# Adversarial Attacks on Neural Network Policies

Sandy Huang[1], Nicolas Papernot[2], Ian Goodfellow[3], Yan Duan[1,3], Pieter Abbeel[1,3]

[1]UC Berkeley, [2]Penn State, [3]OpenAI

**BAIR** BERKELEY ARTIFICIAL INTELLIGENCE RESEARCH

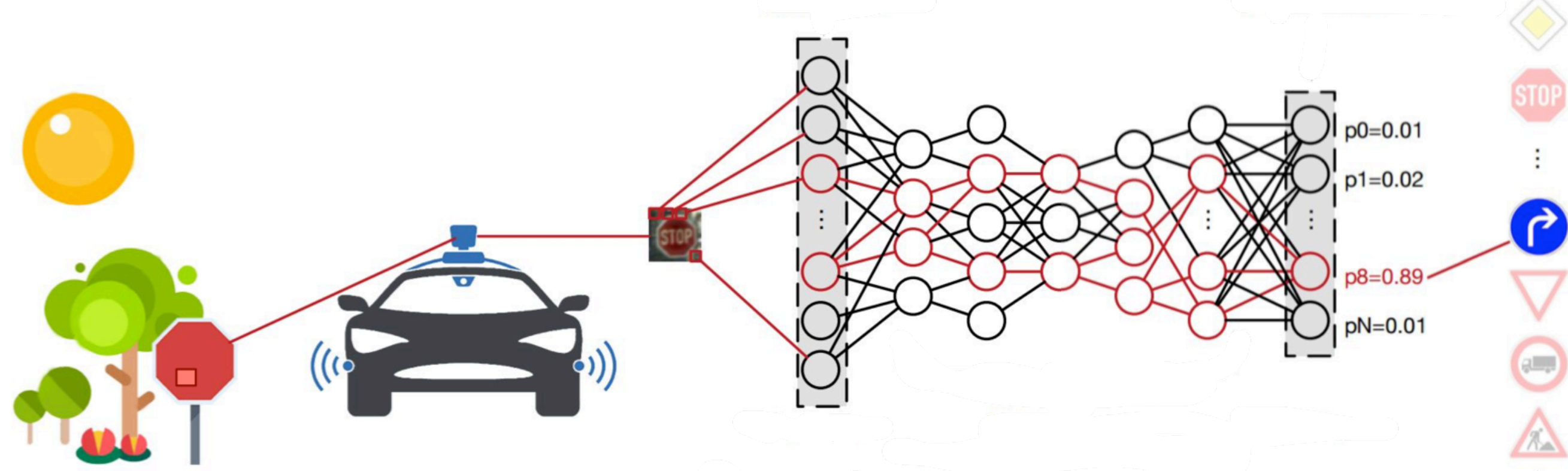PennState College of Engineering | ELECTRICAL ENGINEERING AND COMPUTER SCIENCE

OpenAI

## Motivation

*Adversarial example*: small worst-case perturbation that forces a machine learning model to mishandle an input

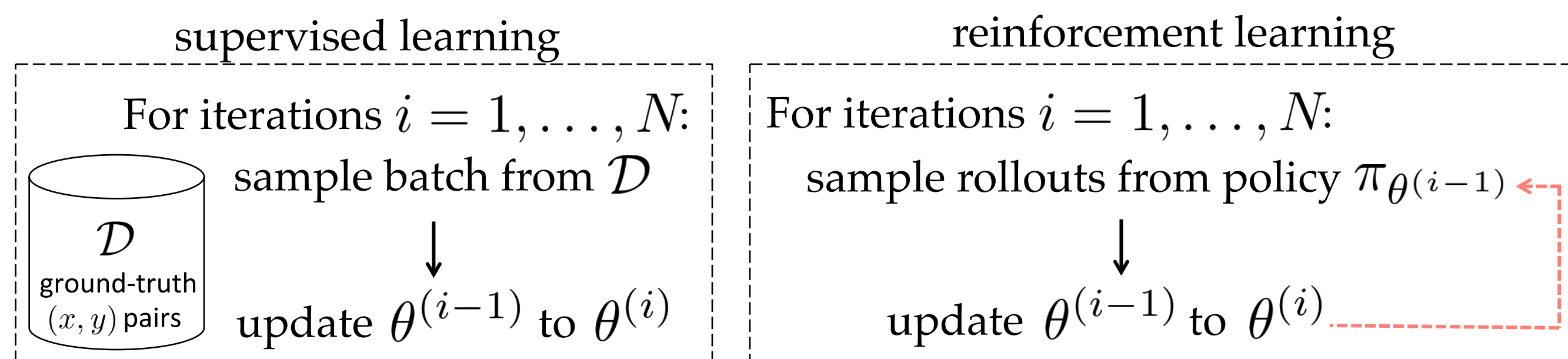- exists for image classification [1] and in the real world [2]



p0=0.01
p1=0.02
p8=0.89
pN=0.01

**Key finding**: Adversaries can degrade test-time behavior of policies trained with reinforcement learning, even when they do not have access to the policies

## Threat Model

white-box adversary: $\theta' = \theta$

black-box adversary:

- adversary trains its own policy $\pi_{\theta'}$ for the task

- requires *transferability*: can adversarial ex. designed to fool one policy also fool others trained for the same task?

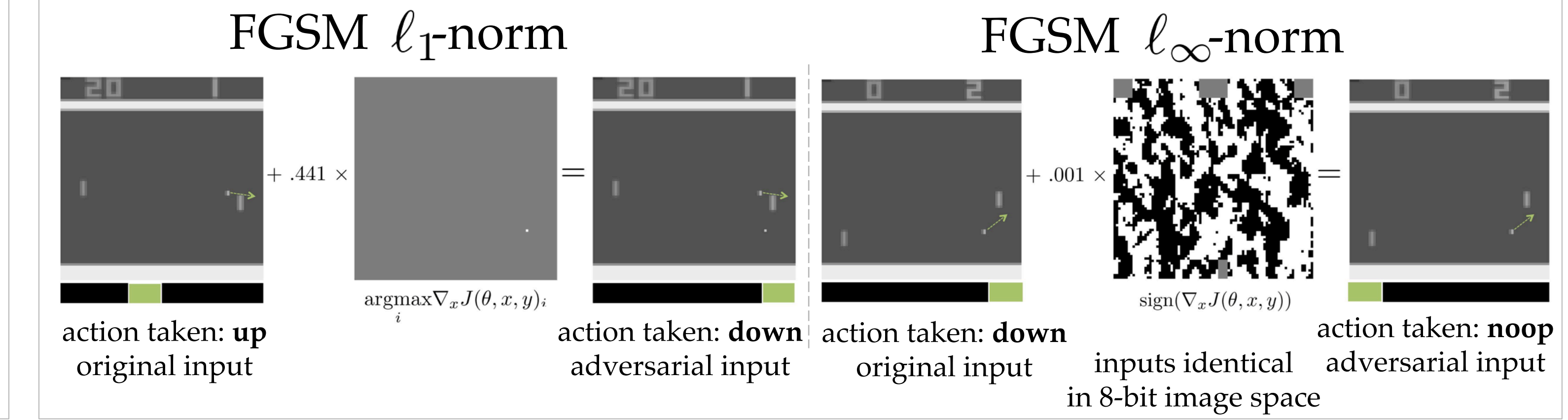- challenging because training data drawn from different distributions:



fully-trained

$x_t \rightarrow + \rightarrow \pi_\theta \rightarrow \tilde{a}_t$

$A_{\pi_{\theta'}} \rightarrow \eta_t$

possibly recurrent

supervised learning

For iterations $i = 1, \ldots, N$:
sample batch from $\mathcal{D}$
$\downarrow$
$\mathcal{D}$ ground-truth $(x, y)$ pairs
update $\theta^{(i-1)}$ to $\theta^{(i)}$

reinforcement learning

For iterations $i = 1, \ldots, N$:
sample rollouts from policy $\pi_{\theta(i-1)}$
$\downarrow$
update $\theta^{(i-1)}$ to $\theta^{(i)}$

## Dormant Adversarial Examples

We introduce dormant attacks (on recurrent policies):

| Time: | t | t+1 | ... | t+k-1 | t+k | |
|---|---|---|---|---|---|---|
| w/o adversary | ✔ | ✔ | ... | ✔ | ✔ | □ perturbation |
| w/ adversary | ✖ | ✖ | ... | ✖ | ✖ | ✔ optimal |
| w/ *dormant* adversary | ✔ | ✔ | ... | ✔ | ✖ | ✖ suboptimal |

## Examples of Adversarial Perturbations

FGSM $\ell_1$-norm



$+ .441 \times$ $\underset{i}{\mathrm{argmax}} \nabla_x J(\theta,x,y)_i$ $=$

action taken: **up** original input

action taken: **down** adversarial input

FGSM $\ell_\infty$-norm



$+ .001 \times$ $\mathrm{sign}(\nabla_x J(\theta,x,y))$ $=$

action taken: **down** original input

inputs identical in 8-bit image space

action taken: **noop** adversarial input

## Crafting Adversarial Examples for Policies

Optimal perturbation $\eta$, given loss $J(x)$: $\underset{\eta}{\mathrm{argmax}}\, J(x + \eta)$

$$J(\theta, x, y) = -\sum_i y_i \log \pi_\theta(x)_i = -\log \underset{i}{\mathrm{argmax}}\, \pi_\theta(x)_i$$

Fast gradient sign method (FGSM) [3] computes the optimal $\eta$ for the linear approximation of $J(x)$

Original version of FGSM constrains $\|\eta\|_\infty$
Instead, we might want to constrain sparsity / magnitude

$$\eta = \begin{cases} \epsilon\, \mathrm{sign}(\nabla_x J(\theta, x, y)) & \text{for } \|\eta\|_\infty \leq \epsilon \\ \epsilon\sqrt{d}\, \frac{\nabla_x J(\theta,x,y)}{\|\nabla_x J(\theta,x,y)\|_2} & \text{for } \|\eta\|_2 \leq \|\epsilon \mathbf{1}_d\|_2 \\ \text{maximally perturb dimensions with budget } \epsilon d \\ \qquad\qquad \text{for } \|\eta\|_1 \leq \|\epsilon \mathbf{1}_d\|_1 \end{cases}$$
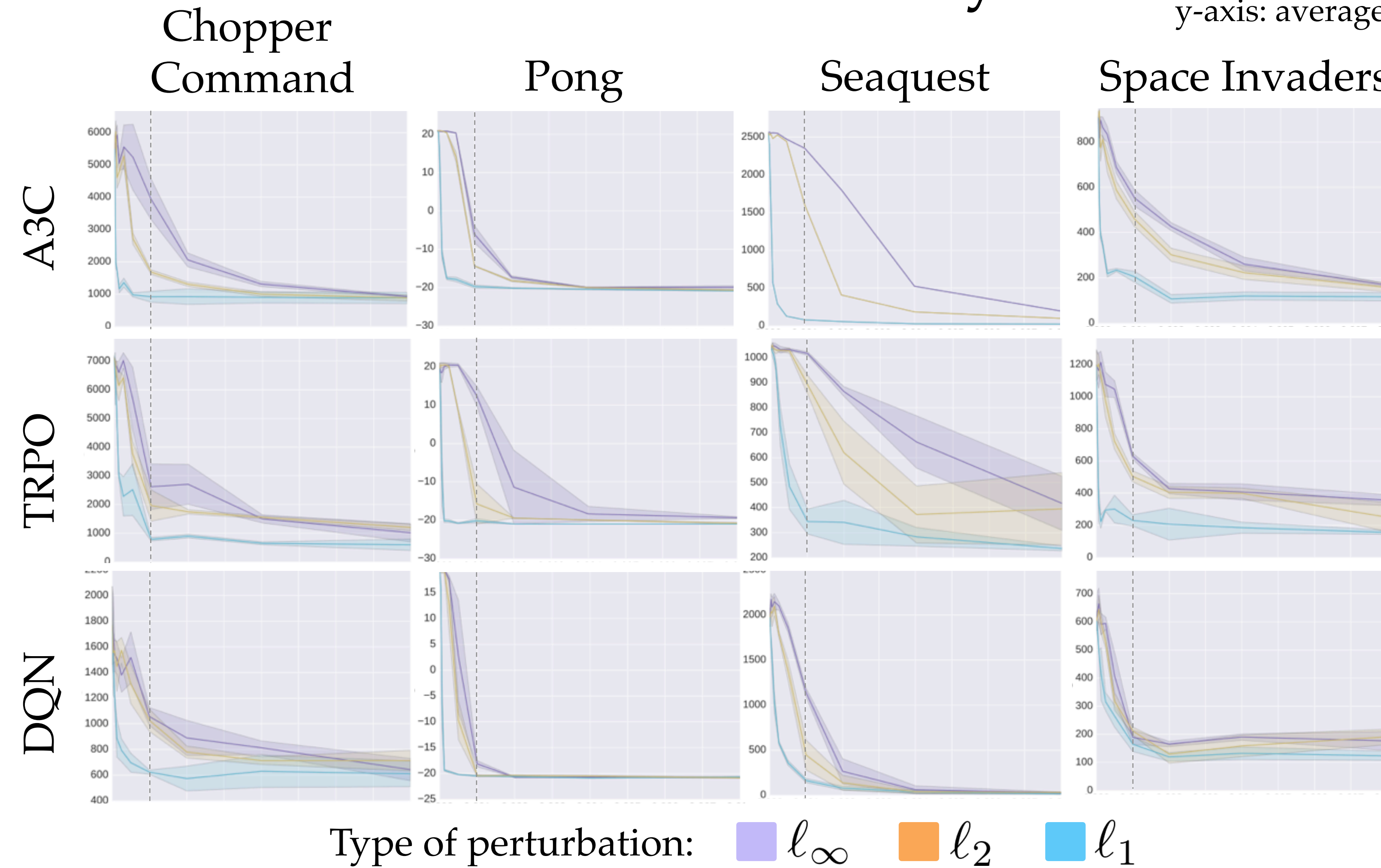
## Crafting Dormant Adversarial Examples

$\underset{\eta}{\mathrm{argmax}} \quad J(\theta, x_t + \eta, y_t)$

subject to $\quad d_i(\eta; \theta, x_{0:t}) = 0 \ \text{ for } \ i = 0, \ldots, k - 1,$
$\qquad\qquad d_k(\eta; \theta, x_{0:t}) = 1, \|\eta\| \leq \epsilon$

Optimize with dual ascent:
$$\eta^{(j)} = -\epsilon \left( \sum_{i=0}^{k-1} \lambda_i^{(j)} \mathrm{sign}(\nabla_x J(\theta, x_{t+i}, y_{t+i})) \right)$$
$$+ \epsilon\, \lambda_k^{(j)} \mathrm{sign}(\nabla_x J(\theta, x_{t+k}, y_{t+k})).$$
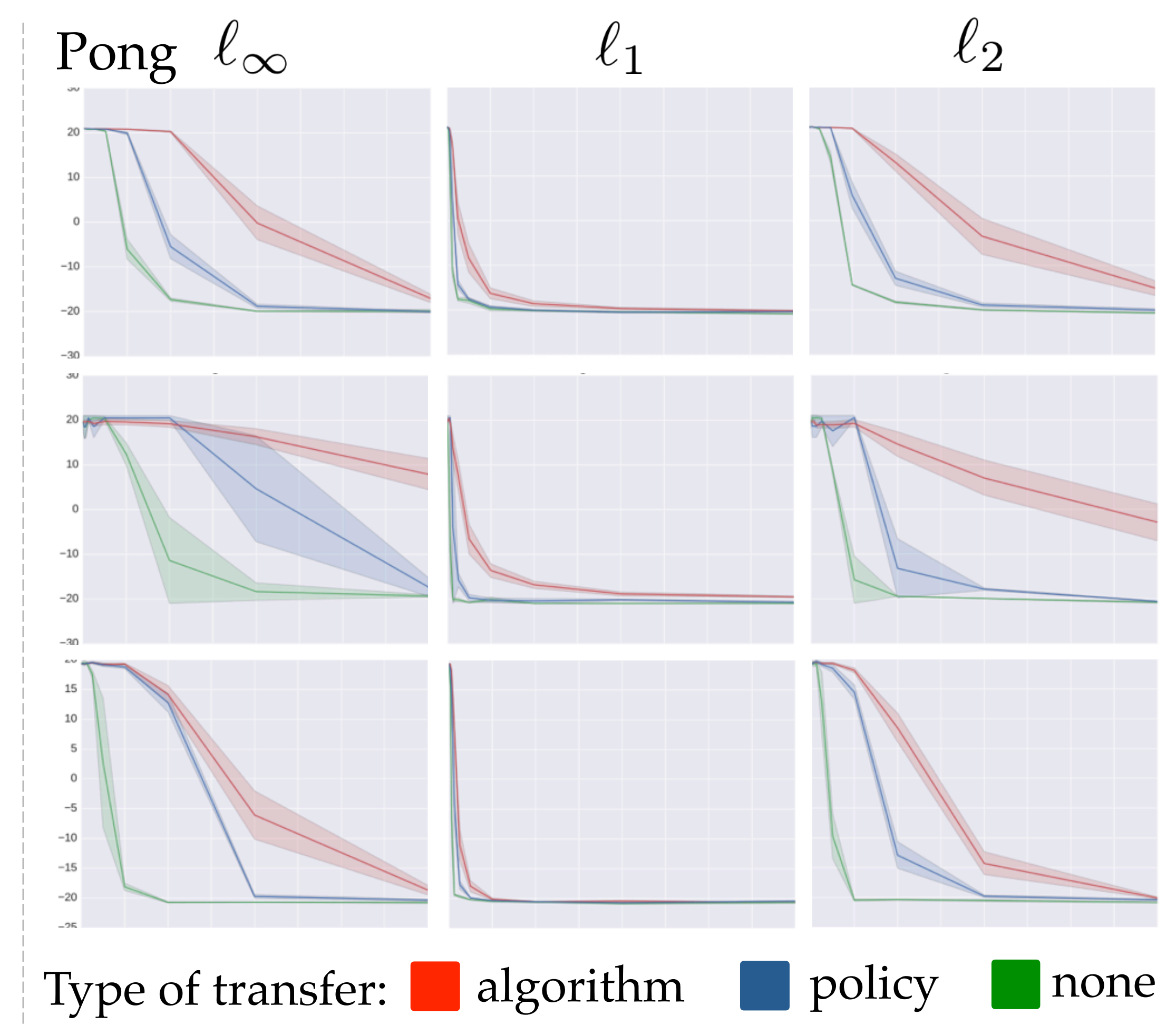$$\lambda_i^{(j+1)} = \lambda_i^{(j)} + \alpha^{(j)}\, d_i(\eta^{(j)}; \theta, x_{0:t}) \ \text{ for } \ i = 0, \ldots, k$$

## White-Box Adversary

x-axis: $\epsilon \in [0, 0.008]$
y-axis: average return

Chopper Command | Pong | Seaquest | Space Invaders



A3C

TRPO

DQN

Type of perturbation: $\ell_\infty$ $\ell_2$ $\ell_1$

- Across all games, adversarial perturbations significantly decrease performance, even for small $\epsilon$
- Policies trained with DQN tend to be more vulnerable
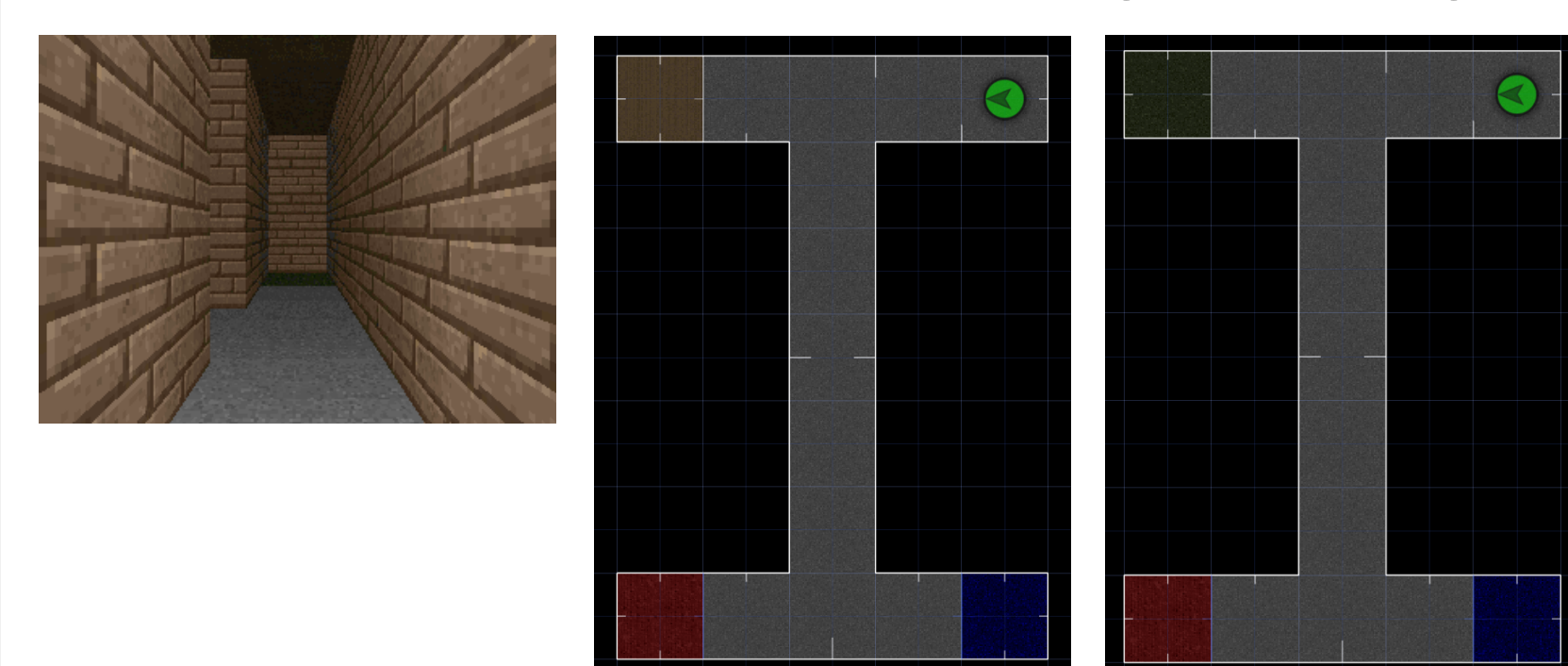
## Black-Box Adversary

Pong $\ell_\infty$ | $\ell_1$ | $\ell_2$



Type of transfer: algorithm policy none

- $\ell_1$-norm is particularly transferable, even across training algorithms

## Dormant Adversary

Task: Navigate I-maze [4]

in VizDoom [5]

yellow marker: go to red goal

green marker: go to blue goal



dormant adversarial examples computed through dual-ascent (k = 7)

$\ell_\infty$-norm



original input action at t+k: **turn right**

unscaled $\eta$

adversarial input action at t+k: **turn left**

$\ell_2$-norm



original input action at t+k: **turn right**

unscaled $\eta$

adversarial input action at t+k: **turn left**

$\ell_1$-norm



original input action at t+k: **turn right**

unscaled $\eta$

adversarial input action at t+k: **turn left**

[1] Szegedy et al. Intriguing properties of neural networks. ICLR 2014
[2] Kurakin et al. Adversarial examples in the physical world. arXiv 2016.
[3] Goodfellow et al. Explaining and harnessing adversarial examples. ICLR 2015.
[4] Oh et al. Control of memory, active perception, and action in Minecraft. ICML 2016.
[5] Kempka et al. ViZDoom. arXiv 2016.