# Adversarial Examples



**Input**

Picture of a Cat

**Adv. crafting**

Careful perturbations of the input

**ML model**

A highly accurate Cats vs Dogs classifier

**Output**

Wrong prediction

# Important Milestones

**2014**
**Gradient-based Attacks**

**2016**
**Defences against Gradient Attacks**

**2017**
**Carlini & Wagner**

**2017**
**Adversarial training**

**2019**
**Certified Defenses**

**2021**
*Morphence*
*(this work)*

FGSM[1], BIM[2] and PGDM[3]

Defensive Distillation[4], Gradient Masking, etc

C&W[5]

Adv training[6]

Certified Defenses[7], PixelDP[8]

*Moving Target Defense*

[1] https://arxiv.org/abs/1412.6572

[2] https://arxiv.org/pdf/1611.01236.pdf

[3] https://arxiv.org/pdf/1706.06083.pdf

[4] https://arxiv.org/abs/1511.04508

[5] https://arxiv.org/abs/1608.04644

[6] https://arxiv.org/abs/1611.01236

[7] https://arxiv.org/abs/1705.07204

[8] https://arxiv.org/abs/1801.09344

[9] https://arxiv.org/abs/1802.03471

# Why Moving Target Defense?

## Fixed Model

- Highly vulnerable to model approximation

- Given enough time, the adversary will eventually find a way to evade it

- Repeated Attack: Once successful, it is always successful

## Moving Model

- Fitting the target model could be harder

- The defender is always one step ahead
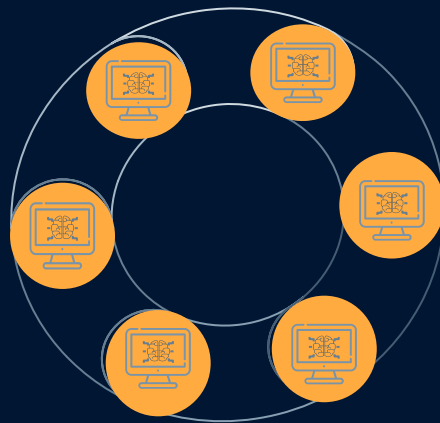
- An attack can succeed only once

# A Moving Model



**Input**

Picture of a Cat

**Adv. crafting**

Careful perturbations of the input

n = 6

**Moving Model**

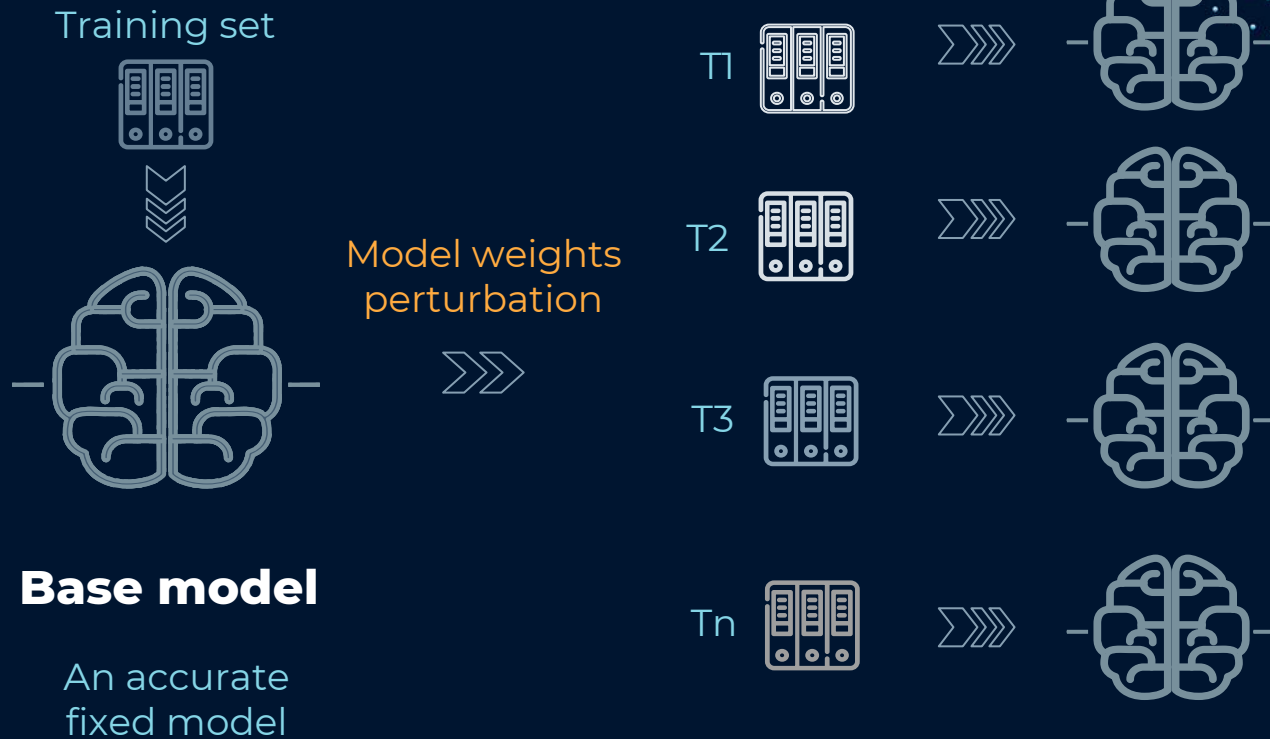A highly accurate Cats vs Dogs classifier

**Output**

Correct prediction

# Morphence

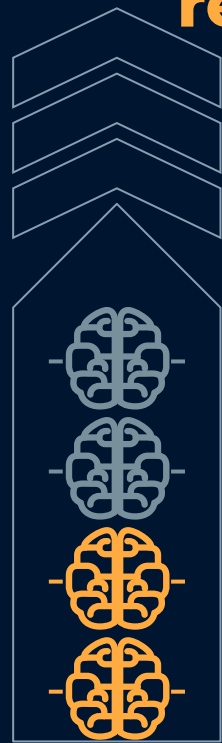Towards Moving Target Defenses against Adversarial Examples ...
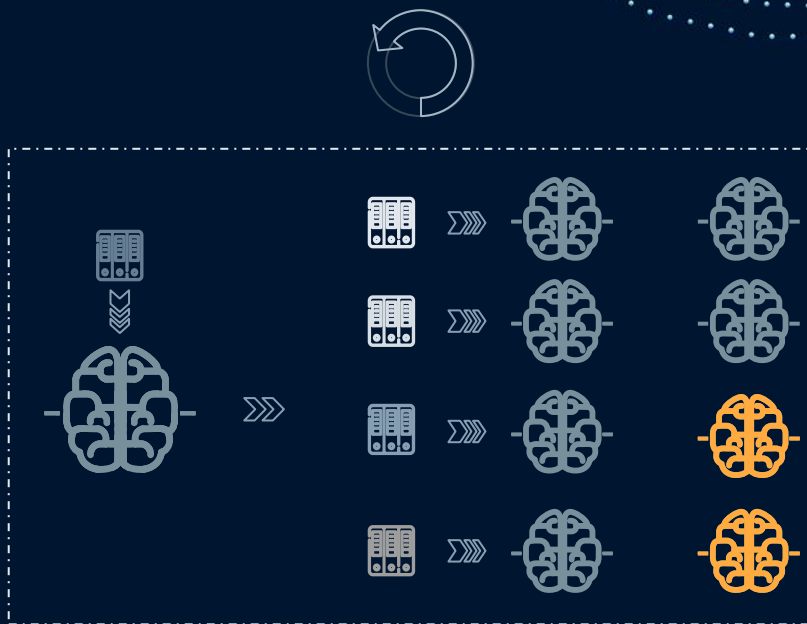
# Model Pool Generation

Training set

Model weights perturbation

**Base model**

An accurate fixed model

T1

T2

T3

Tn
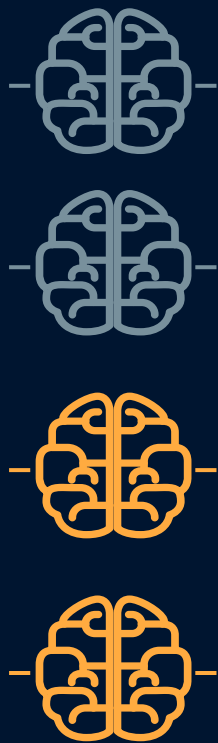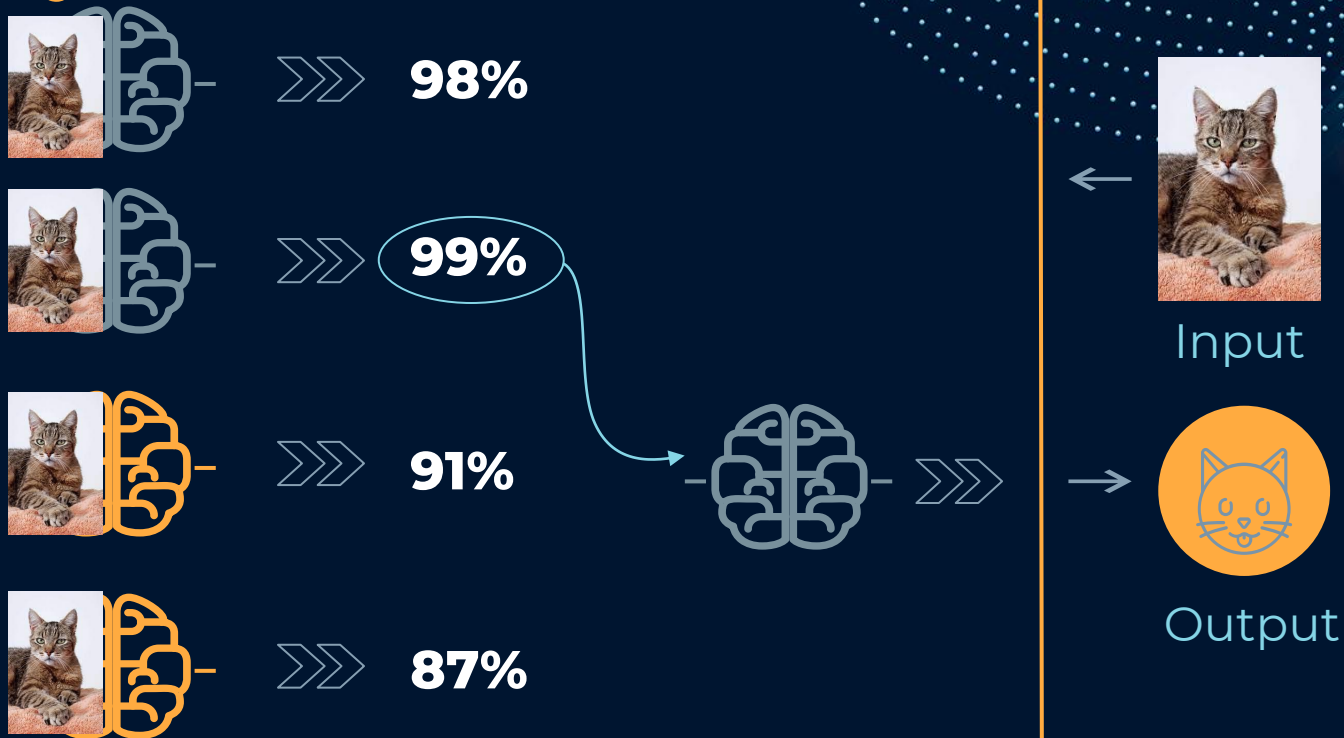
Another Layer of Robustness
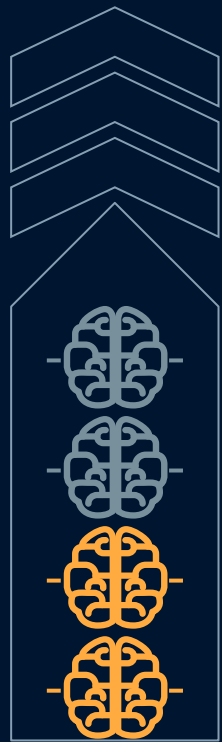
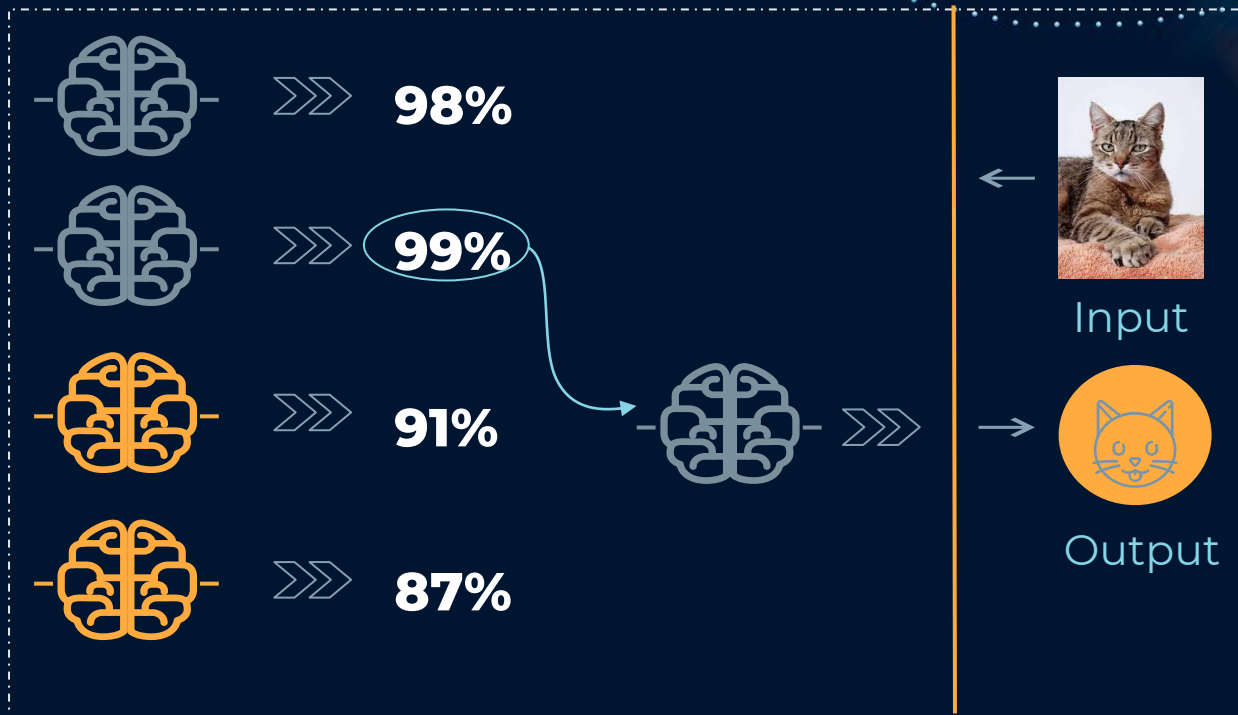# Local storage and repetitive generation

Stack

# Deployment

98%

99%

91%

87%

n models

Input

Output

# Model Pool Renewal

# queries = Qmax



98%

99%

91%

87%

Input

Output

Stack

# Results on MNIST

| | Undefended | Adversarially-trained | *Morphence* | |
|---|---|---|---|---|
| **No Attack** | 99.72% | 97.17% | **99.04%** | Do not sacrifice accuracy on benign data |
| **FGSM** | 9.98% | 42.38% | **71.43%** | Significant increase compared to adv training |
| **C&W** | 0.0% | 0.0% | **97.75%** | Overcomes C&W |
| **SPSA [10]** | 29.04% | 59.43% | **97.77%** | Robust against iterative-query attacks |

[10] https://arxiv.org/abs/1802.05666

# Results on CIFAR10

| | Undefended | Adversarially-trained | *Morphence* | |
|---|---|---|---|---|
| **No Attack** | 83.63% | 75.37% | **84.64%** | Can improve accuracy on benign data |
| **FGSM** | 9.98% | 36.62% | **38.78%** | Improvement compared to adv training |
| **C&W** | 1.25% | 1.34% | **44.50%** | Significant improvement on C&W |
| **SPSA** | 38.96% | 59.43% | **62.83%** | Higher robustness against iterative-query attacks |

# Robustness Against Repeated Attacks



Does successful attacks on pool-1 remain evasive on different pool of models?

# Detailed Results in the paper
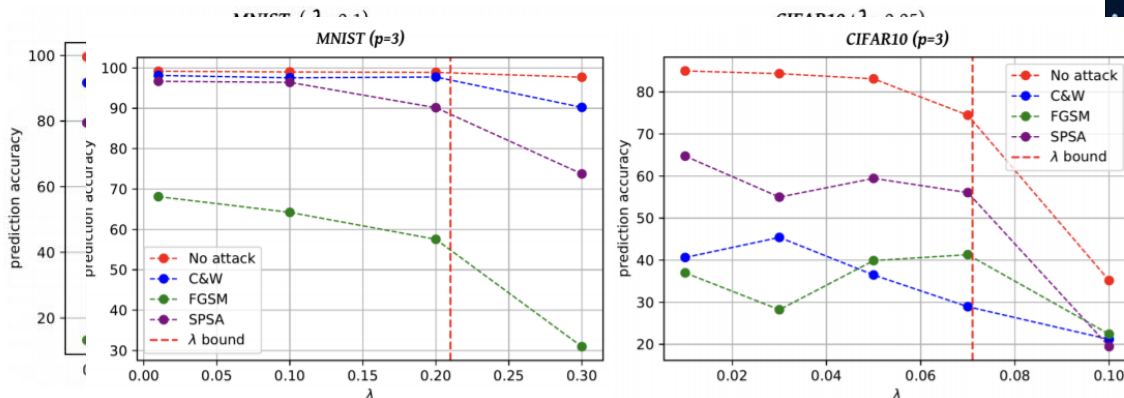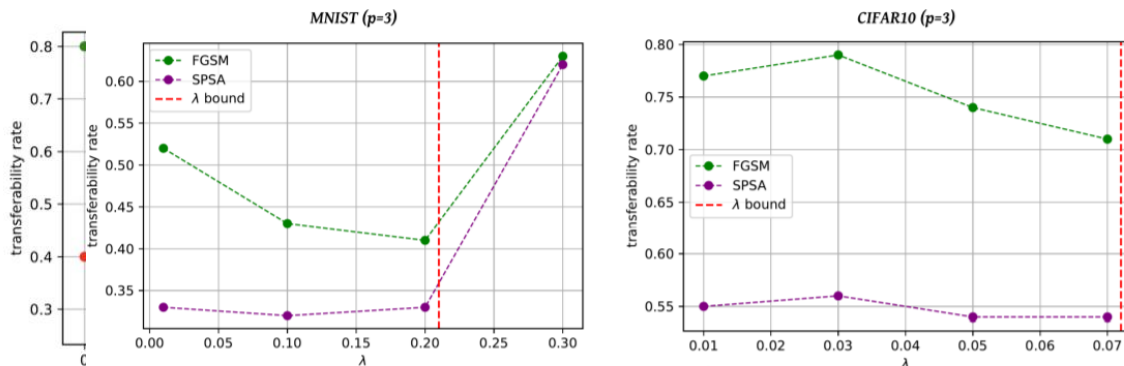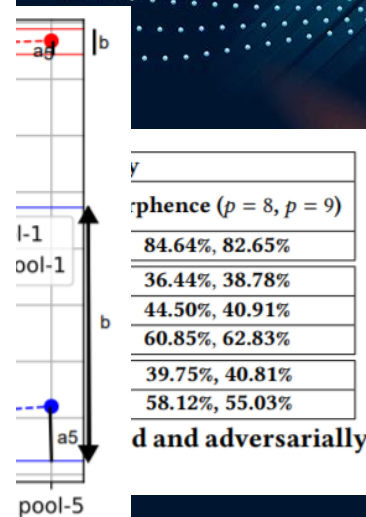


Fig. 5: Noise scale $\lambda$ vs. accuracy.

Fig. 6: Noise scale $\lambda$ vs. average transferability rate.

Fig. 4: # adversarially trained models $p$ vs. average transferability rate.
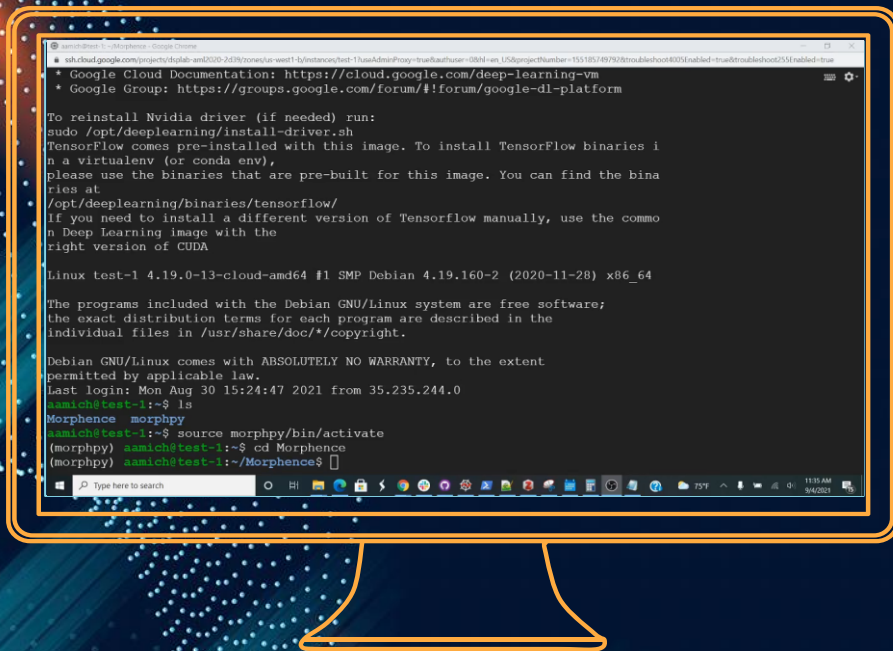
# Conclusions

**Morphence**

Moving Target Defense Against Adversarial Examples

- A Moving target model is more robust than the best fixed model defense.

- A Moving target model can prevent falling to the same attack multiple times.

- Iteratively querying a moving target model is not effective to optimize adversarial perturbations.

- We hope that Morphence will be used as a new benchmark for robustness against evasion attacks

# Available Artifact

https://github.com/um-dsp/Morphence



# THANKS!

## Do you have any questions?

 aamich@umich.edu

 @AbderrahmenAmi2

https://abderrahmen-amich.netlify.app/