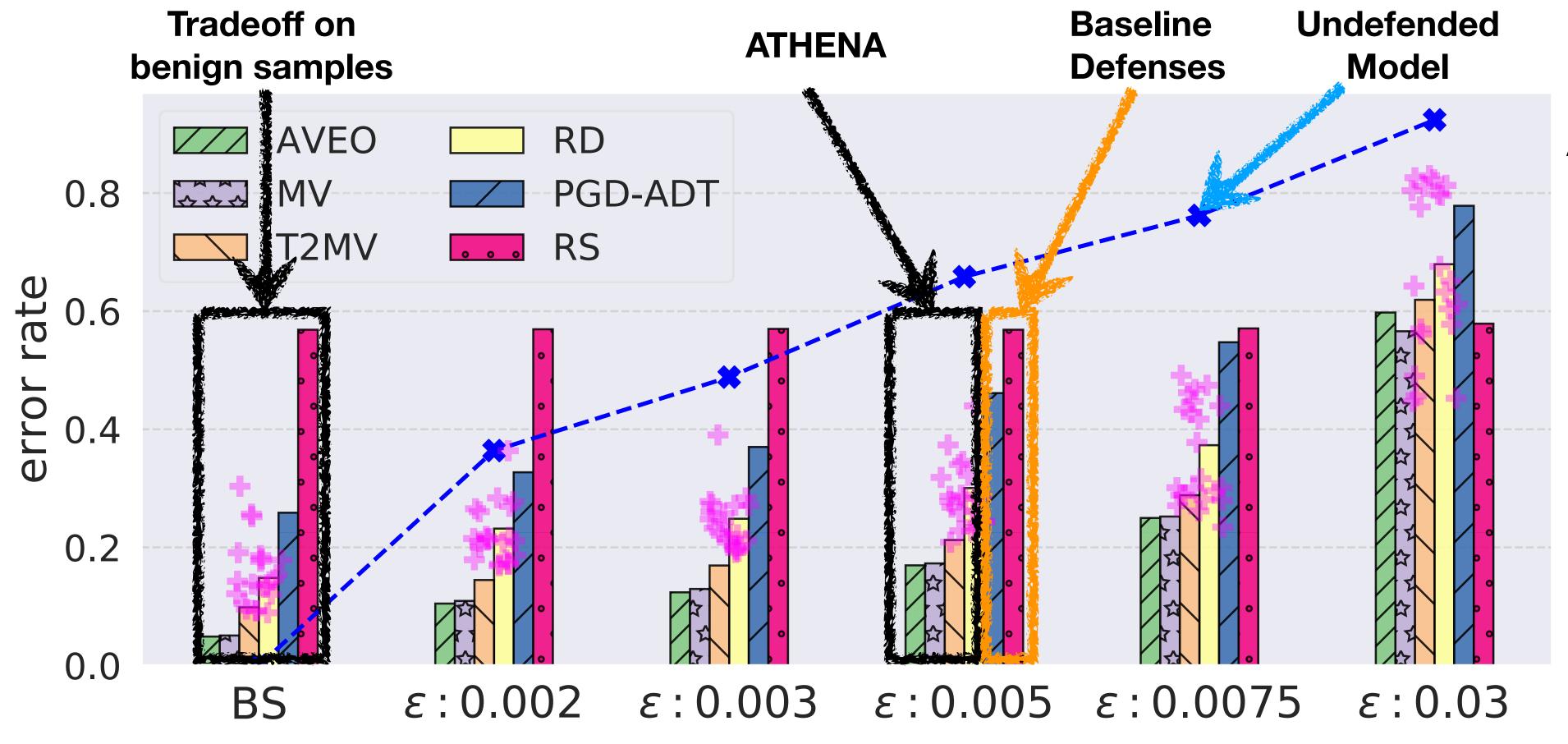
Weak

Ensemble

	Classifier	Defense	Defenses	Strategy
Zero-knowledge				
Blackbox				
Greybox				
Whitebox				

Existence of

Although the effectiveness of each weak defense varies, ATHENA is able to decrease the error rate effectively



ATHENA (ensemble strategy):

- MV: Majority Voting
- T2MV: Top-2 MV
- AVEO: Average of Output
- RD: Random Defense

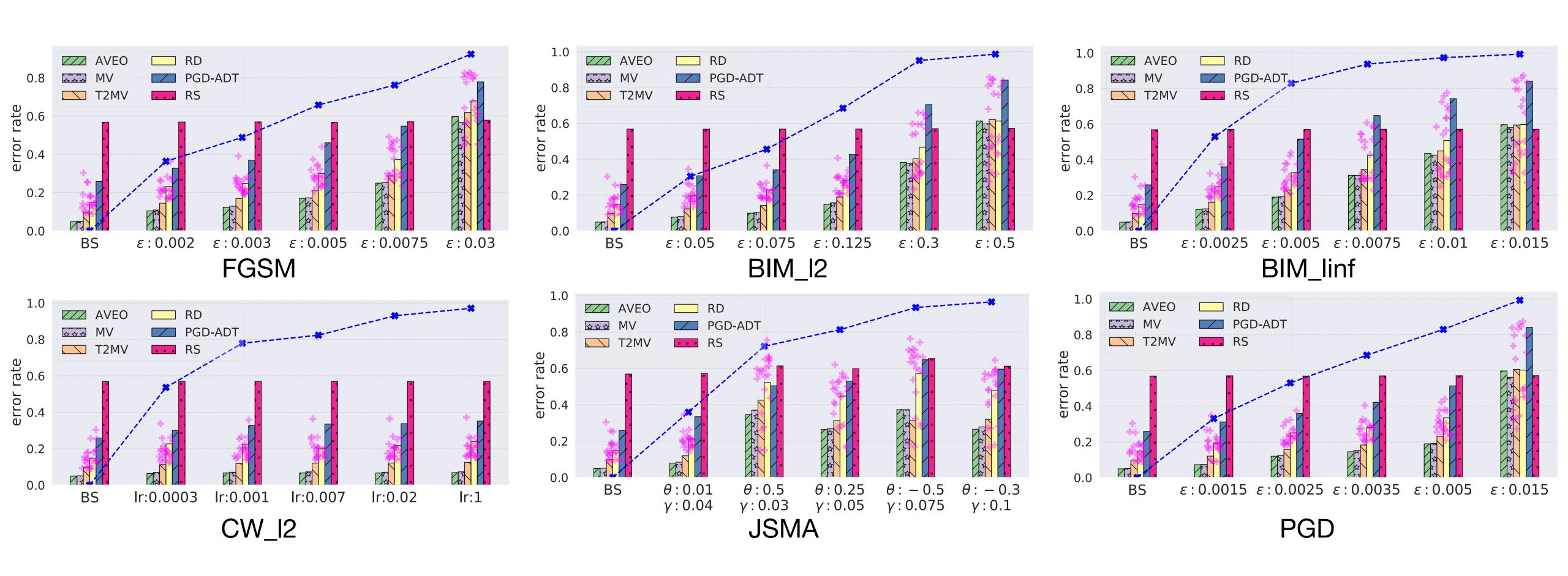
Baseline Defense:

PGD-ADT: Adversarial Training RS: Randomized Smoothing

Adversarial Attack: FGSM **Model**: 28×10 Wide ResNet

Dataset: CIFAR100

Although the effectiveness of each weak defense varies, ATHENA is able to decrease the error rate effectively



Model: 28×10 Wide ResNet

Dataset: CIFAR100

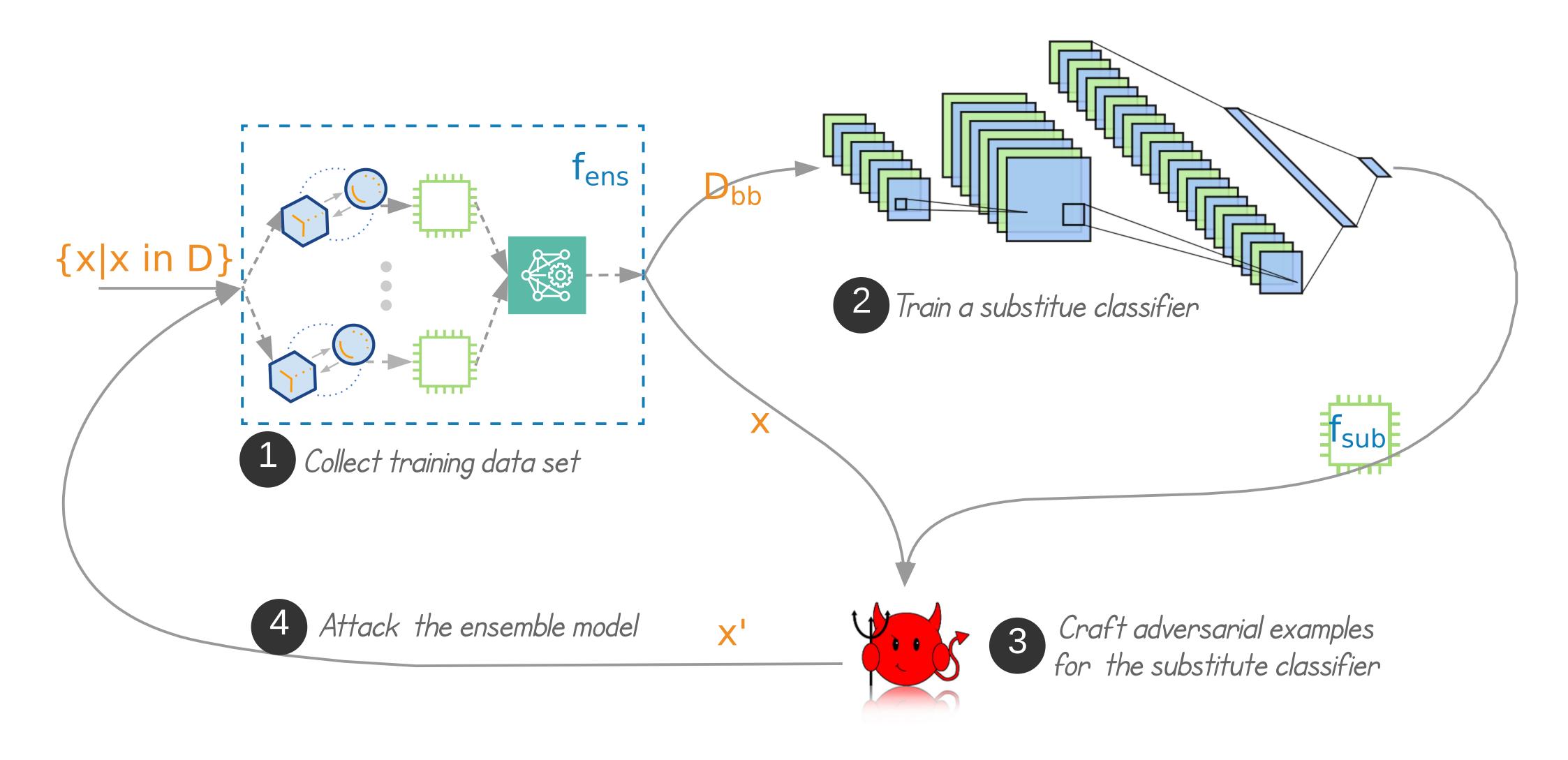
Threat model



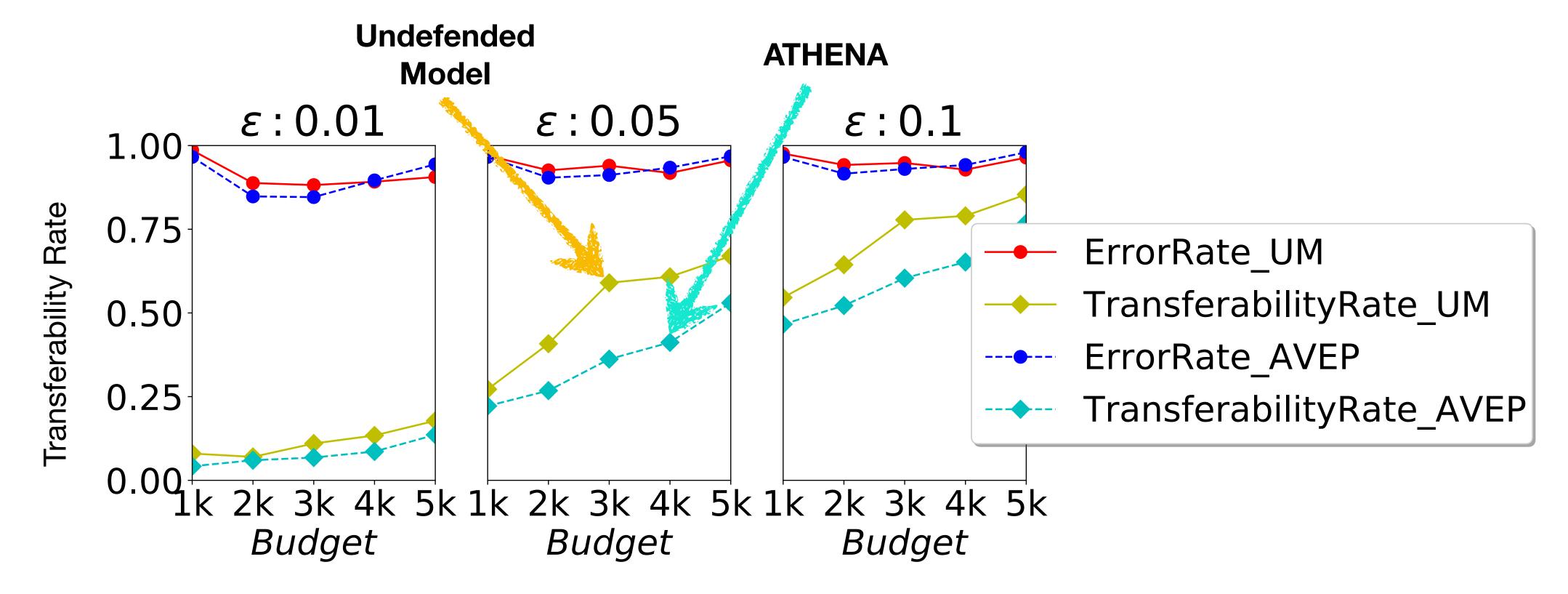
Knows the parameters of

Ensemble Existence of Weak Target Classifier Defense Defenses Strategy Zero-knowledge Blackbox Greybox Whitebox

Blackbox attack: The transferability-based approach



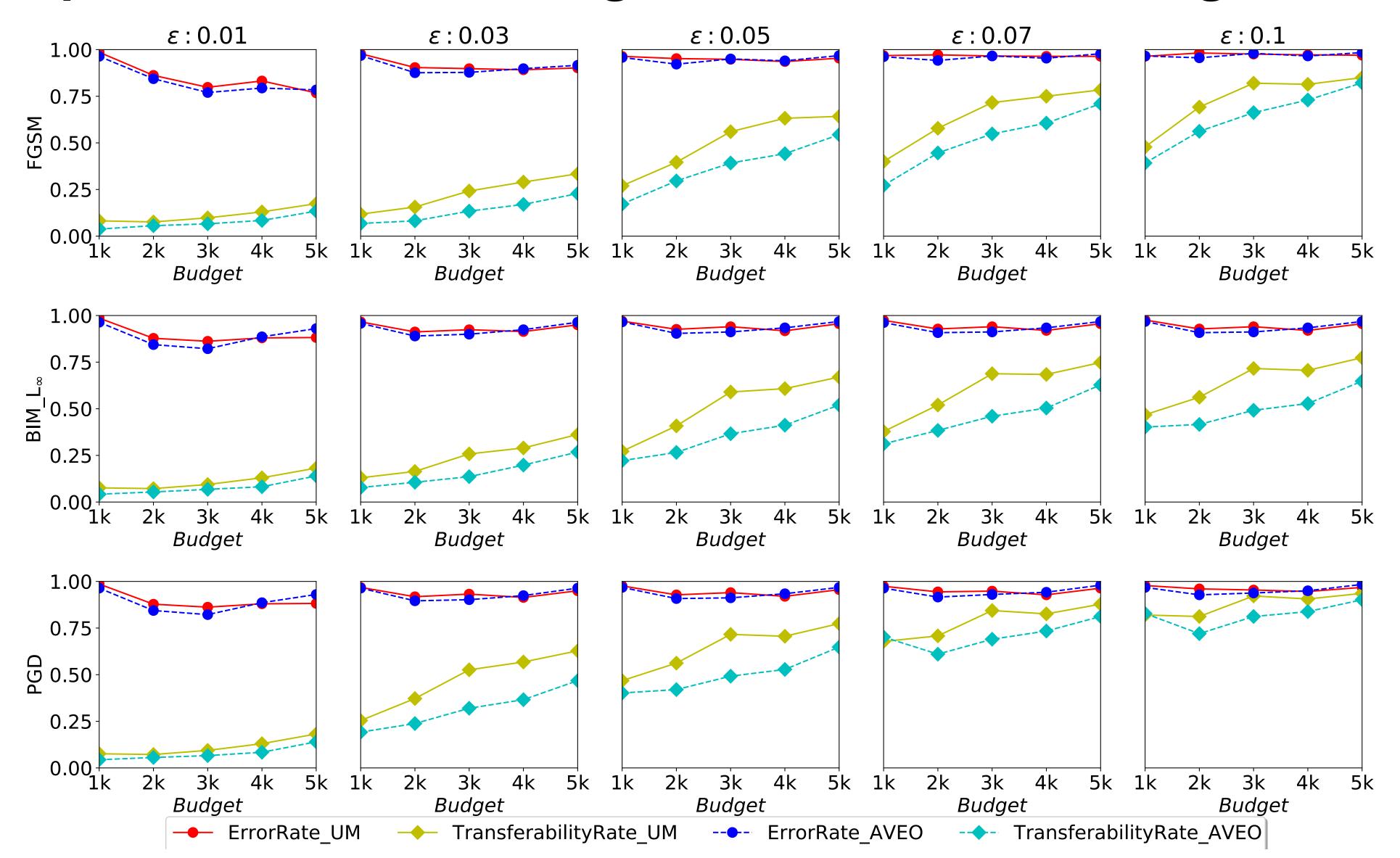
ATHENA lowered the "transferability" of adversarial examples from the surrogate model to the target model



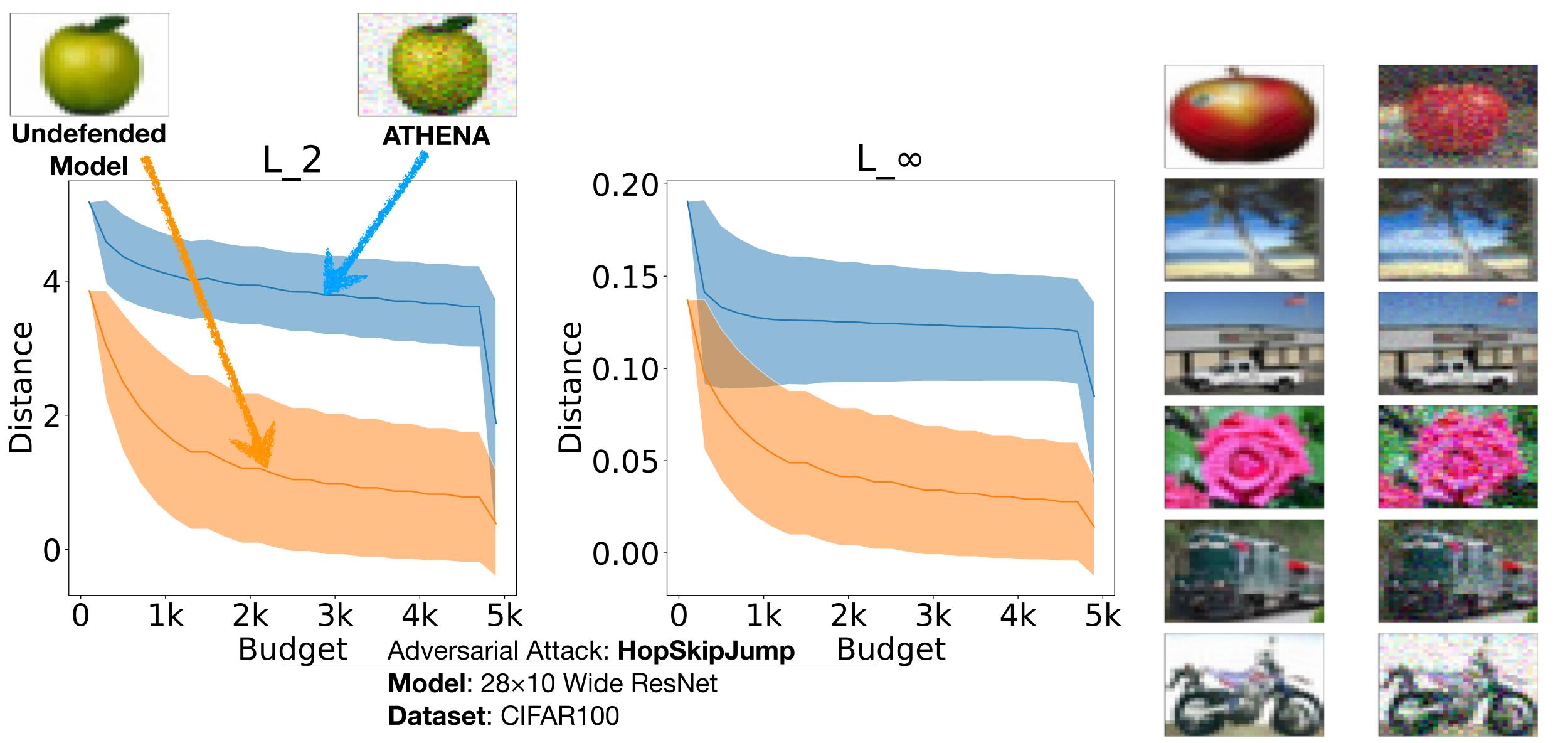
Adversarial Attack: BIM_linf **Model**: 28×10 Wide ResNet

Dataset: CIFAR100

ATHENA lowered the transferability of adversarial examples from the surrogate model to the target model



ATHENA forces the "optimization-based" blackbox attack to generate adversarial examples with larger perturbation



Threat model



Knows the parameters of

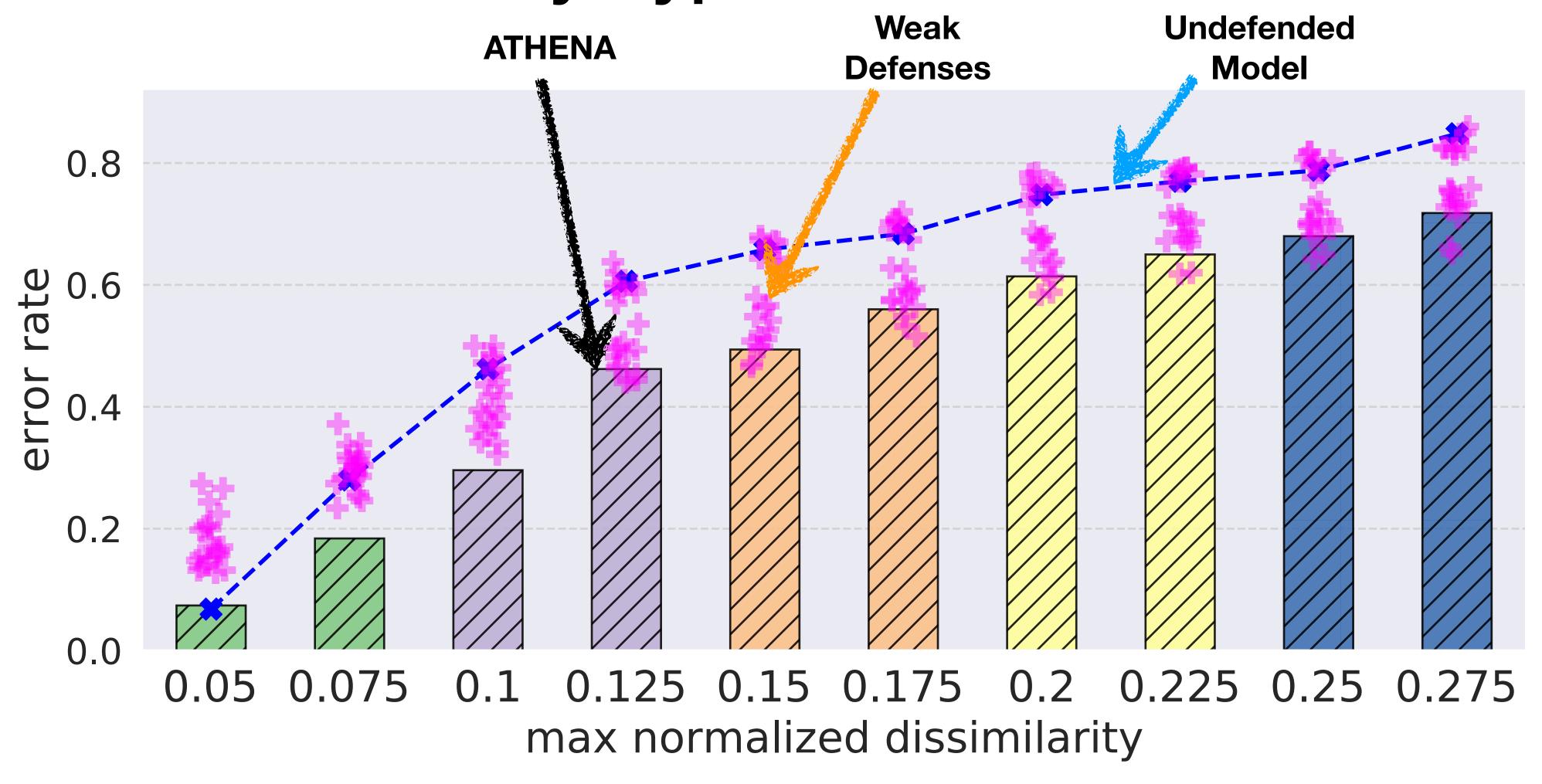
	Target Classifier	Existence of Defense	Weak Defenses	Ensemble Strategy
Zero-knowledge				
Blackbox				
Greybox				
Whitebox	THE RESERVENCE OF THE PERSON O	medical constants of the second constant of the second constants of the second constant of the second constants of the second constant of the second constants of the second constants of the second constants of the second constants of the second constant of the second co	Marken de la	Marches and and the second and the s

Greedy White-box Attack

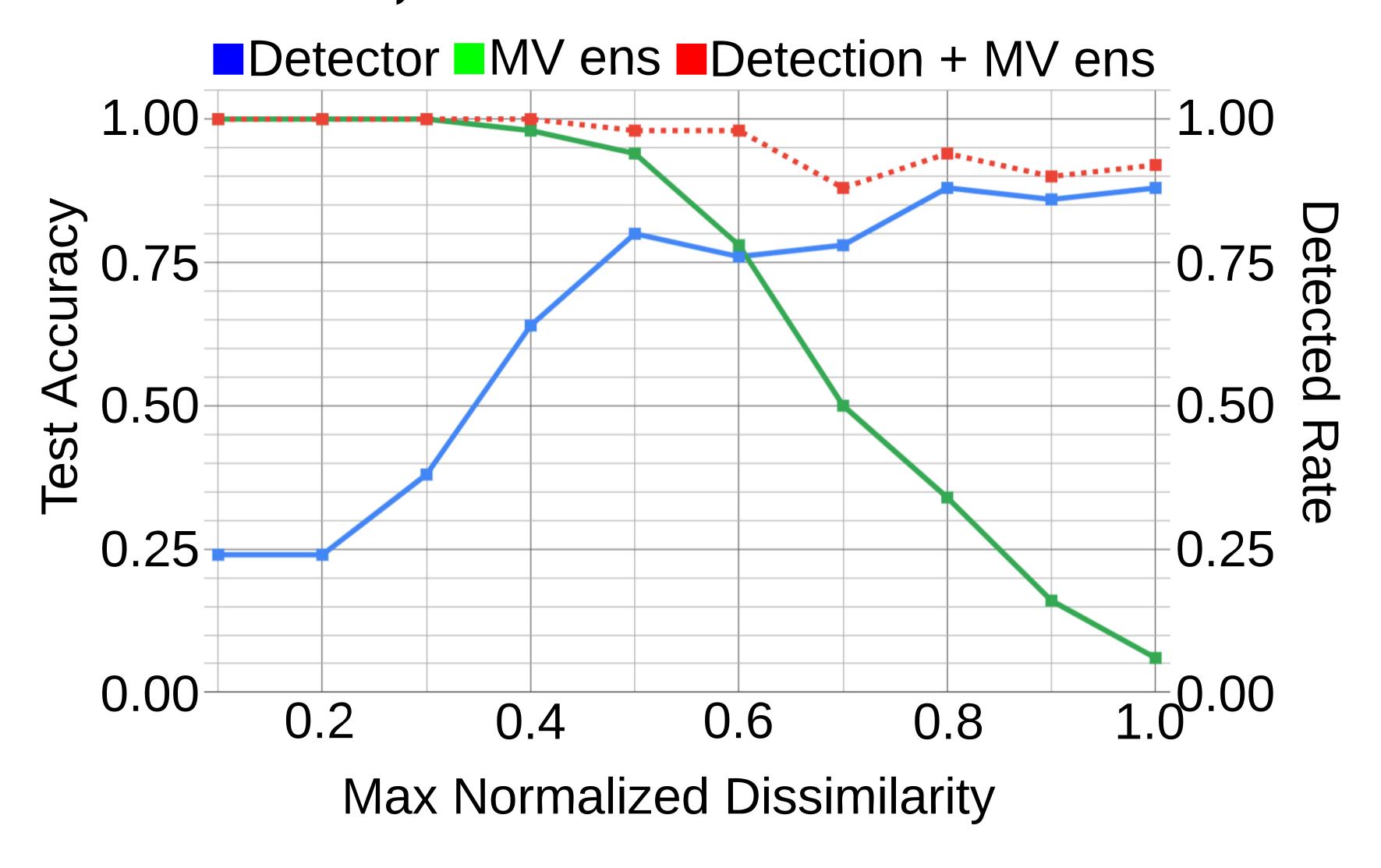
Algorithm 1: Crafting white-box AEs (Greedy)

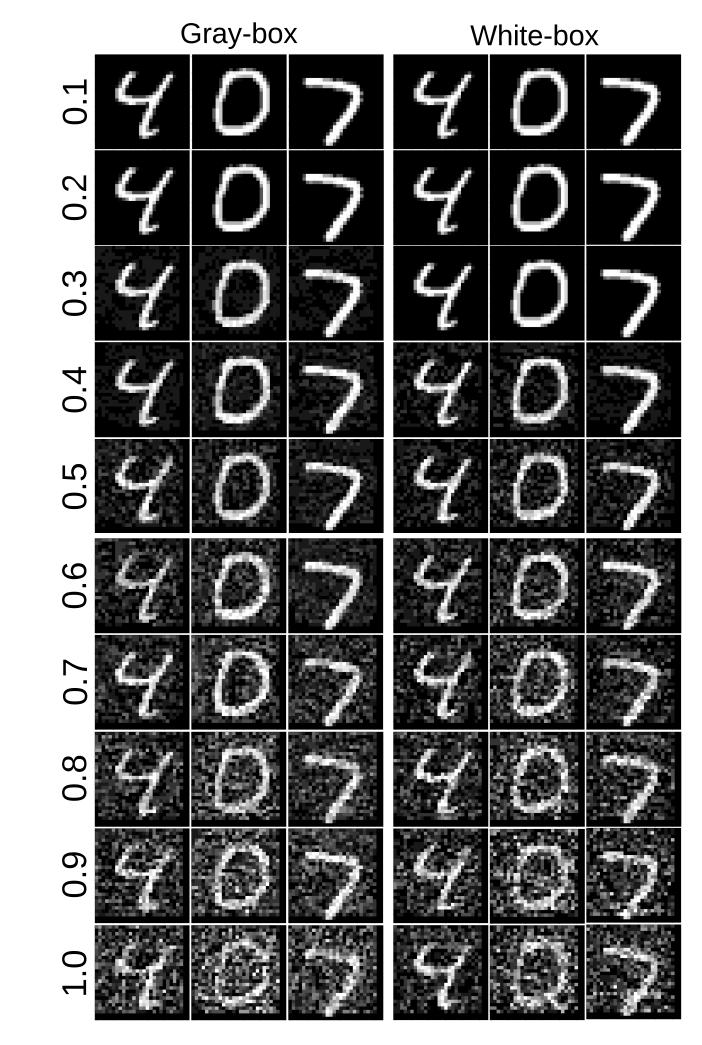
```
: x, y, attacker, N, max_dissimilarity
    input
 1 F_{fooled} \leftarrow \{\};
 2 F_{cand} \leftarrow all weak defenses;
 \mathbf{x}' \leftarrow \mathbf{x};
 4 while size(F_{fooled}) < N do
           f_{target} \leftarrow pickTarget(F_{cand}, strategy);
           // getPerturbation(\mathbf{x}) returns \|\mathbf{x} - \mathbf{x}'\|_2
           perturbation \leftarrow attacker.getPerturbation(f_{target}, \mathbf{x}');
           \mathbf{x}'_{tmp} \leftarrow \mathbf{x}' + \text{perturbation}
 8
           // dissimilairity(\mathbf{x}',\mathbf{x}) returns the normalized l_2
 9
                 dissimilarities between \mathbf{x}' and \mathbf{x}
           if dissimilarity(\mathbf{x}'_{tmp}, \mathbf{x}) > max\_dissimilarity then
10
                  break;
11
           end
12
           for f_{t_i} in F_{cand} do
13
                 if \mathbf{y} \neq f_{t_i}(\mathbf{x}'_{tmp}) then
14
                        addModel(F_{fooled}, f_{t_i});
15
                        removeModel(F_{cand}, f_{t_i});
16
                  end
17
           end
18
        \mathbf{x}' \leftarrow \mathbf{x}'_{tmp};
20 end
21 return x';
```

A strong adaptive white-box adversary may be able to successfully bypass the defense

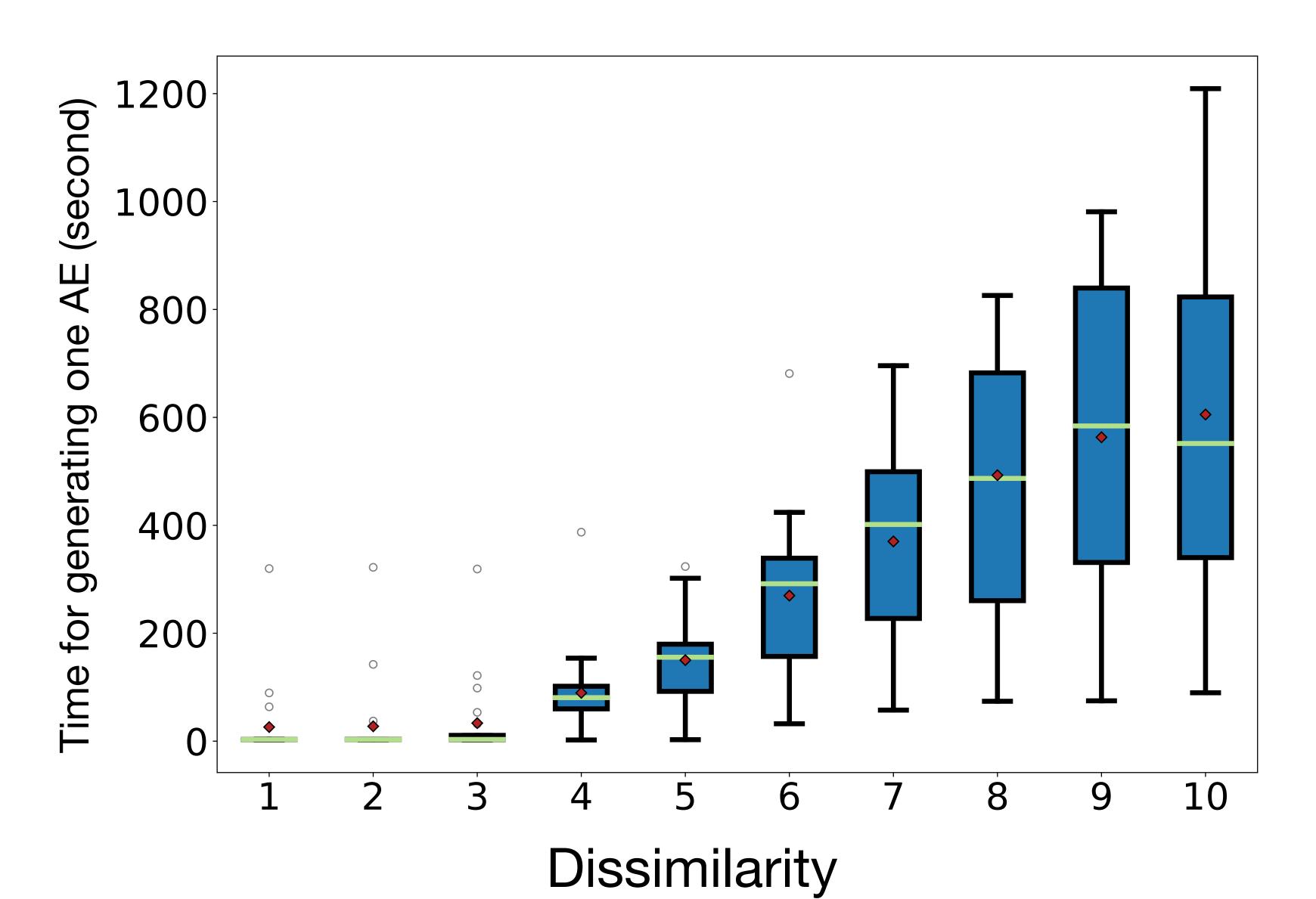


However, it becomes very easy to "detect" such attacks, so a defense+detection would be robust





Also, it comes with a high cost



An Adaptive Adversary for ATHENA

Standard Adversarial Attack

$$max_{\delta} \mathcal{L}(\theta, x + \delta, y)$$
$$||\delta||_{p} \le \epsilon$$

EOT (Expectation Over Transformation)

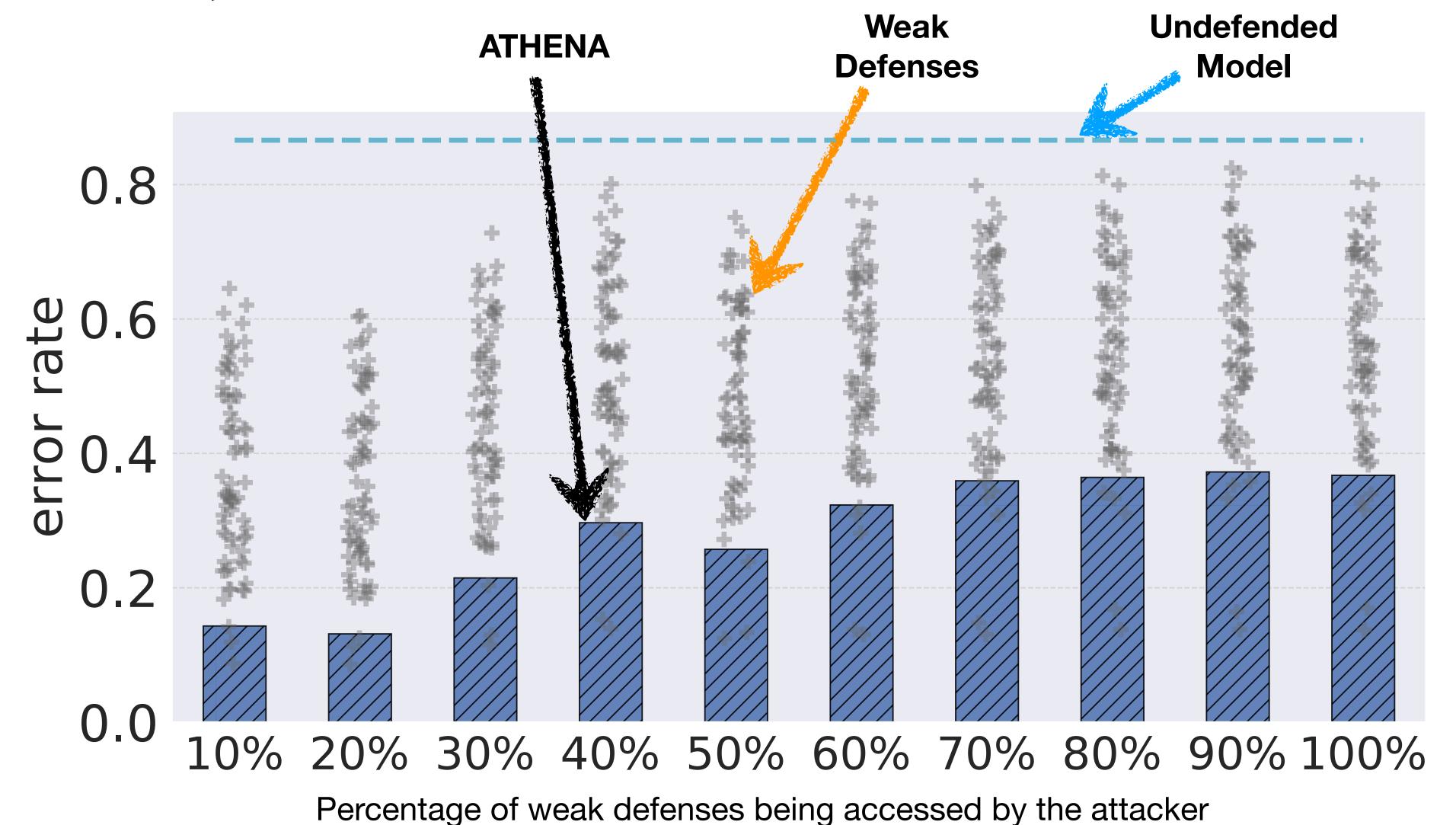
$$x' = \operatorname*{argmax}_{x'} \mathbb{E}_{t \sim T}[\log f(y_t | t(x'))]$$

s.t. $\mathbb{E}[d(t(x'), t(x))] < \epsilon, x \in [0, 1]^d$,

Extending EOT for Ensembles

$$\underset{x'}{\operatorname{argmax}} \frac{1}{K} \Sigma_{i=1}^K \mathbb{E}_{t \sim T} [\log f_i(y_t | t(x')]$$
 s.t. $\mathbb{E}_{t \sim T} [d(t(x'), t(x))] < \epsilon, x \in [0, 1]^d,$

As the adaptive attacker knows more about ATHENA, it can launch more successful attacks

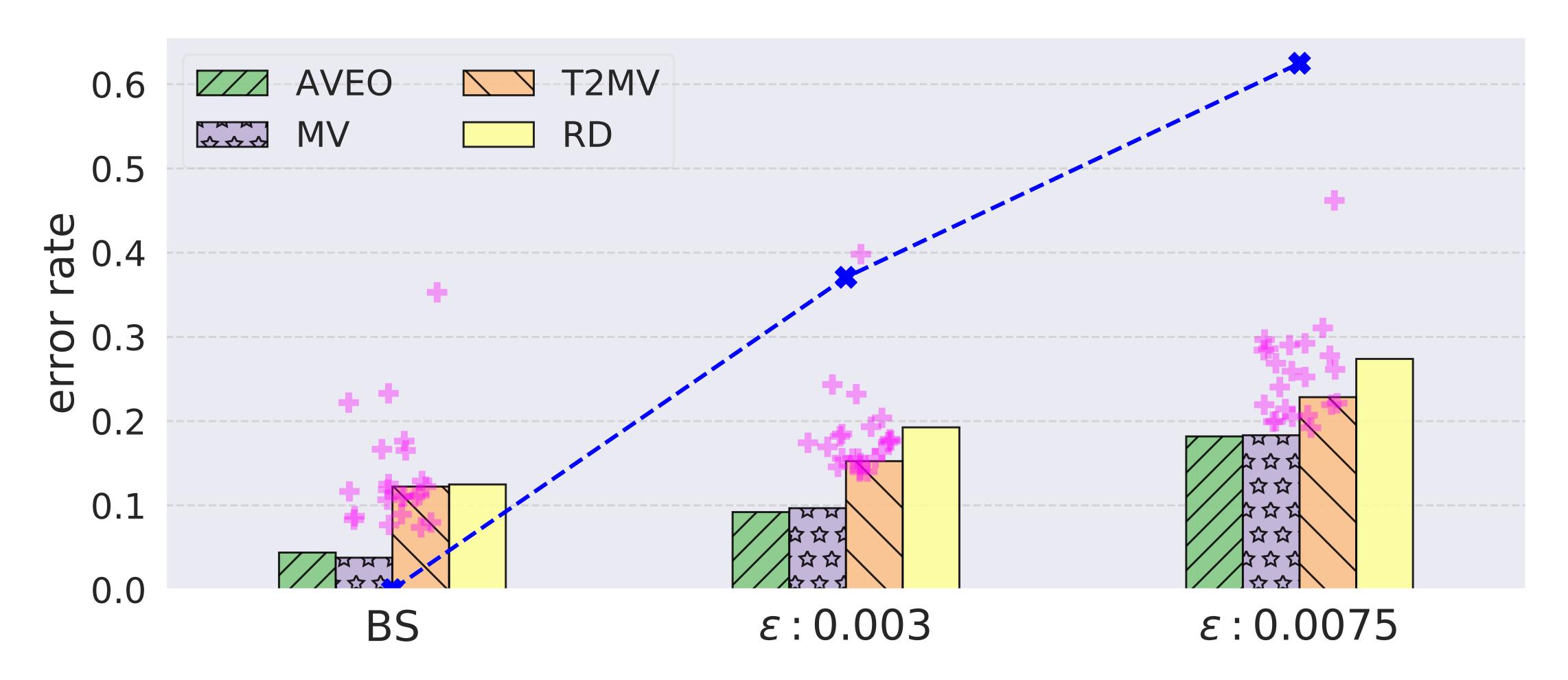


Is ATHENA a general defense?

Will it work with different types of machine learning models?



ATHENA performs similarly well with other types of machine learning models (DNNs, SVMs, RF)



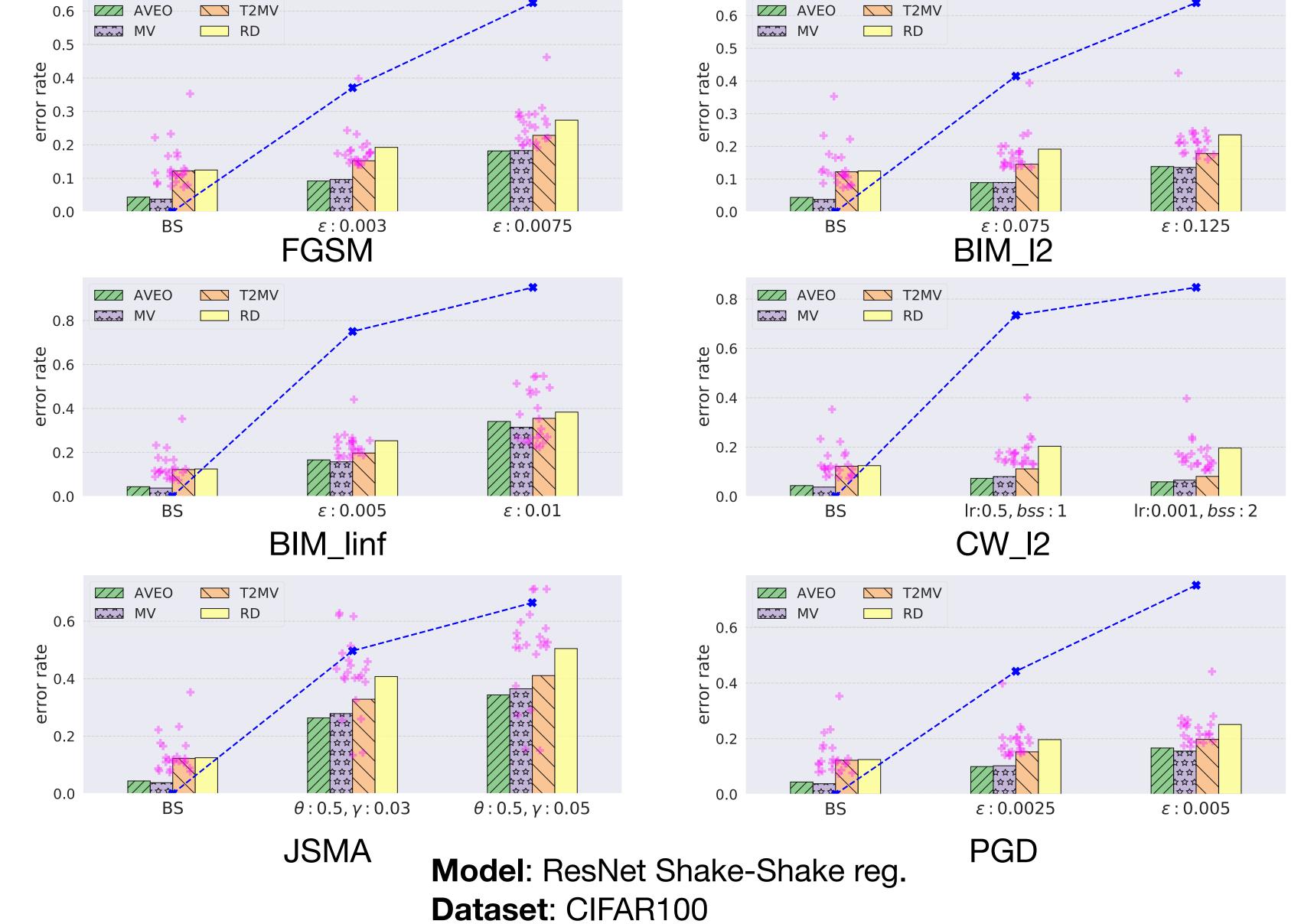
Adversarial Attack: FGSM

Model: ResNet Shake-Shake reg

Dataset: CIFAR100

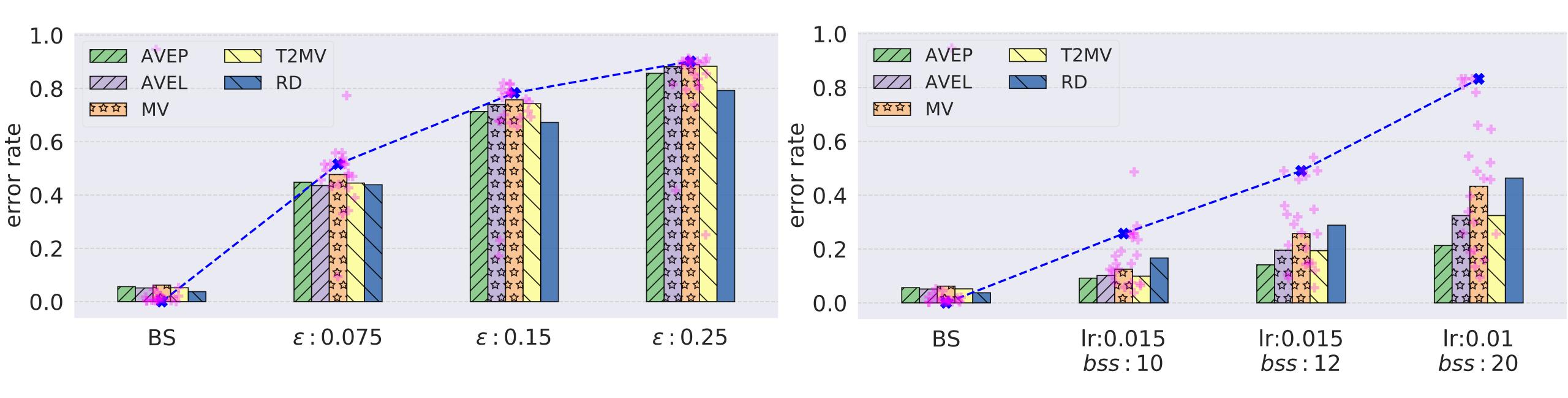
ATHENA is effective similarly with other types of

models



57

However, the effectiveness of defense may vary depending on the type of models



Adversarial Attack: FGSM

Model: SVM

Dataset: MNIST

Adversarial Attack: CW_I2

Model: SVM

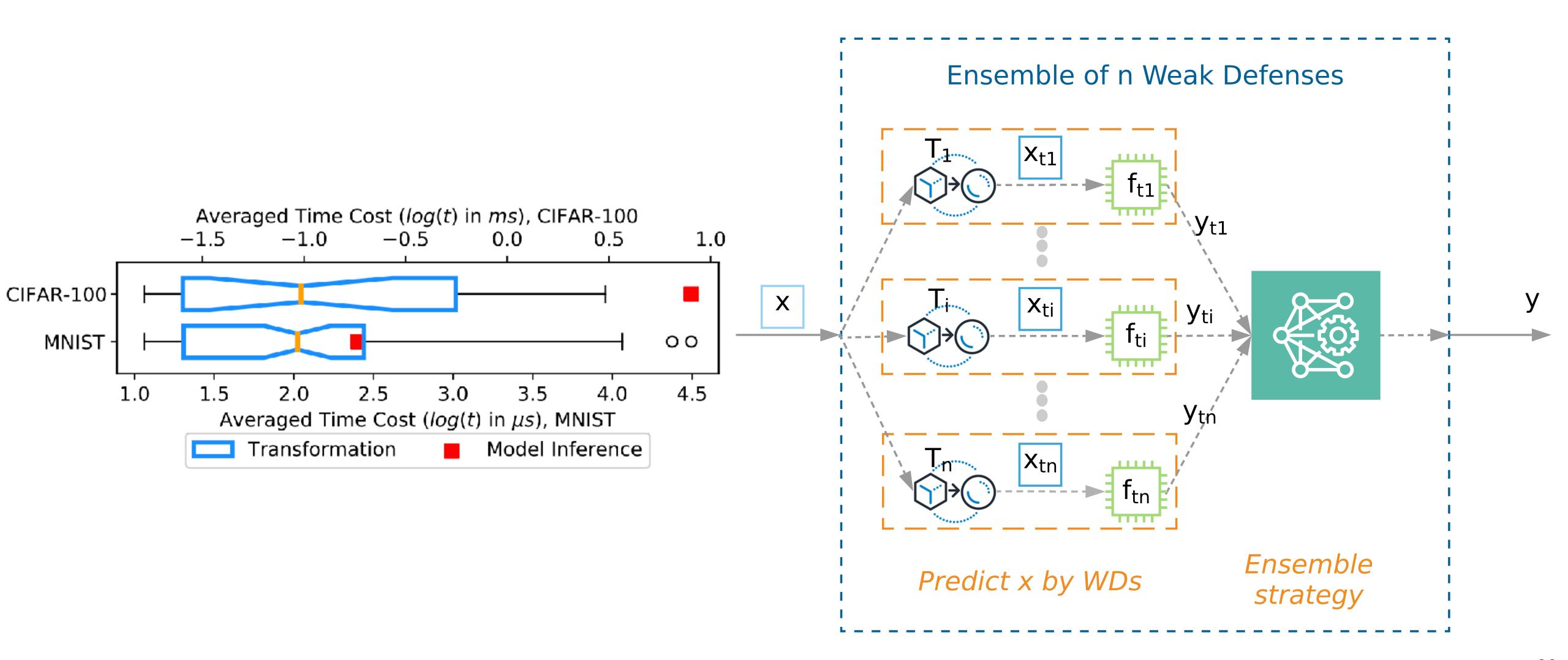
Dataset: MNIST

What is the overhead of ATHENA?

- Memory
- Inference Time



The memory overhead of ATHENA is linear with number of WDs, the inference time is on par with model inference



ATHENA is:

- Flexible
- Extensible
- General
- Moderate overhead





https://softsys4ai.github.io/athena/



A Framework based on Diverse Weak Defenses for Building Adversarial Defense

Ву

Ying Meng, Jianhai Su, Jason M O'Kane, and Pooyan Jamshidi

AlSys









<u>arXiv</u> <u>Preprint</u>

Code GitHub CSCE 585
Project

HelloWorld Tutorial