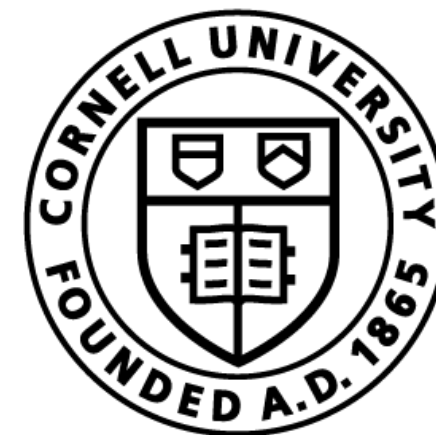


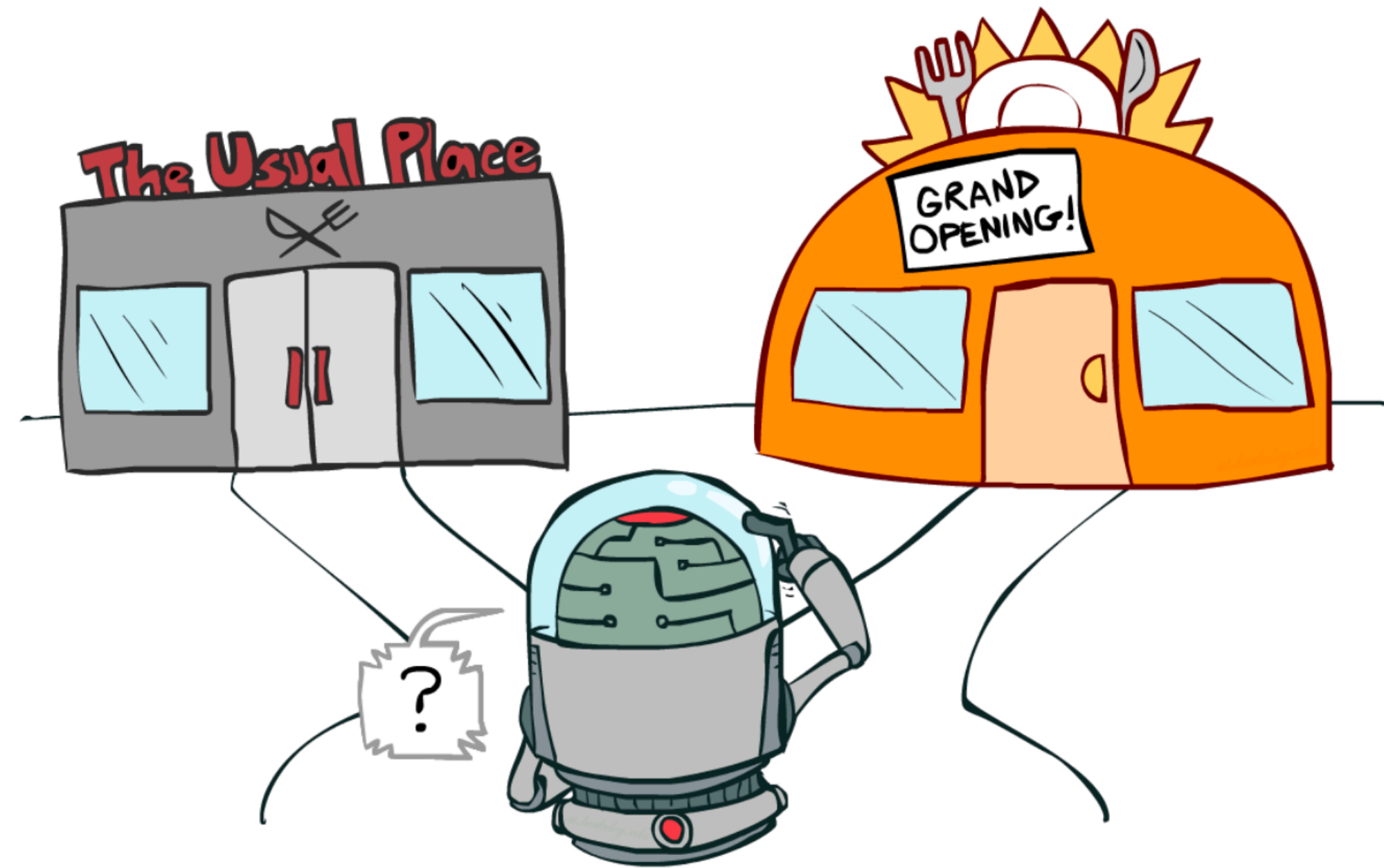
# Black-box vs White-box Policy Optimization

Sanjiban Choudhury

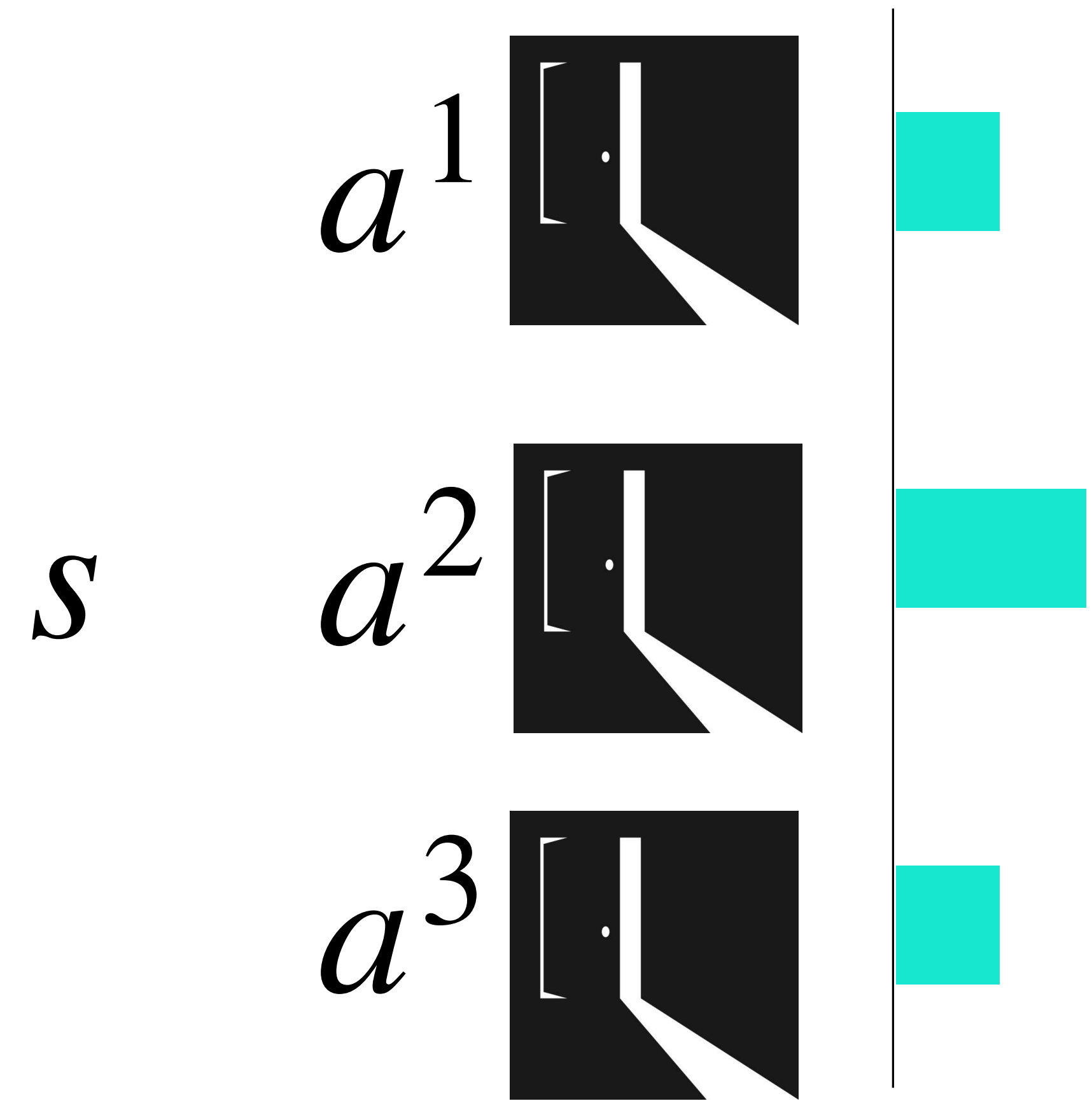


Cornell Bowers CIS  
**Computer Science**

# Recap: Two Ingredients of RL



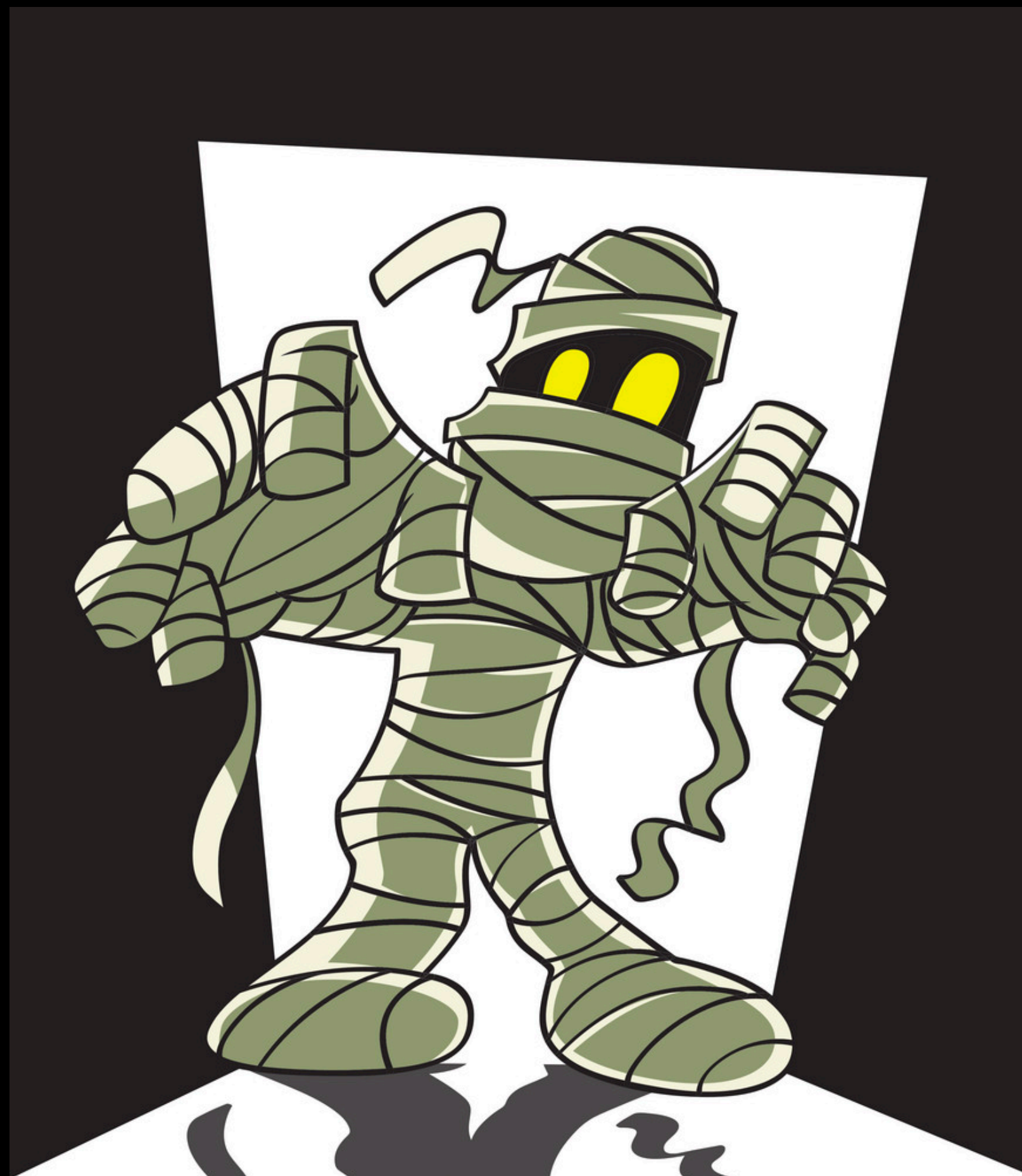
Exploration   Exploitation



Estimate Values  $Q(s, a)$

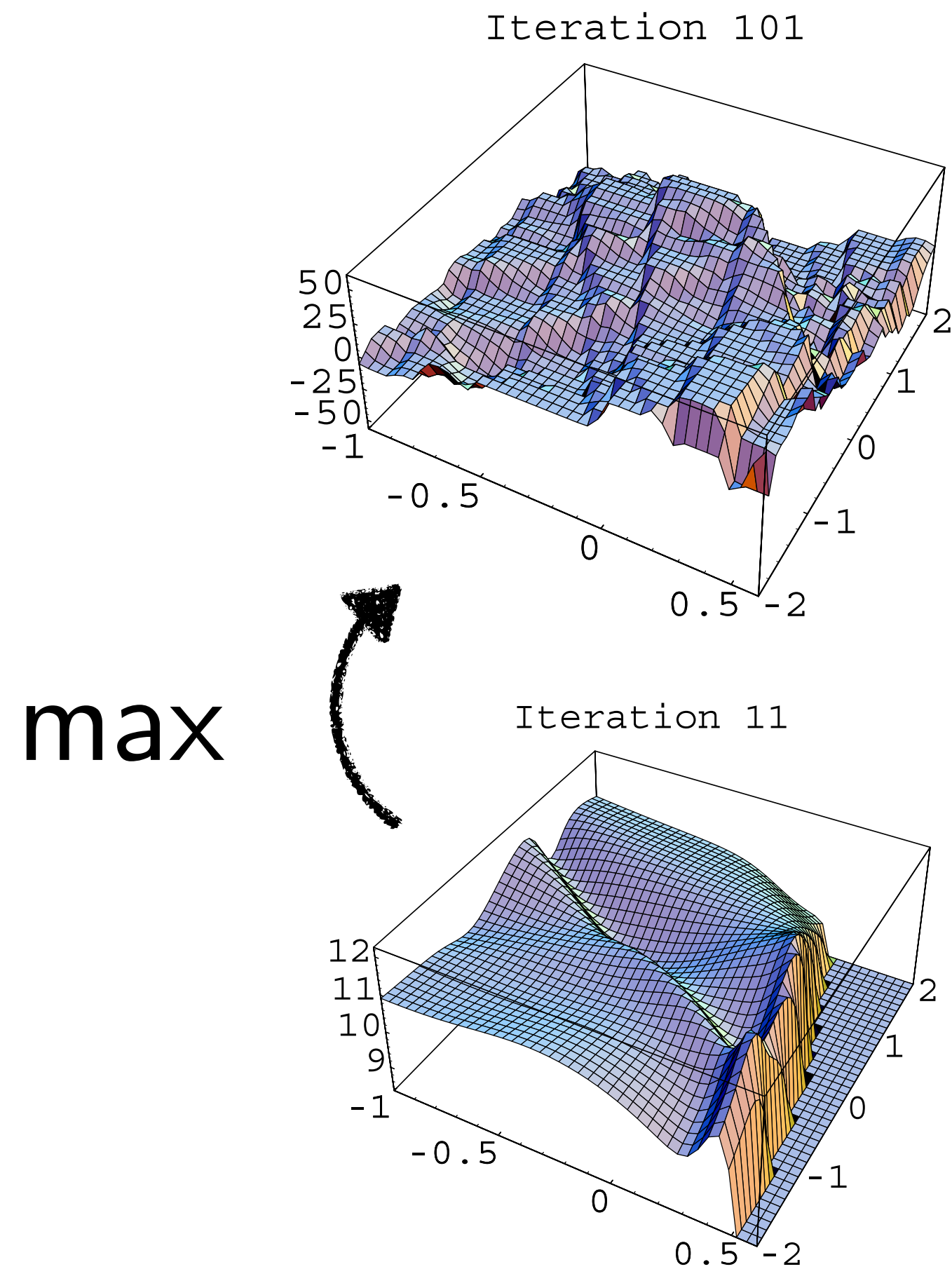


# CURSE OF VALUE APPROXIMATION!

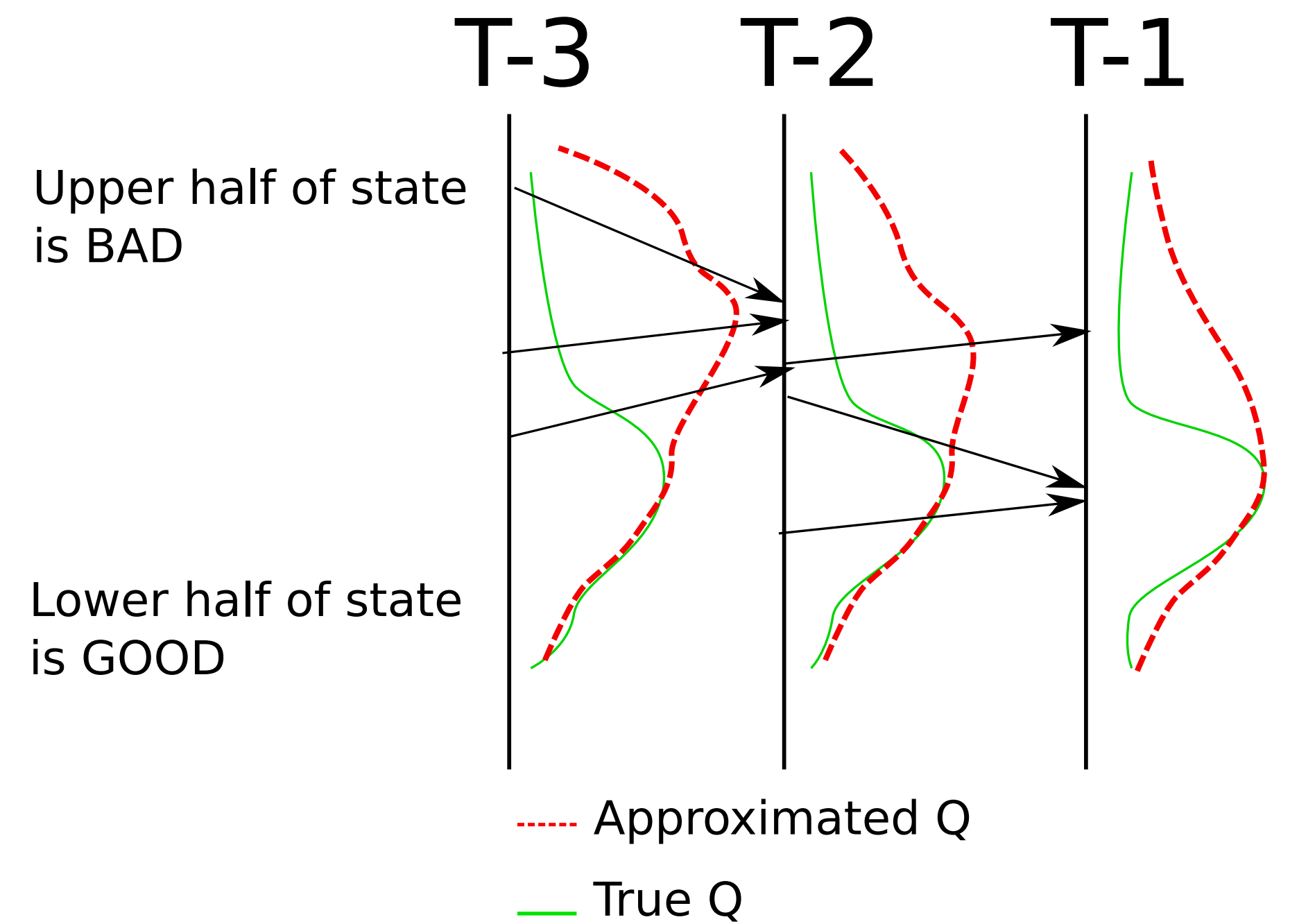


# Two sides of the same coin

## Bootstrapping

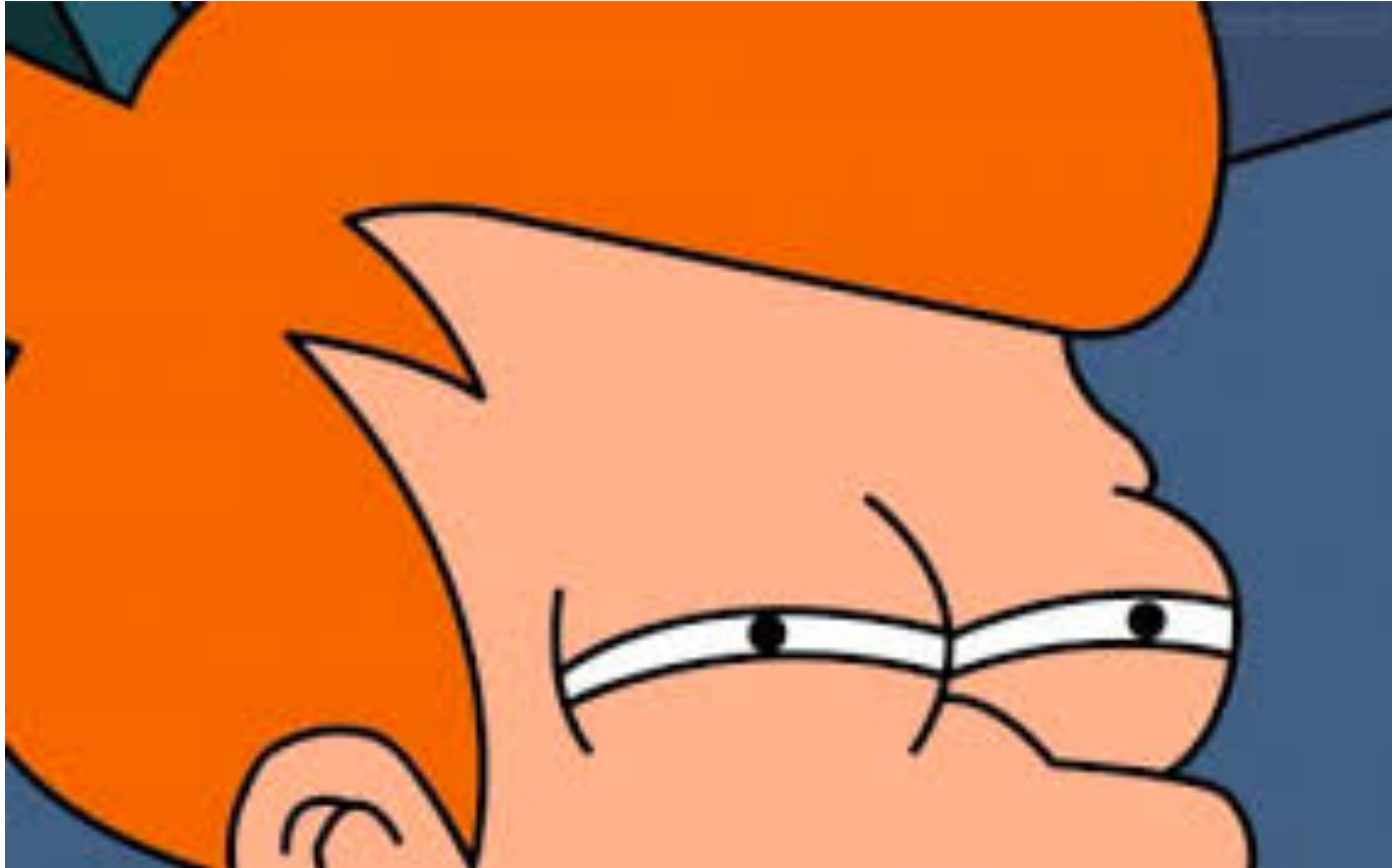


## Distribution shift





# To hell with Value Estimates!



## Trust ONLY actual Returns



Bye Bye Bellman ...

“not to be blinded by the  
beauty of the Bellman  
equation”

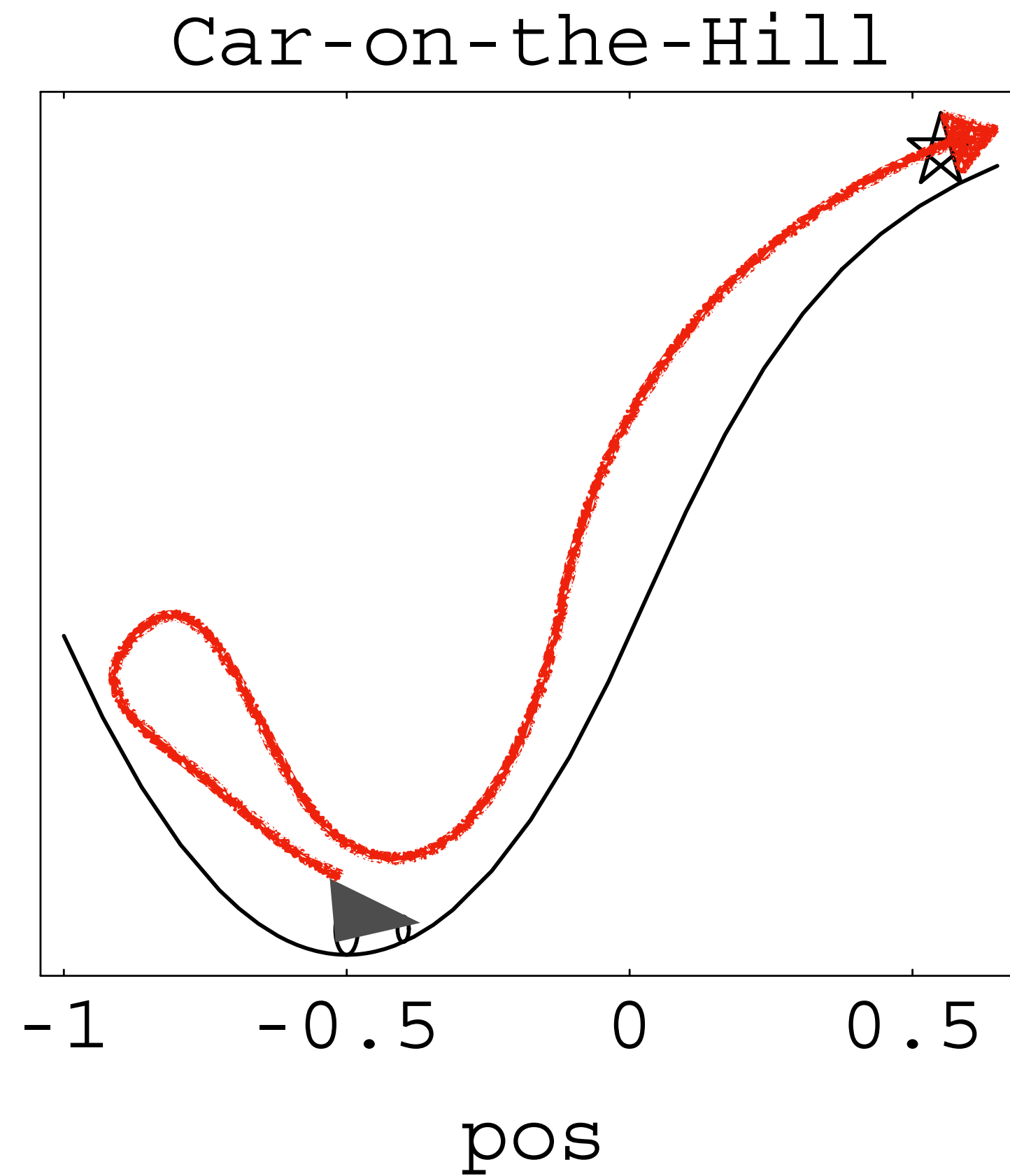
- Andrew Moore



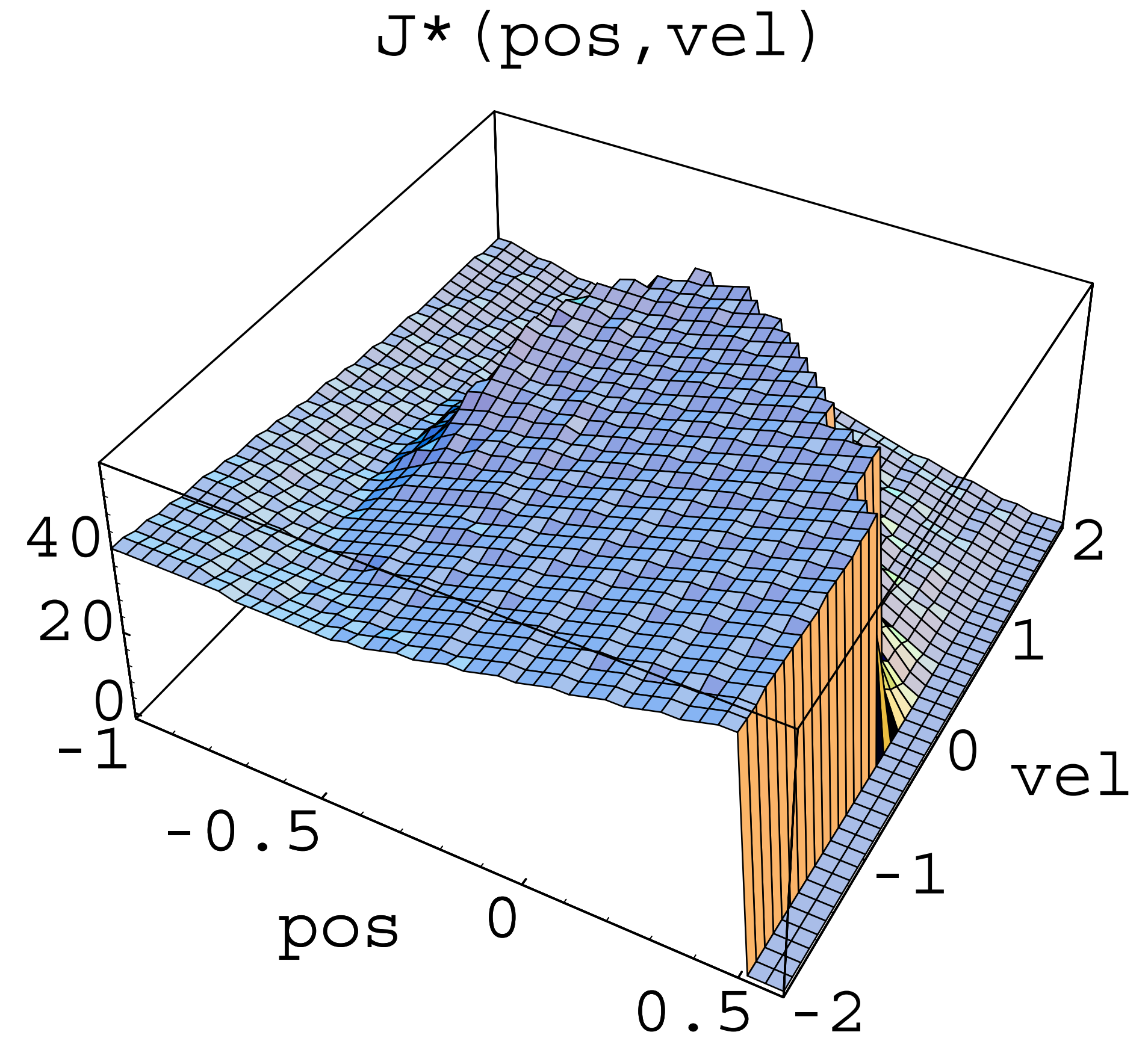
What if we focused on finding good policies ... ?



Sometimes a policy is waaaaay simpler than the value



The Policy!

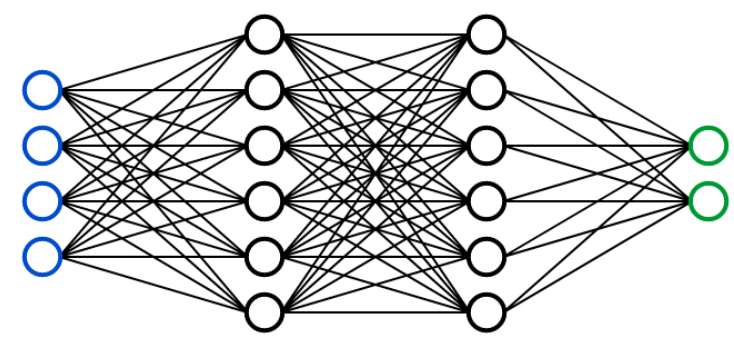


The Value!

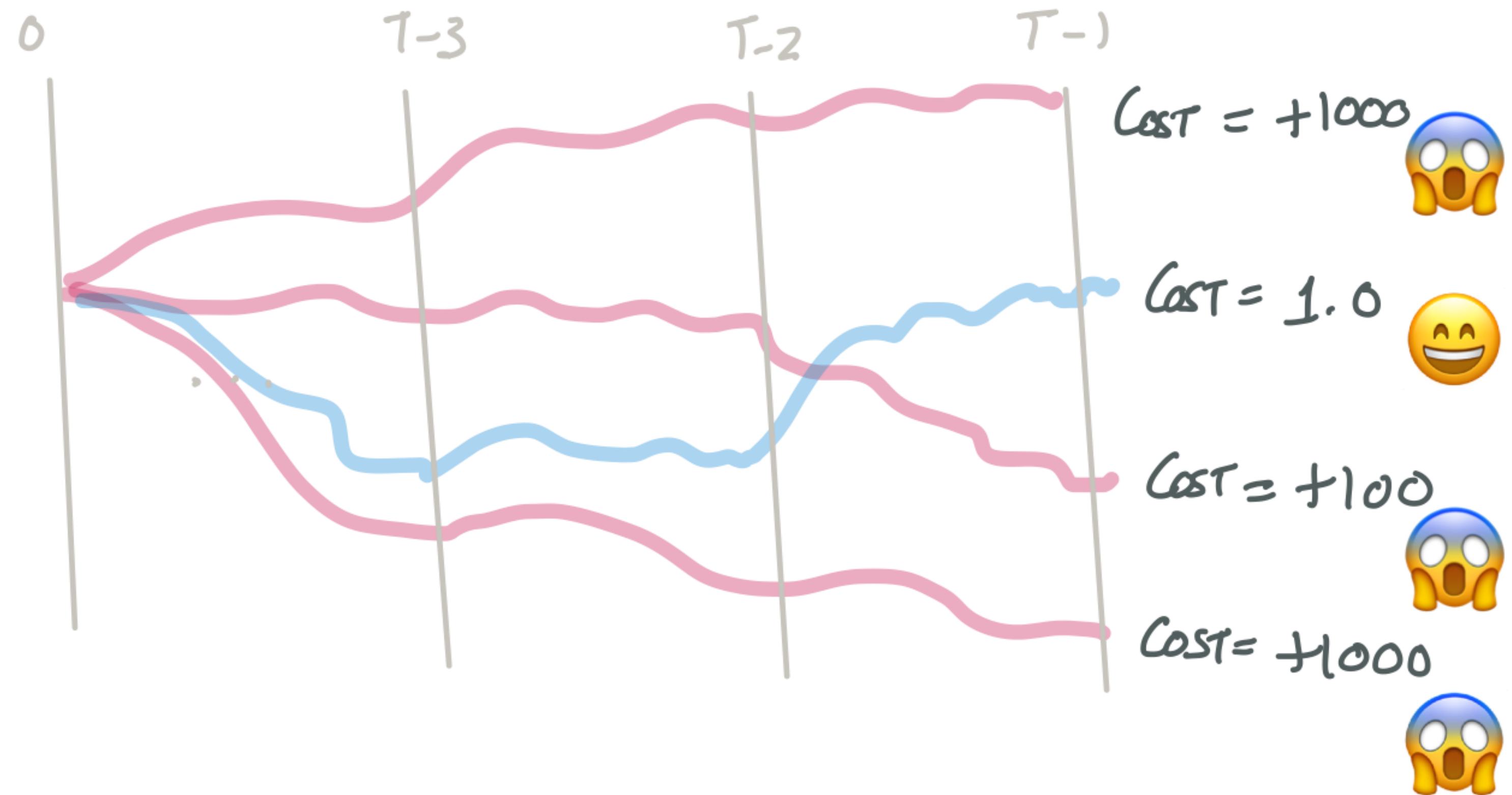


# Can we just focus on finding a good policy?

$$\pi_{\theta} : S_t \rightarrow a_t$$



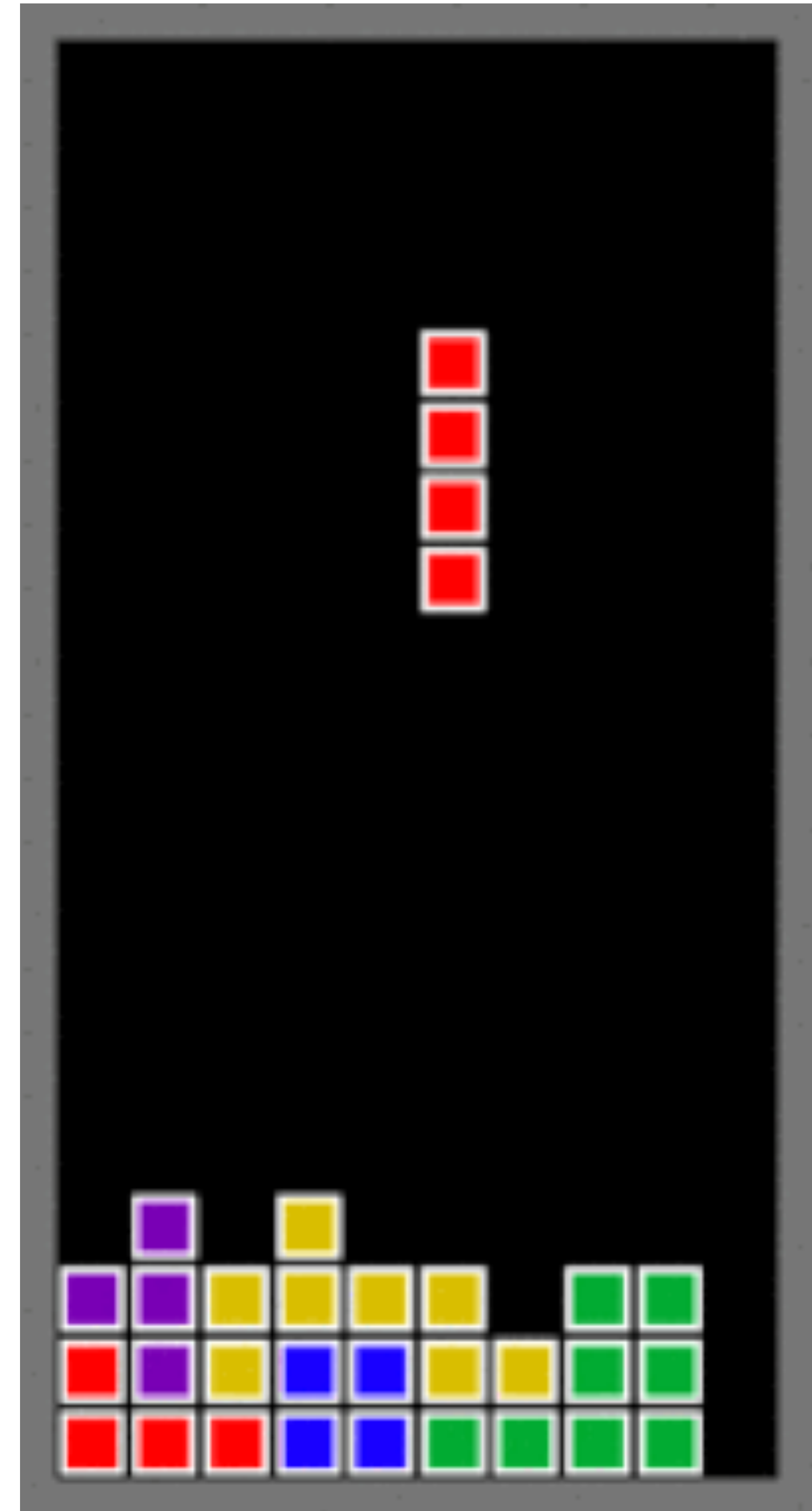
Learn a mapping from states to actions



Roll-out policies in the real-world to estimate value

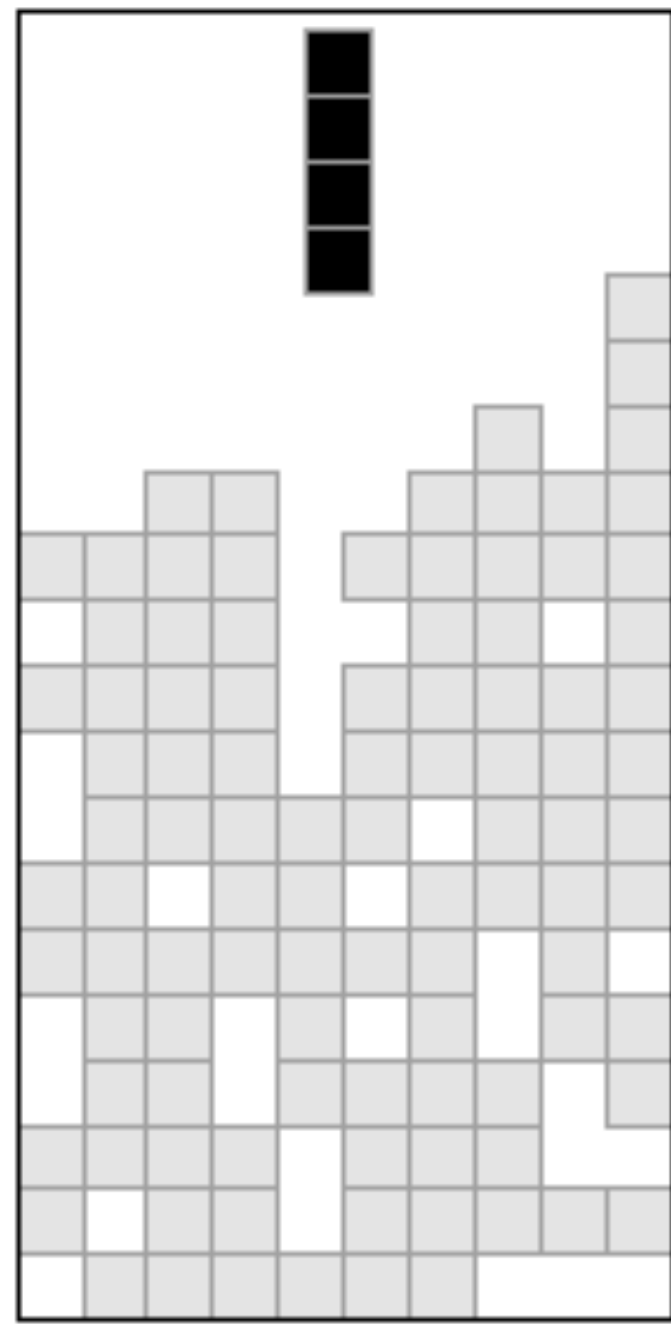


# The Game of Tetris



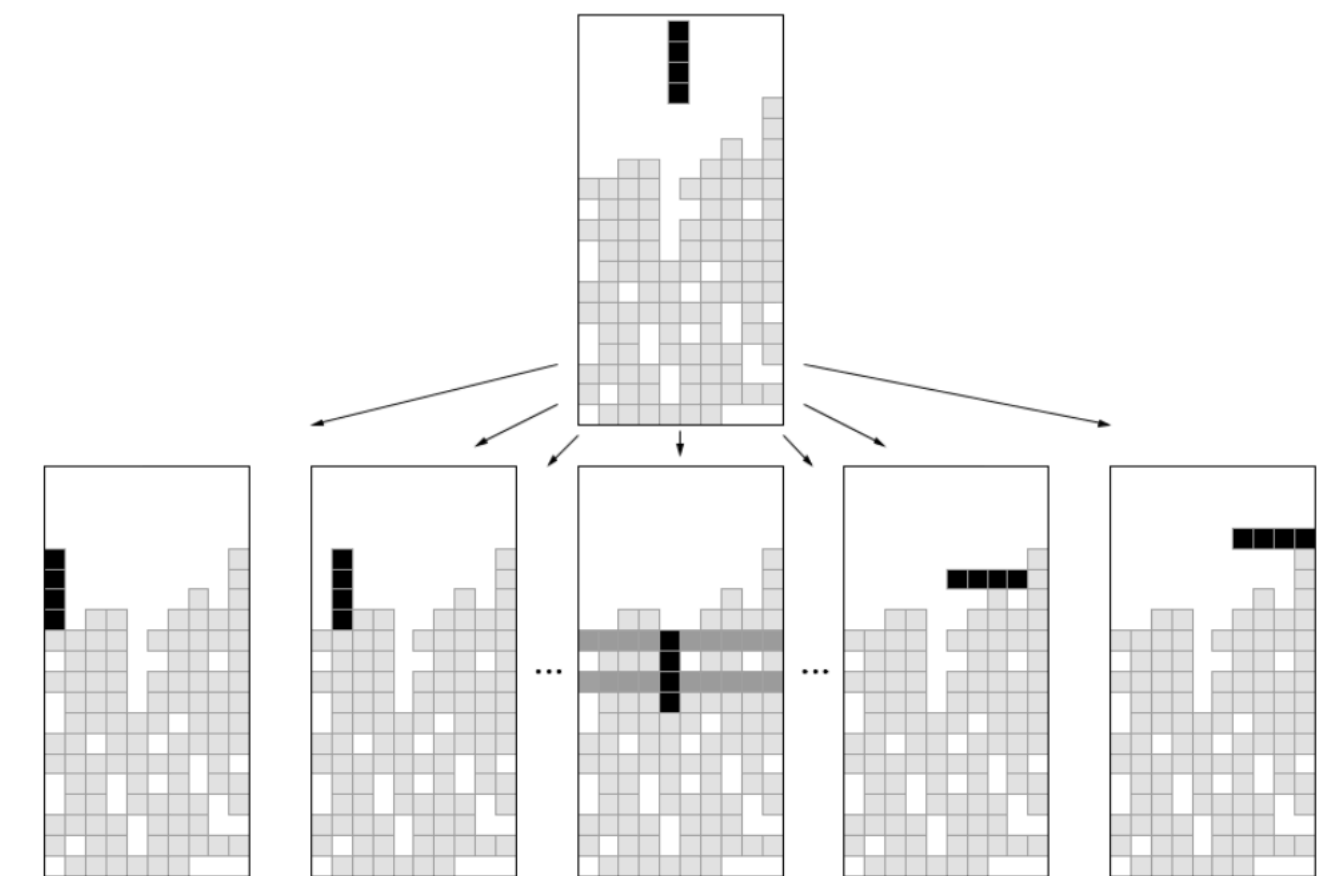
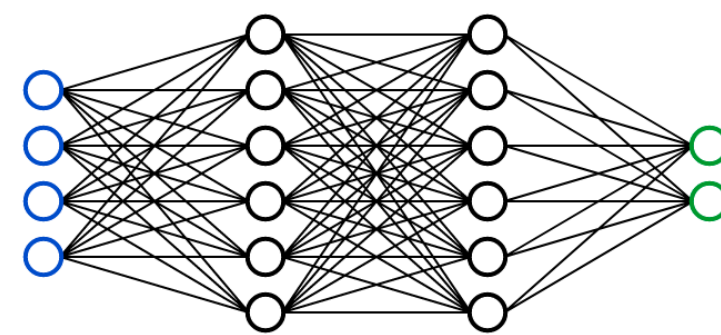
# What's a good policy representation for Tetris?

(4 rotations)\*(10 slots)  
- (6 impossible poses) = 34



State ( $s_t$ )

$$\pi_{\theta} : s_t \rightarrow a_t$$



Action ( $a_t$ )

Activity!

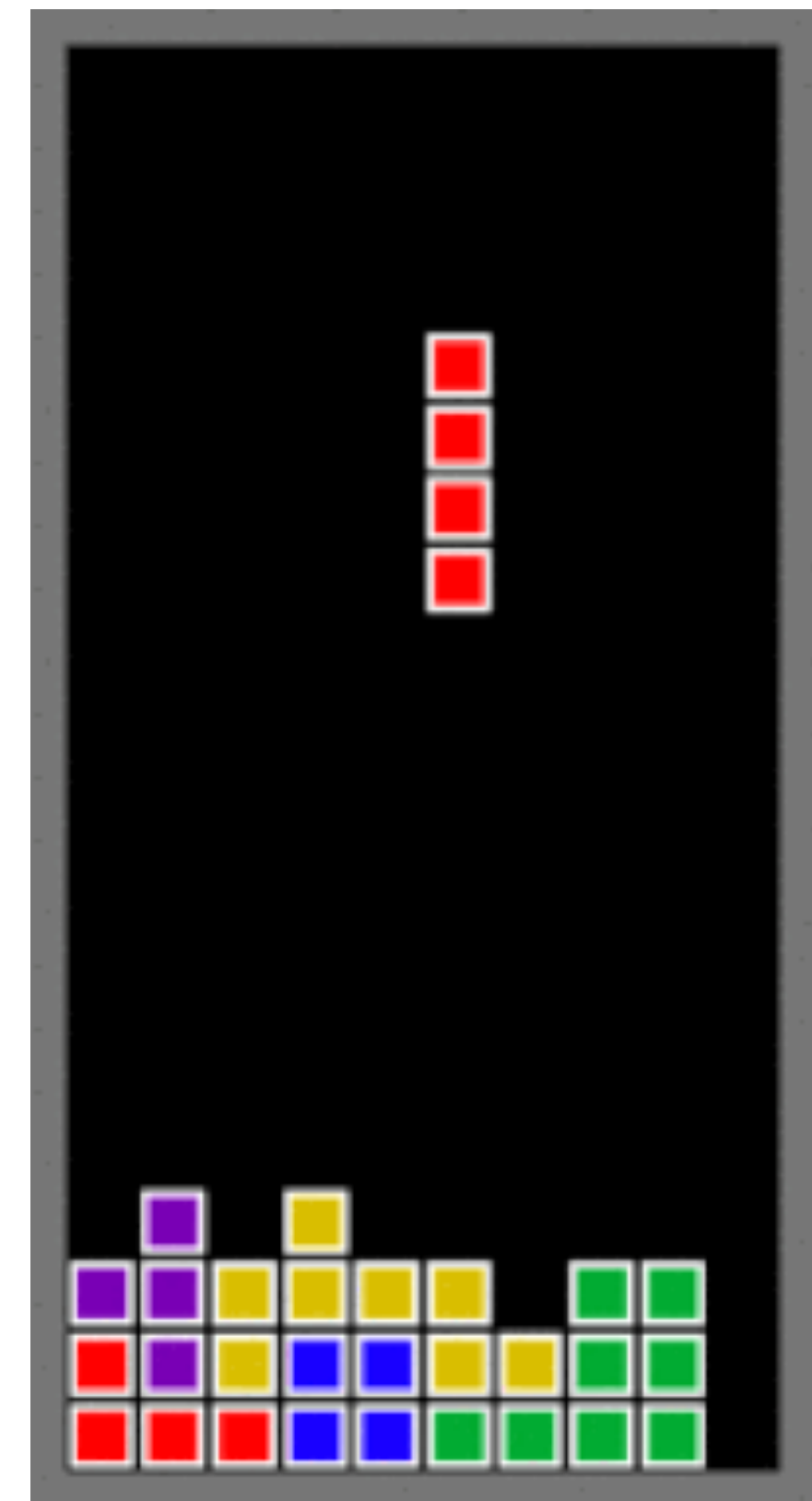


# Think-Pair-Share

Think (30 sec): Ideas for how to represent policy for tetris?

Pair: Find a partner

Share (45 sec): Partners exchange ideas

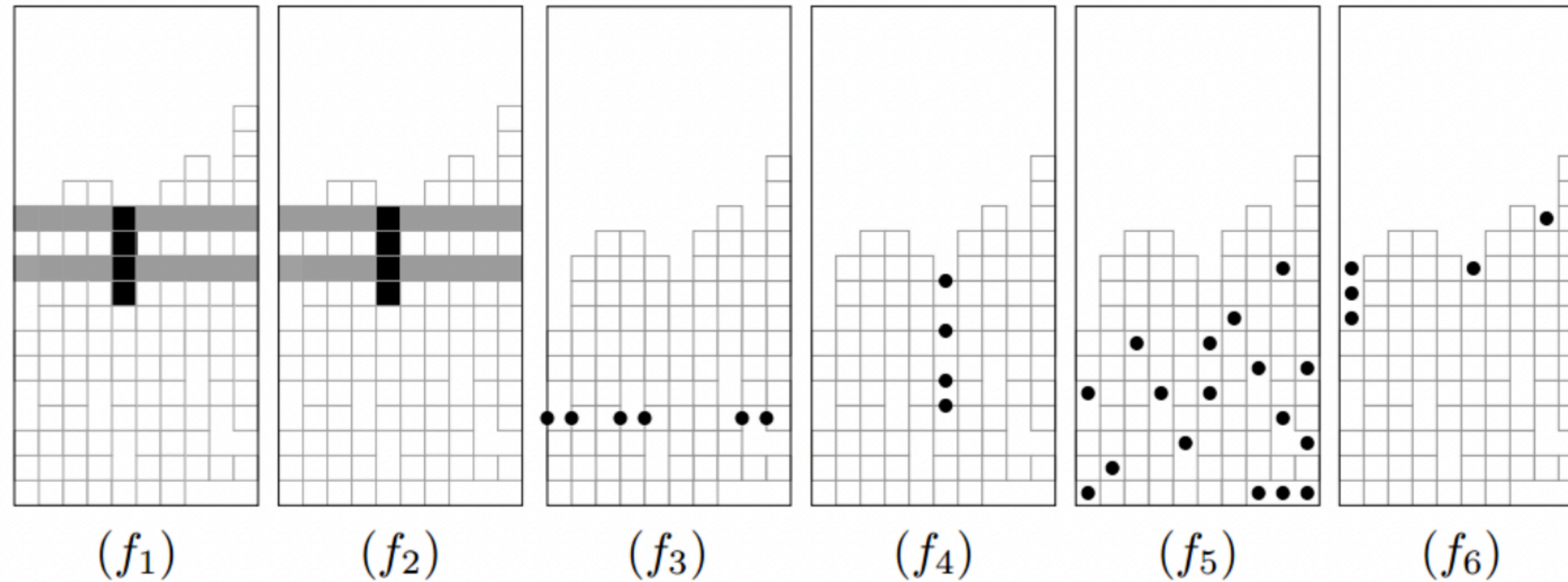


# Some inspiration for Tetris policy

*Until 2008, the best artificial Tetris player **was handcrafted**, as reported by Fahey (2003). Pierre Dellacherie, a self declared average Tetris player, identified six simple features and tuned the weights by trial and error.*



# Dellacherie Features



Landing  
Heights

Eroded  
Cells

Row  
Transitions

Column  
Transitions

Holes

Cumulative  
Wells

*The contribution of the last  
piece to the cleared lines  
time the number of cleared  
lines.*

*The number of filled cells  
adjacent to the empty cells  
summed over all rows*

*A well is a succession of  
empty cells and the cells to  
the left and right are  
occupied*

# A *magic* formula ?!?

- $4 \times$  holes – *cumulative wells*
- *row transitions* – *column transitions*
- *landing height* + *eroded cells*



# *A magic formula ?!?*

- 4 × holes – cumulative wells*
- row transitions – column transitions*
- landing height + eroded cells*

*This linear evaluation function cleared an **average of 660,000 lines** on the full grid ...*

*... In the simplified implementation used by the approaches discussed earlier, the games would have continued further, until every placement would overflow the grid. Therefore, this report underrates this simple linear rule compared to other algorithms.*

Can YOU do better  
than Dellacherie?



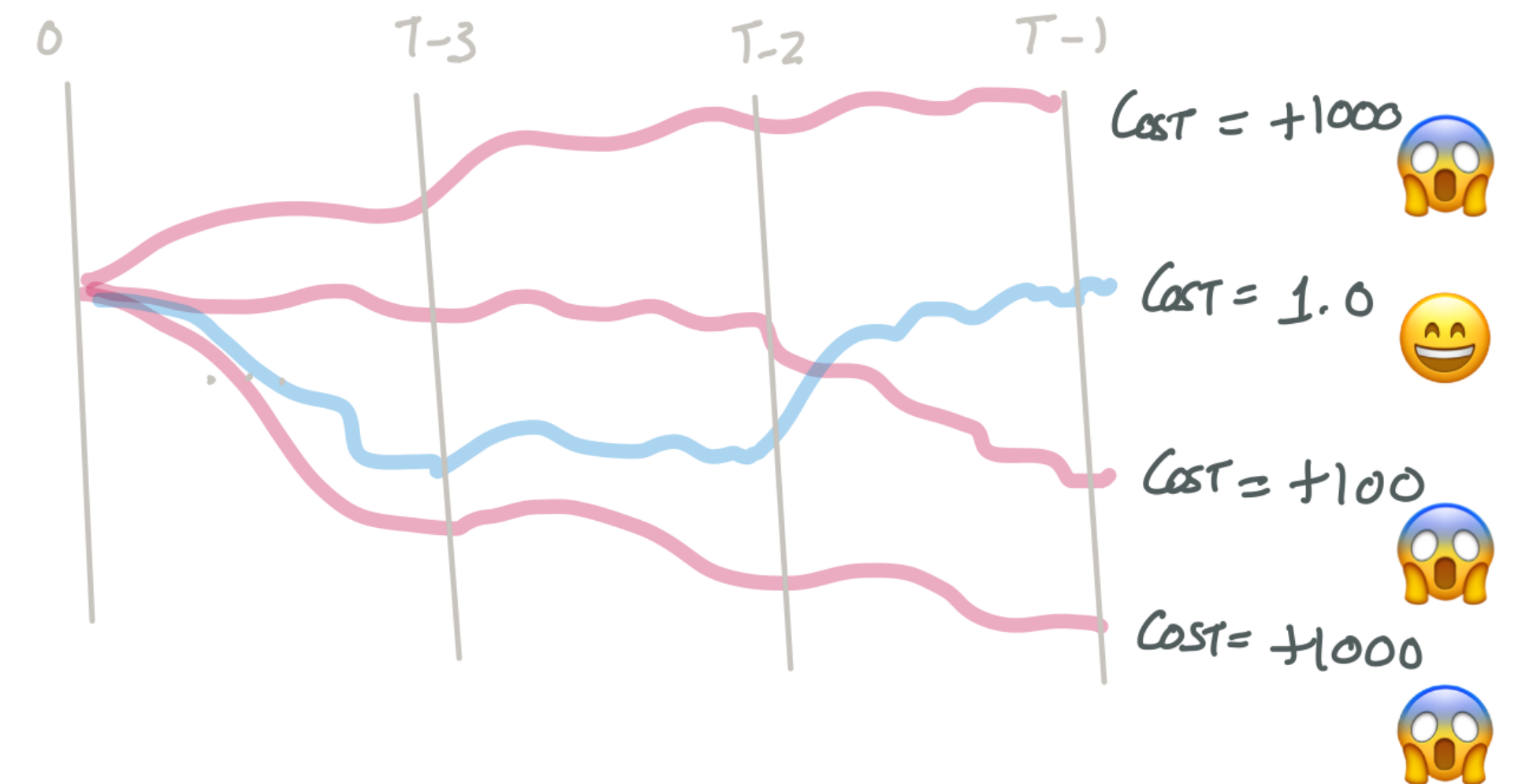
# The Goal of Policy Optimization

$$\pi_{\theta}(s) = \arg \min_a \theta^T f(s, a)$$

#Think of  $f(s,a)$  being dellacherie features

$$\min_{\theta} J(\theta) = \sum_{t=0}^{T-1} \mathbb{E}_{\pi_{\theta}} c(s_t, a_t)$$

#Think of  $c(s,a)$  as  
-num\_rows\_cleared



# Is Policy Optimization a good idea?



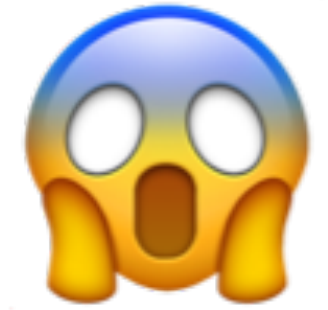
## Pros

A policy is simpler than value function

Easy to bake in engineering knowledge

Easy to code up!

## Cons



Careful feature engineering

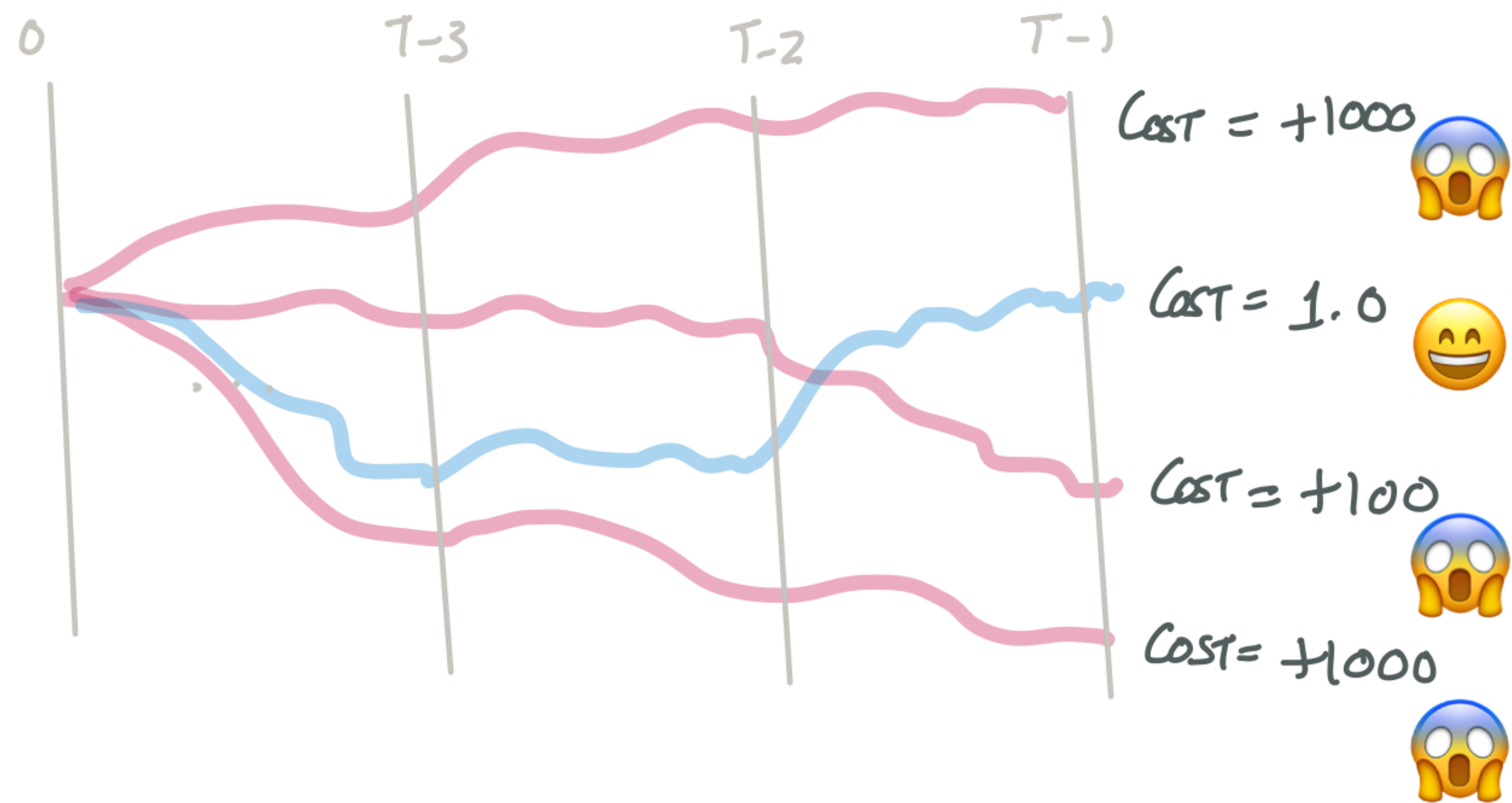
Exploration is difficult in this setting

Ignoring Markov structure of states and costs



# What are some ways to solve this optimization?

$$\min_{\theta} J(\theta) = \sum_{t=0}^{T-1} \mathbb{E}_{\pi_{\theta}} c(s_t, a_t)$$



Activity!



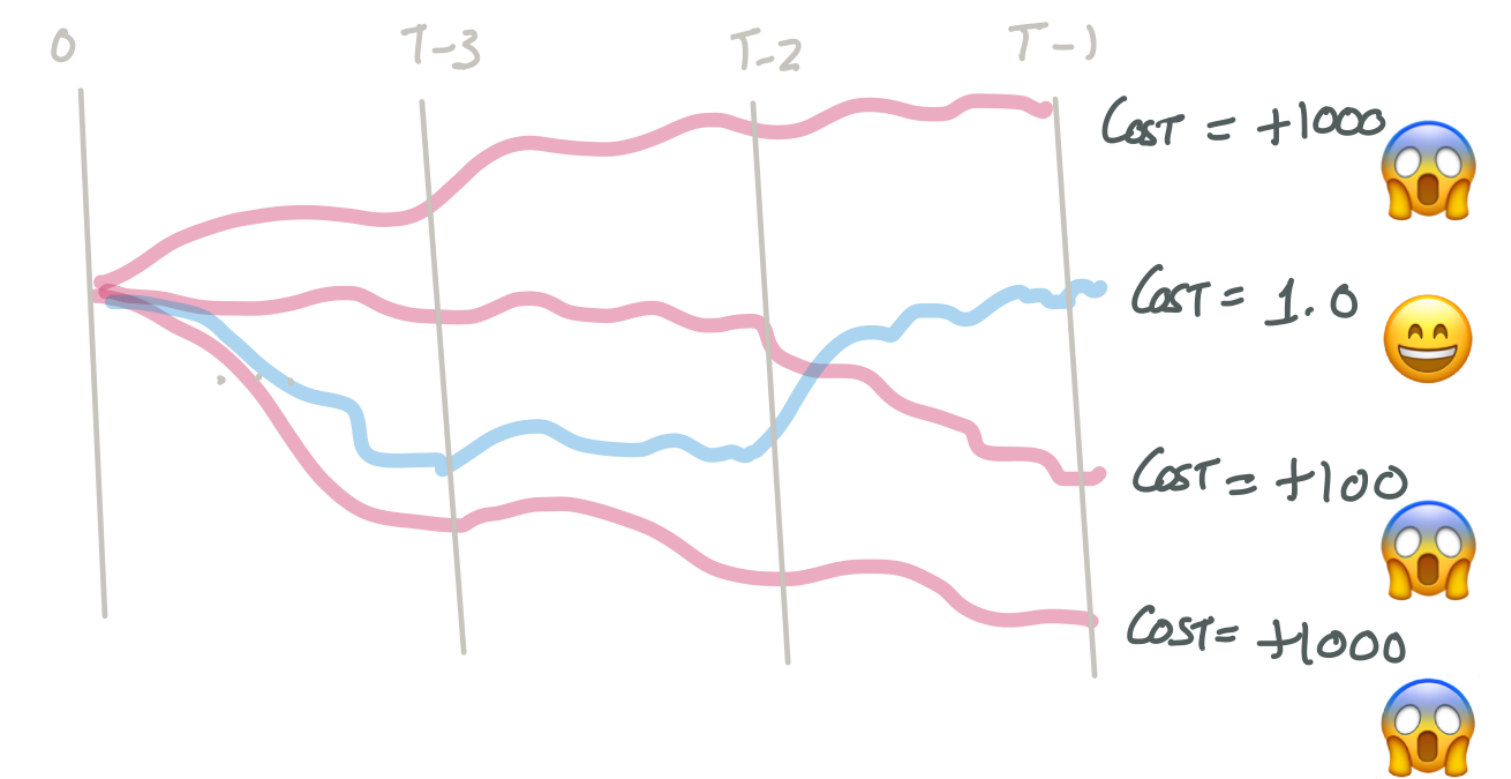
# Think-Pair-Share

Think (30 sec): What are some of the techniques we can use to solve the policy optimization?

Pair: Find a partner

Share (45 sec): Partners exchange ideas

$$\min_{\theta} J(\theta) = \sum_{t=0}^{T-1} \mathbb{E}_{\pi_{\theta}} c(s_t, a_t)$$



# Option 1: Gradient Descent

$$\theta^+ = \theta - \eta \nabla_{\theta} J(\theta)$$

What could go wrong?

1. Is  $J(\theta)$  differentiable?
2. Is  $J(\theta)$  noisy?



# Option 2: Zeroth Order Methods

- Nelder Mead
- Cross Entropy
- Simulated Annealing
- Genetic Algorithm
- Response Surface Methods
- Coordinate Descent





A photograph of a computer monitor sitting on a sandy beach. The monitor's screen displays the words "Cross Entropy" in green text. A white cable is connected to the back of the monitor and extends across the sand to a connector. The background shows the ocean and a hazy sky.

Cross  
Entropy

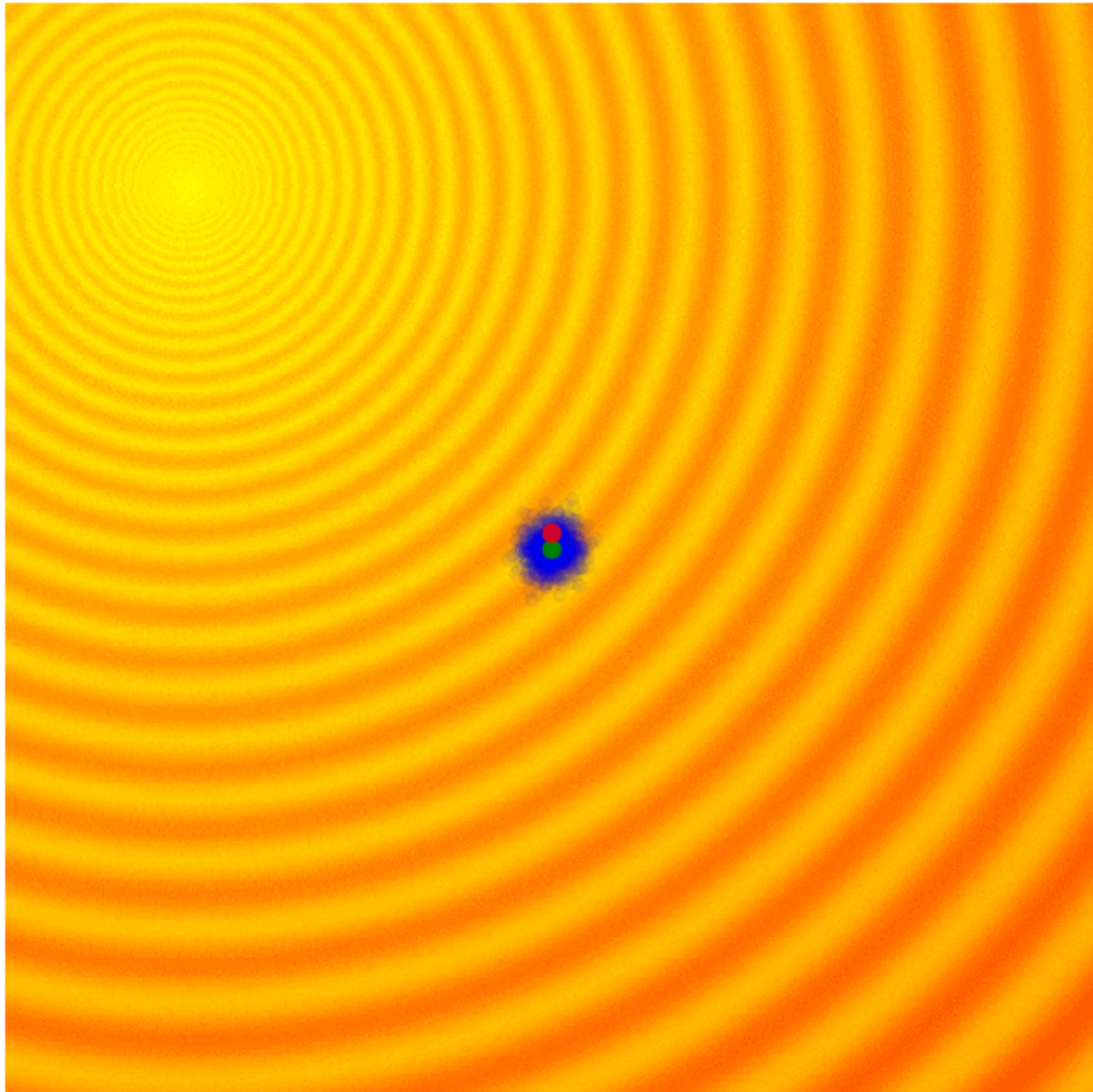
If you were ever  
stranded on an  
island ...



Let's build  
some intuition!







Credit: [https://  
blog.otoro.net/  
2017/10/29/visual-  
evolution-strategies/](https://blog.otoro.net/2017/10/29/visual-evolution-strategies/)

# The Cross Entropy Algorithm



$I_{\text{NIT}}$

$D_{\theta}$

# The Cross Entropy Algorithm



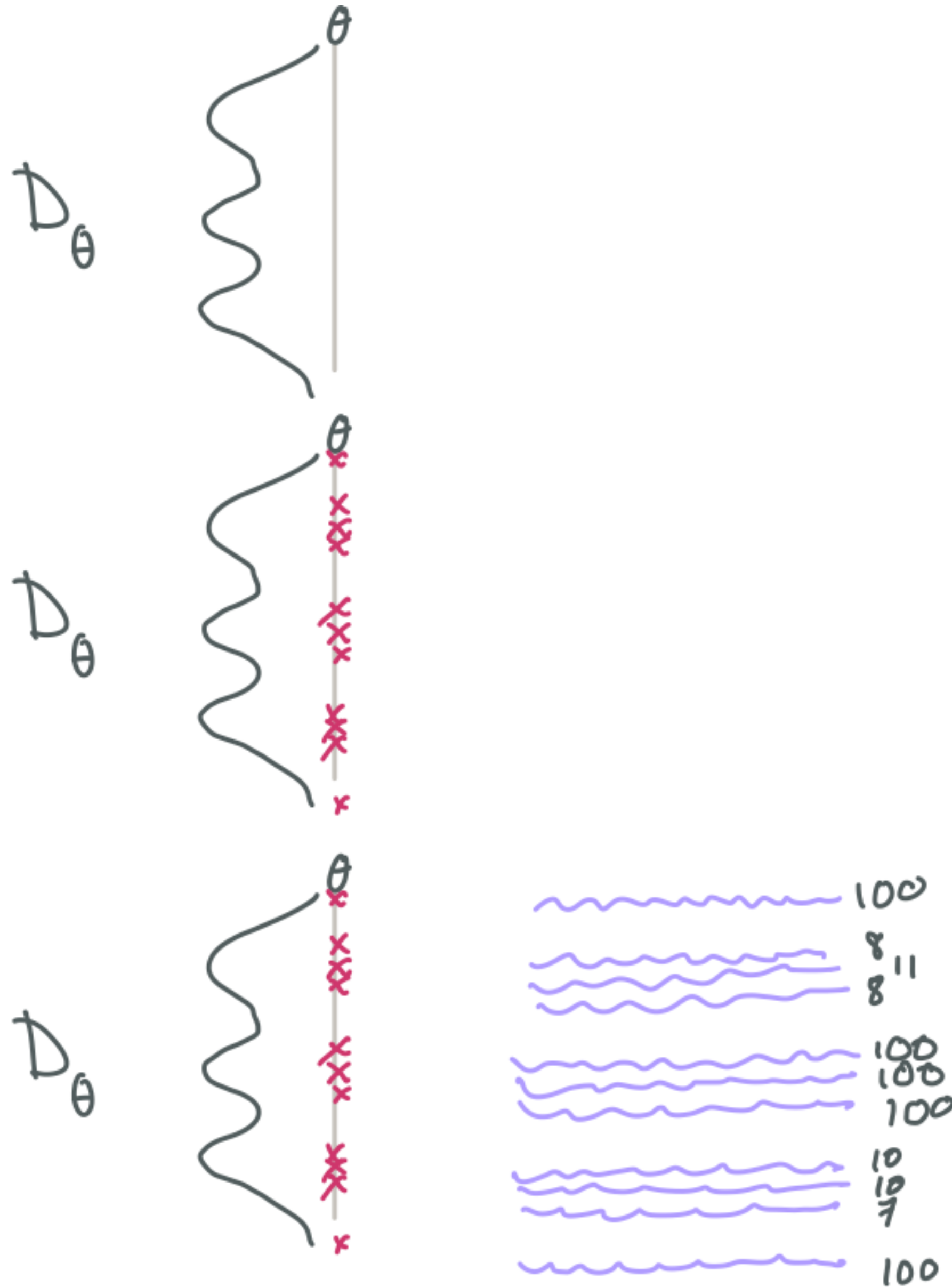
$I_{\text{INIT}}$

$D_\theta$

SAMPLE  $k$  TIMES  
to get  $\{\theta_i\}_{i=1}^k$



# The Cross Entropy Algorithm



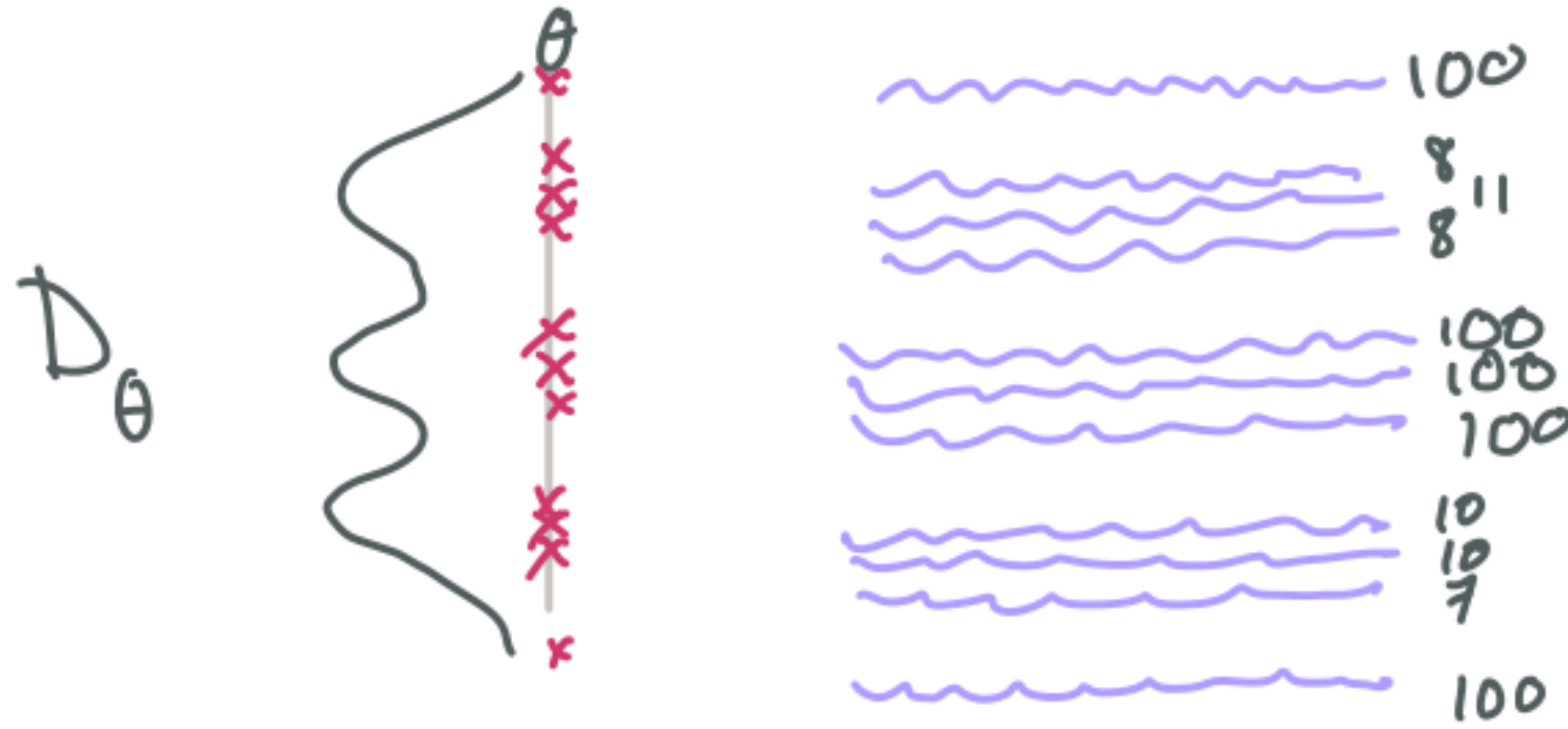
$I_{INIT}$       $D_\theta$

SAMPLE  $k$  TIMES  
to get  $\{\theta_i\}_{i=1}^k$

EVALUATE EACH  $\theta_i$

- EXECUTE POLICY MULTIPLE TIMES

# The Cross Entropy Algorithm

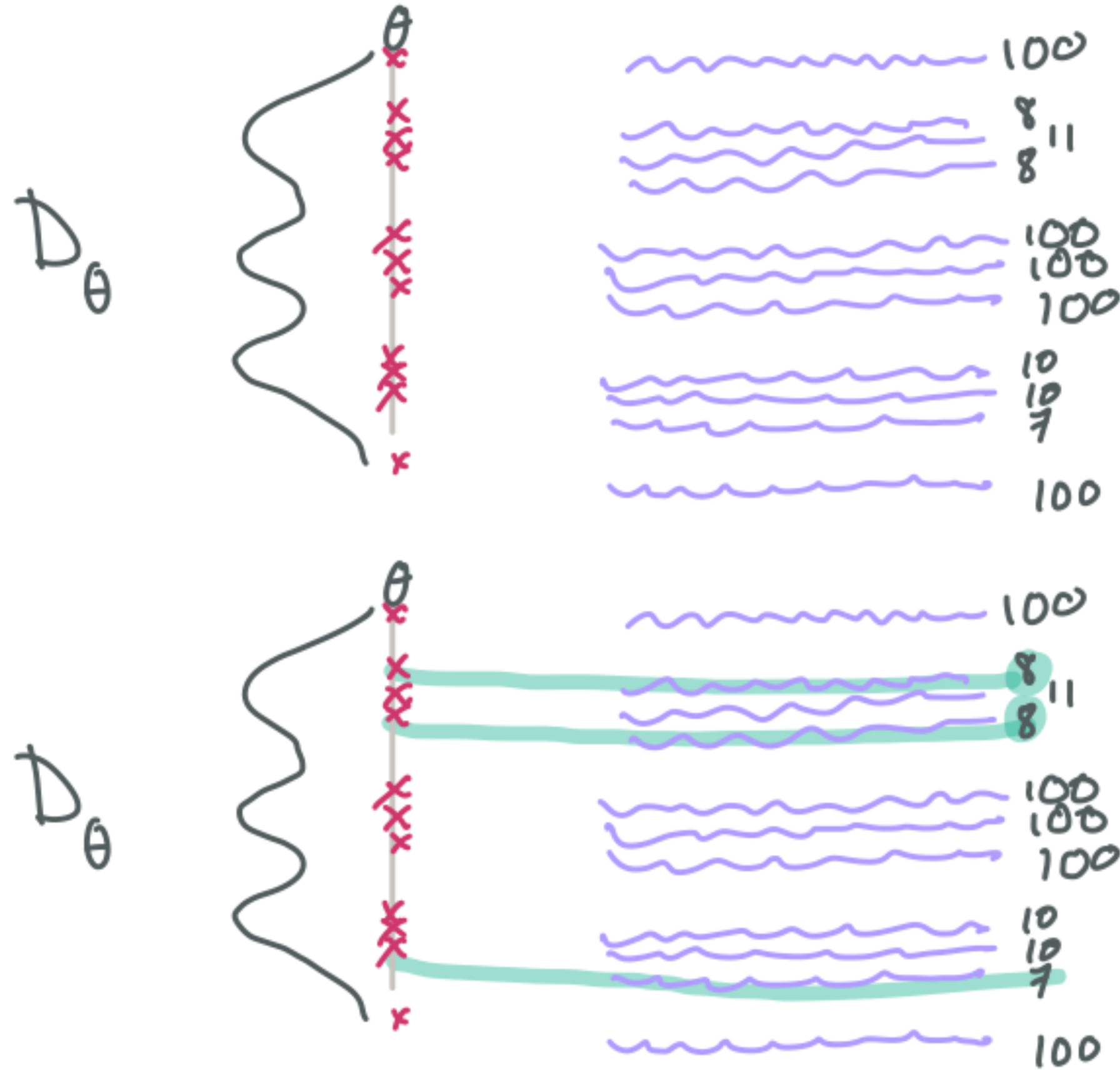


EVALUATE EACH  $\theta_i$

- EXECUTE POLICY MULTIPLE TIMES



# The Cross Entropy Algorithm

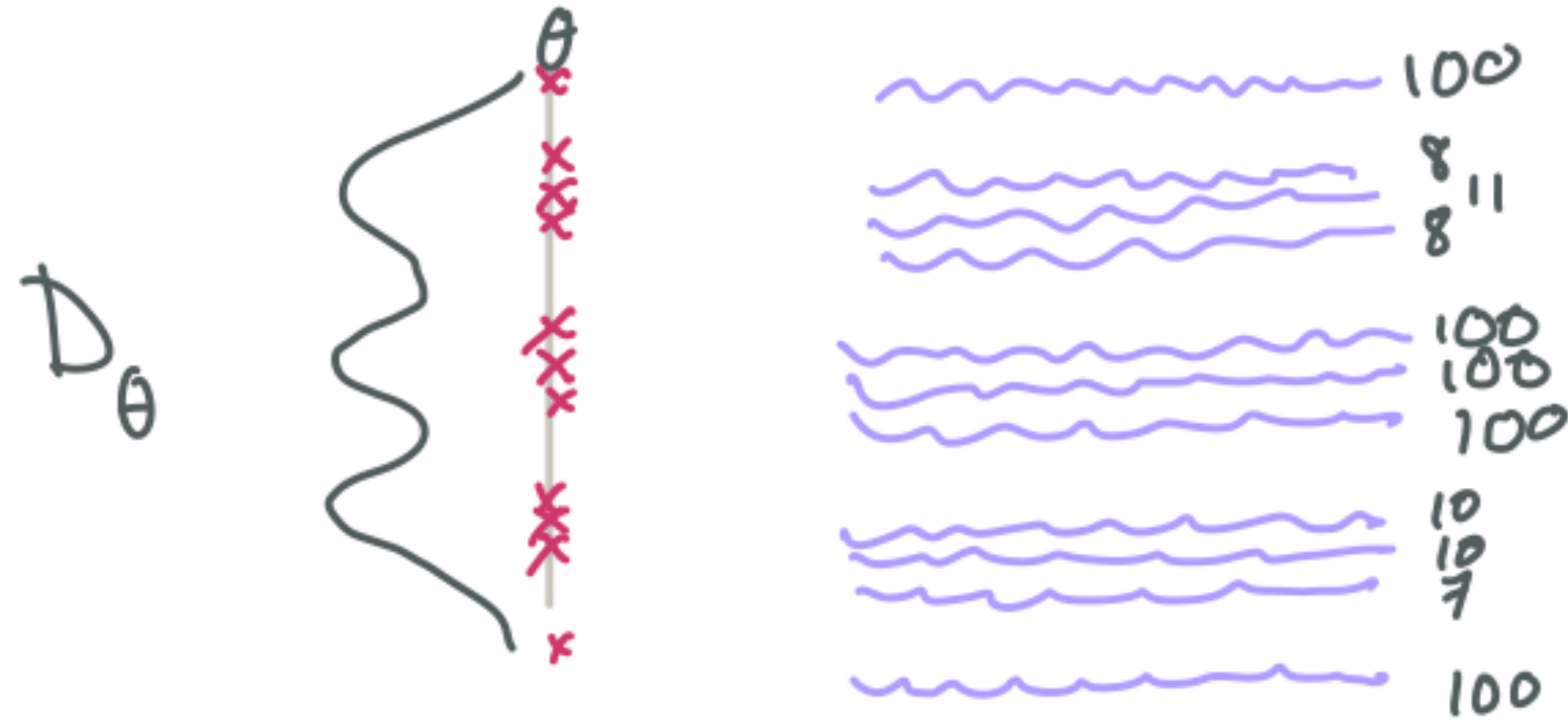


EVALUATE EACH  $\theta_i$

- EXECUTE POLICY MULTIPLE TIMES

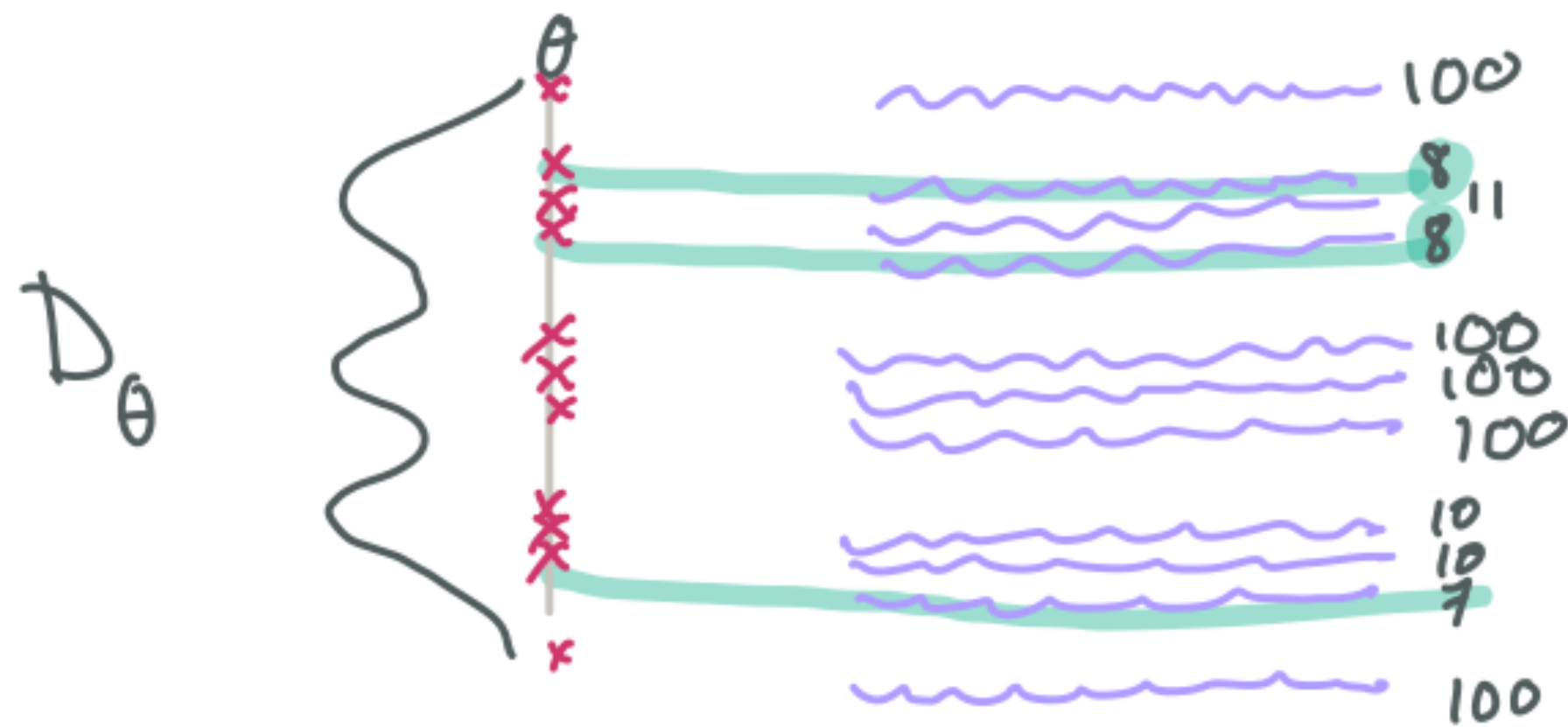
FIND TOP 'E' ELITES  
(e.g. 25%)

# The Cross Entropy Algorithm

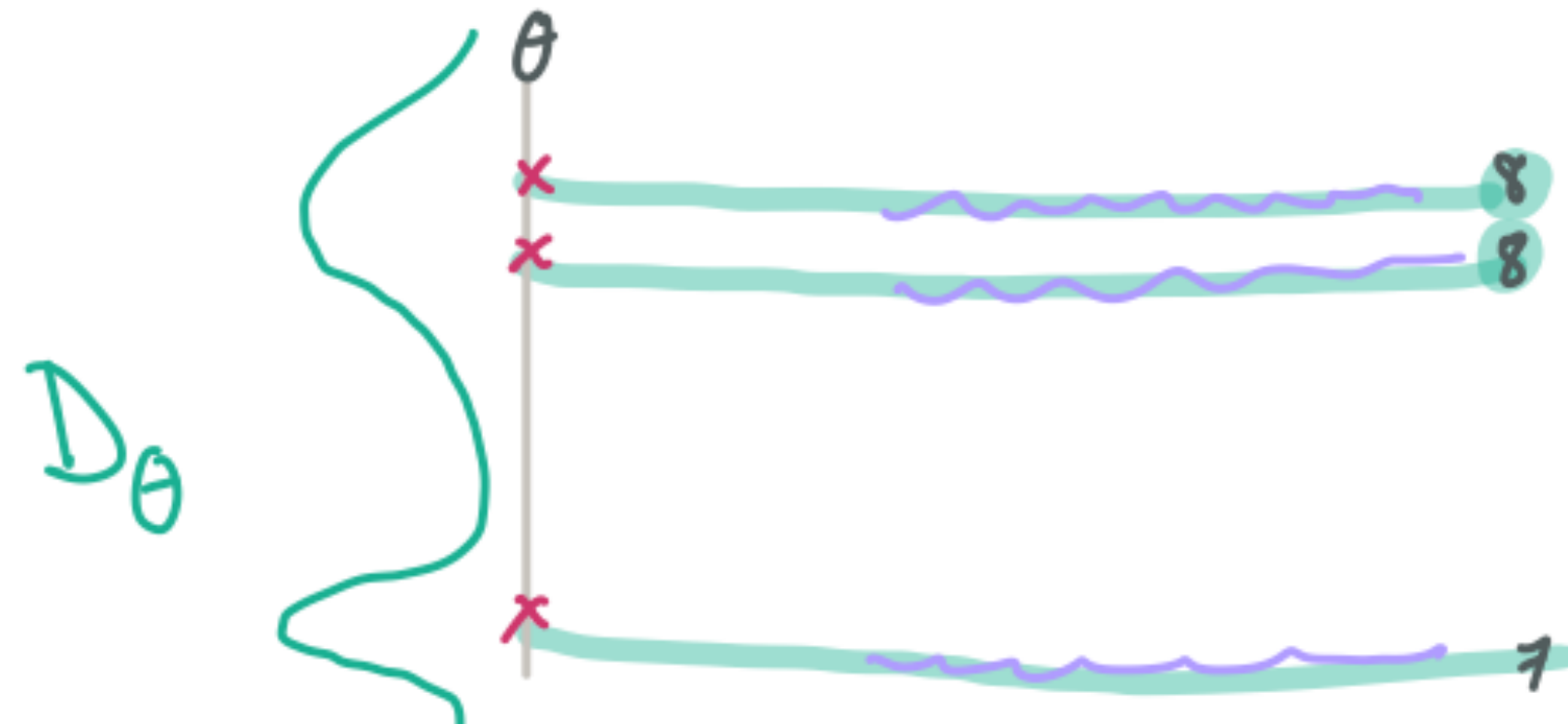


EVALUATE EACH  $\theta_i$

- EXECUTE POLICY MULTIPLE TIMES



FIND TOP 'E' ELITES  
(e.g. 25%)



FIT A NEW DISTRIBUTION

$D_\theta$



# Cross Entropy for Gaussian

Gaussian Distribution

$$D_{\theta} := \mathcal{N}(\mu, \Sigma)$$

Mean

$$\mu^t = \frac{1}{e} \sum_{i=1}^e \theta_i$$

Variance

$$\Sigma^t = \frac{1}{e} \sum_{i=1}^e (\theta_i - \mu^t)^2$$

# Does it work?

## Learning Tetris Using the Noisy Cross-Entropy Method

**István Szita**

*szityu@eotvos.elte.hu*

**András Lőrincz**

*andras.lorincz@elte.hu*

*Department of Information Systems, Eötvös Loránd University, Budapest, Hungary  
H-1117*

**The cross-entropy method is an efficient and general optimization algorithm. However, its applicability in reinforcement learning (RL) seems to be limited because it often converges to suboptimal policies. We apply noise for preventing early convergence of the cross-entropy method, using Tetris, a computer game, for demonstration. **The resulting policy outperforms previous RL algorithms by almost two orders of magnitude.****



# Does it work?

	ALGORITHM	GRID SIZE	LINES CLEARED	FEATURE SET USED
TSITSIKLIS & VAN ROY (1996)	APPROXIMATE VALUE ITERATION	16 × 10	30	HOLES AND PILE HEIGHT
BERTSEKAS & TSITSIKLIS (1996)	λ - PI	19 × 10	2,800	BERTSEKAS
LAGOUDAKIS ET AL. (2002)	LEAST-SQUARES PI	20 × 10	≈ 2,000	LAGOUDAKIS
KAKADE (2002)	NATURAL POLICY GRADIENT	20 × 10	≈ 5,000	BERTSEKAS
DELLACHERIE [REPORTED BY FAHEY (2003)]	HAND TUNED	20 × 10	660,000	DELLACHERIE
RAMON & DRIESSENS (2004)	RELATIONAL RL	20 × 10	≈ 50	
BÖHM ET AL. (2005)	GENETIC ALGORITHM	20 × 10	480,000,000 (TWO PIECE)	BÖHM
FARIAS & VAN ROY (2006)	LINEAR PROGRAMMING	20 × 10	4,274	BERTSEKAS
SZITA & LÖRINCZ (2006)	CROSS ENTROPY	20 × 10	348,895	DELLACHERIE
ROMDHANE & LAMONTAGNE (2008)	CASE-BASED REASONING AND RL	20 × 10	≈ 50	
BOUMAZA (2009)	CMA-ES	20 × 10	35,000,000	BCTS
THIERY & SCHERRER (2009A;B)	CROSS ENTROPY	20 × 10	35,000,000	DT
GABILLON ET AL. (2013)	CLASSIFICATION-BASED POLICY ITERATION	20 × 10	51,000,000	DT FOR POLICY DT + RBF FOR VALUE



# Practical Issues and Fixes





# Problem 1: What happens to the variance?

$$\Sigma^t = \frac{1}{e} \sum_{i=1}^e (\theta_i - \mu^t)^2$$

Collapses too quickly!

Simple fix: Add a bit of noise to the variance

$$\Sigma^t = \frac{1}{e} \sum_{i=1}^e (\theta_i - \mu^t)^2 + \Sigma_{noise}$$

Problem 2: What if we have a bad batch of samples?

$$\mu^t = \frac{1}{e} \sum_{i=1}^e \theta_i$$

The elites can be bad, and the mean can slingshot into a bad value

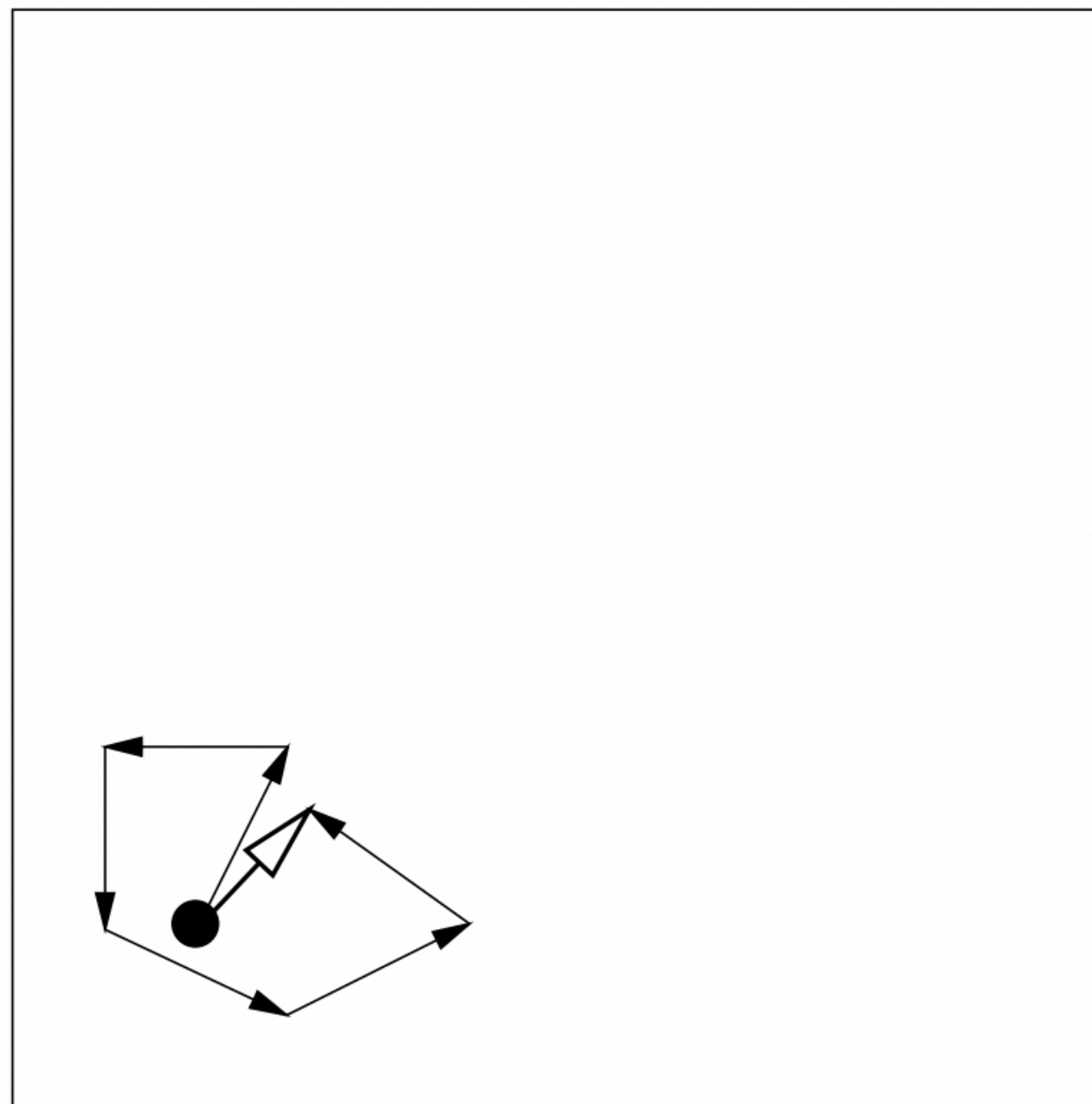
Simple fix: Slowly update mean

$$\mu^t = \mu^{t-1} + \eta \frac{1}{e} \sum_{i=1}^e \theta_i$$

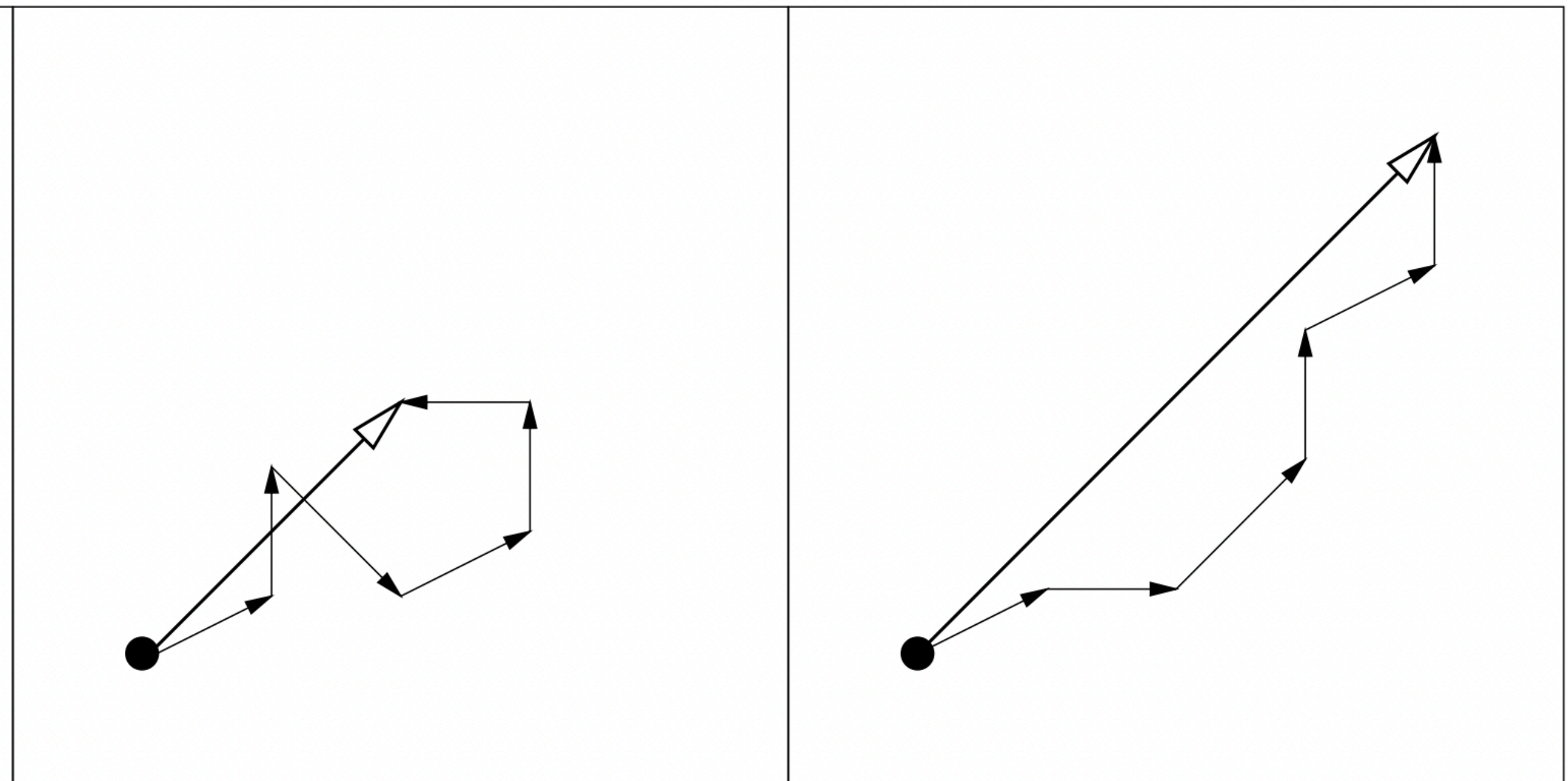


# Problem 3: What if we never converge and do random walks?

Single-steps cancel out  
Use small  $\Sigma$



Progress correlated  
Use large  $\Sigma$



A very fancy version of Cross Entropy: CMA-ES

Tetris is cute...  
But what about *real*  
robots?

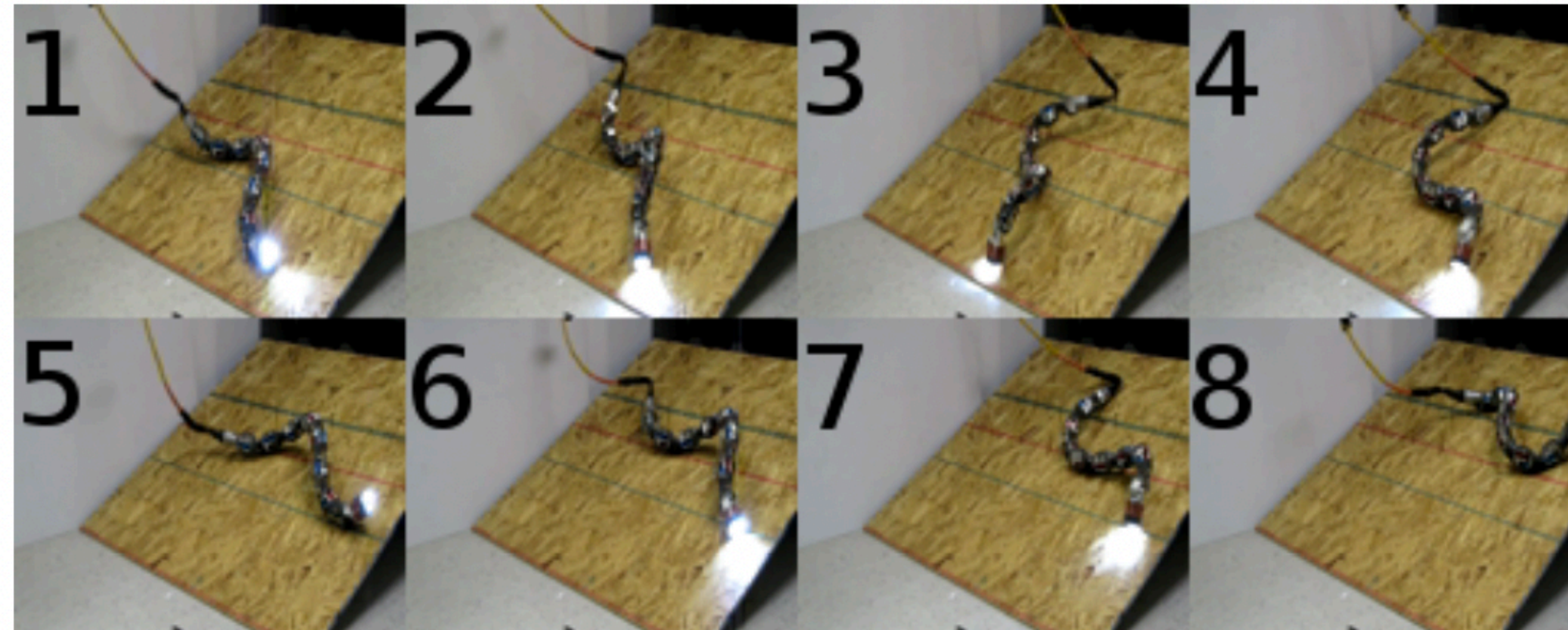




# Cross Entropy for Snake Robot Gaits

Using Response Surfaces and Expected Improvement to Optimize Snake Robot Gait Parameters

Matthew Tesch, Jeff Schneider, and Howie Choset

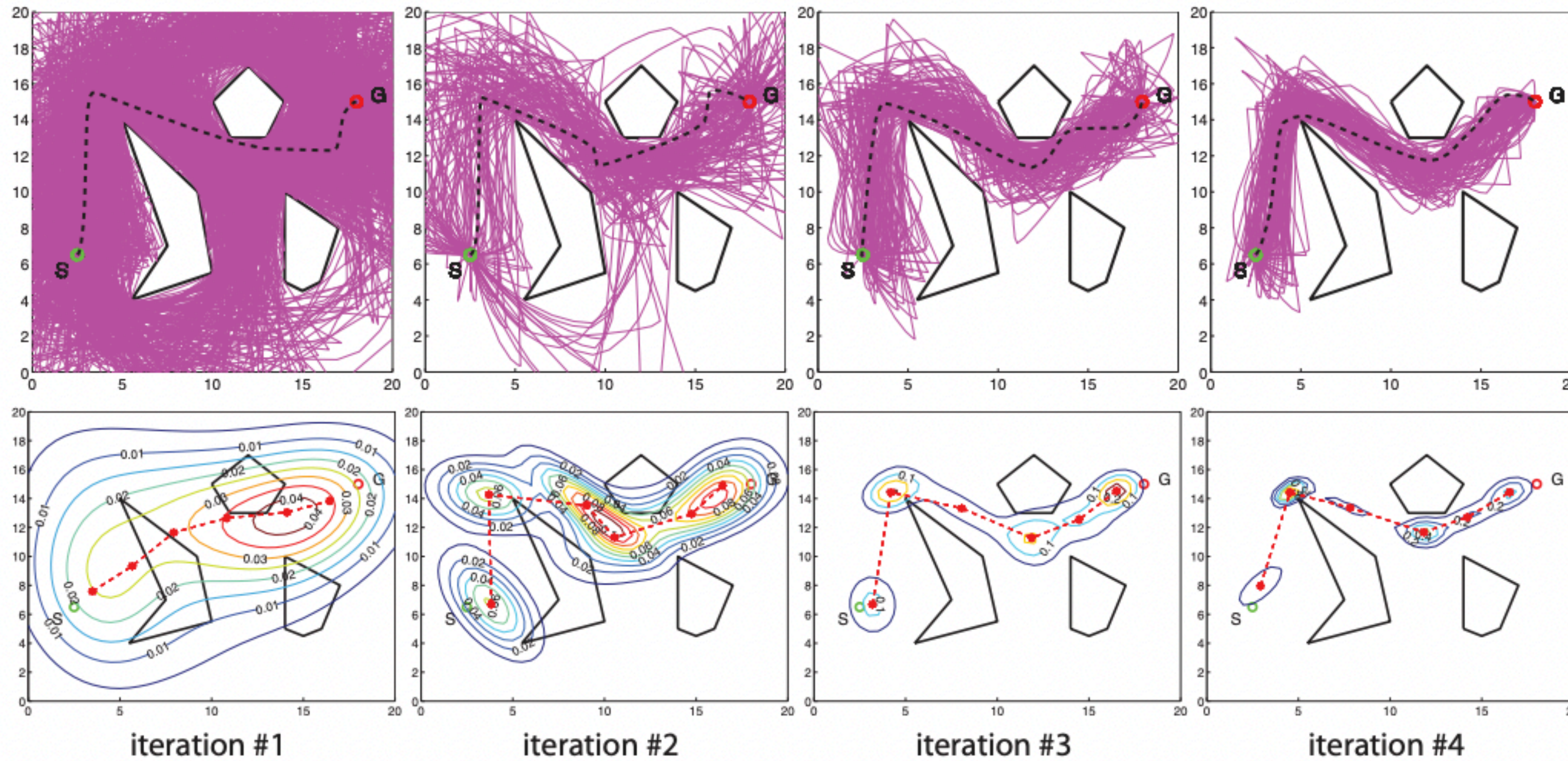


Uses a Gaussian Process to fit a distribution

Prove it can find the optimal gait with *minimal samples*



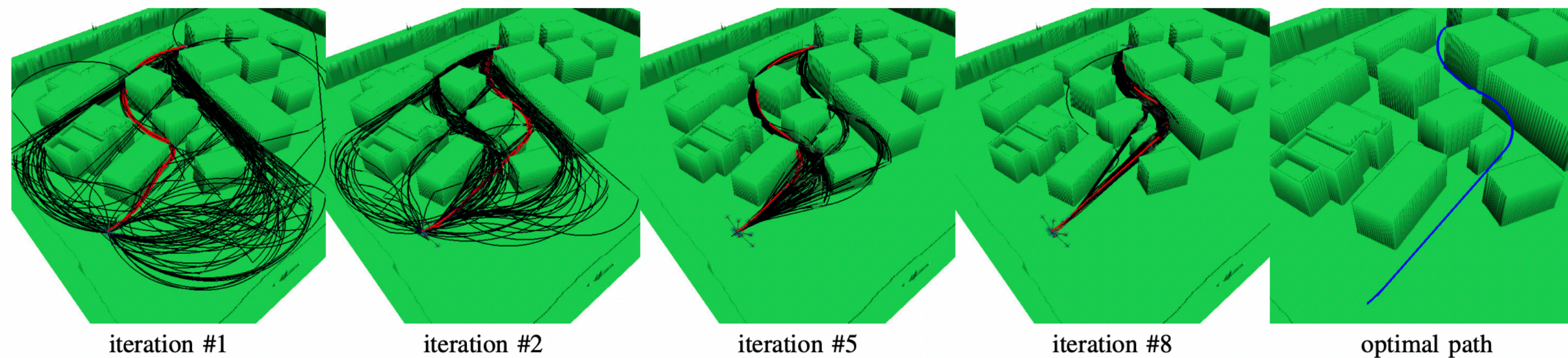
# Cross Entropy Search for Motion Planning



Cross-Entropy Randomized Motion Planning

Marin Kobilarov

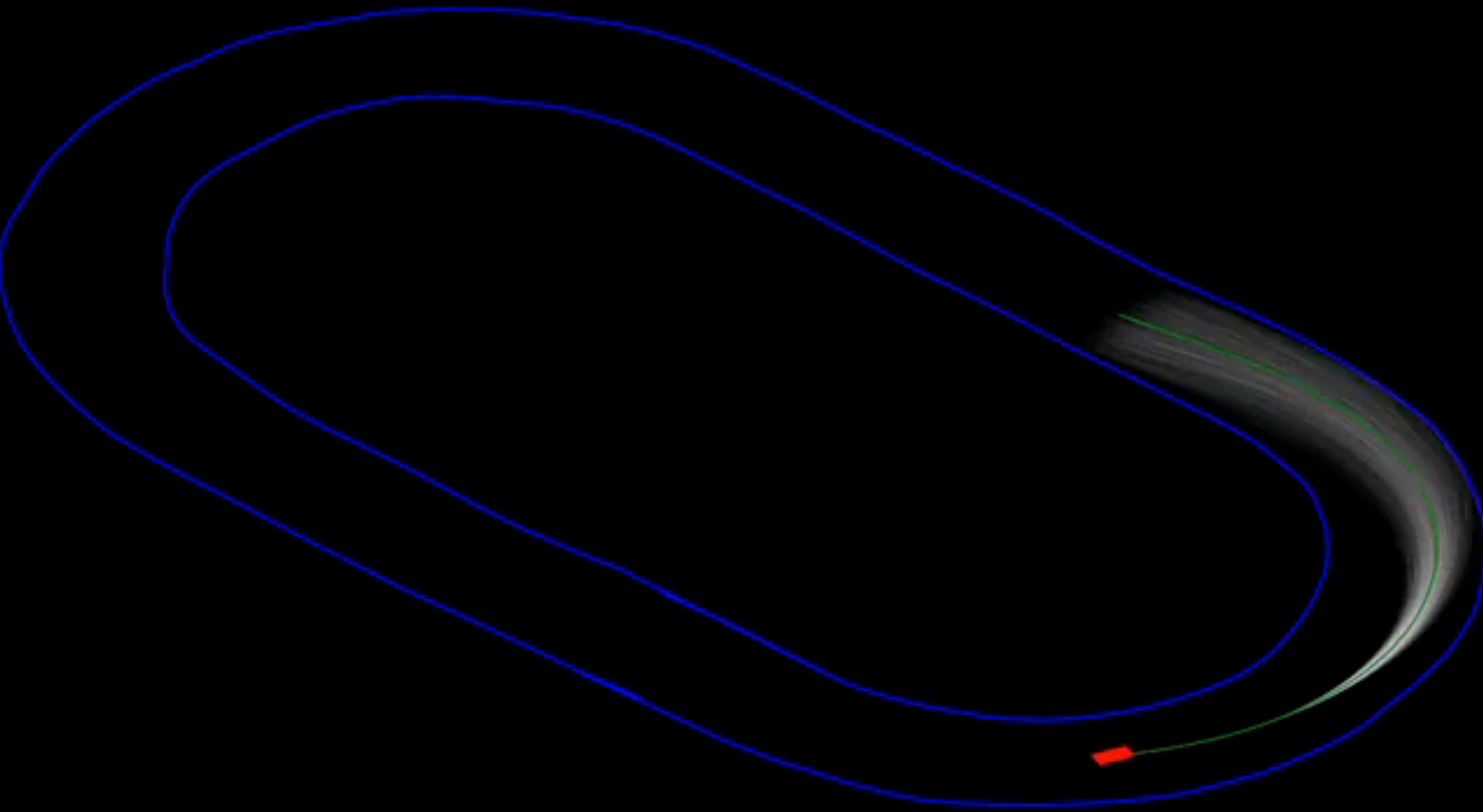
Distribution over  
control trajectories





2560, 2.5 second trajectories sampled  
with cost-weighted average @ 60 Hz

# Cross Entropy for Control



Georgia Tech Auto Rally (Byron Boots lab)



DATE:  
No  
Pr  
P

Form No. 1  
THIS CASE ORIGINATED AT

REPORT MADE AT

TITLE

Mr. ARNOLD GIBBY

SYNOPSIS OF FACTS:

Subject employment as a worker  
chief [redacted] for the [redacted]

DETAILS:

BACKGROUND

Birth [redacted]

Have we redacted too much?

SUBJECT: [redacted]

1. [redacted] is being sent to [redacted] work of [redacted] and biological study at [redacted] in the areas of [redacted] until [redacted]. The project [redacted]

2. This [redacted] into the [redacted] areas, [redacted] effects of [redacted] load continue to be [redacted] and [redacted] will [redacted]

3. The [redacted] The [redacted] estimated to in the [redacted] required, under [redacted]

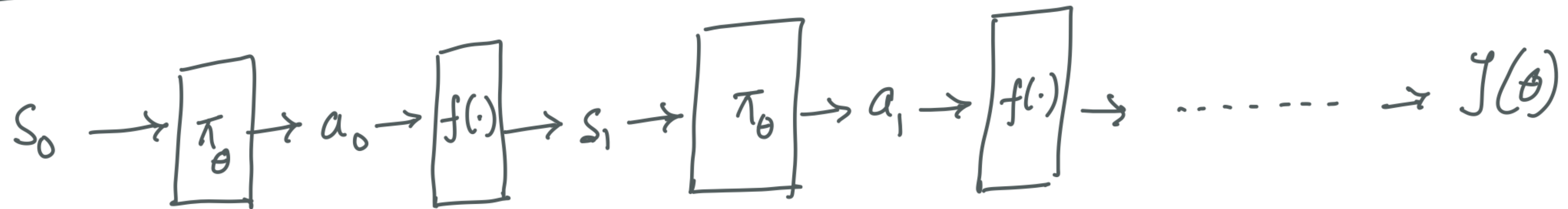


# Black-box vs White-box vs Gray-box

BLACK BOX



WHITE BOX

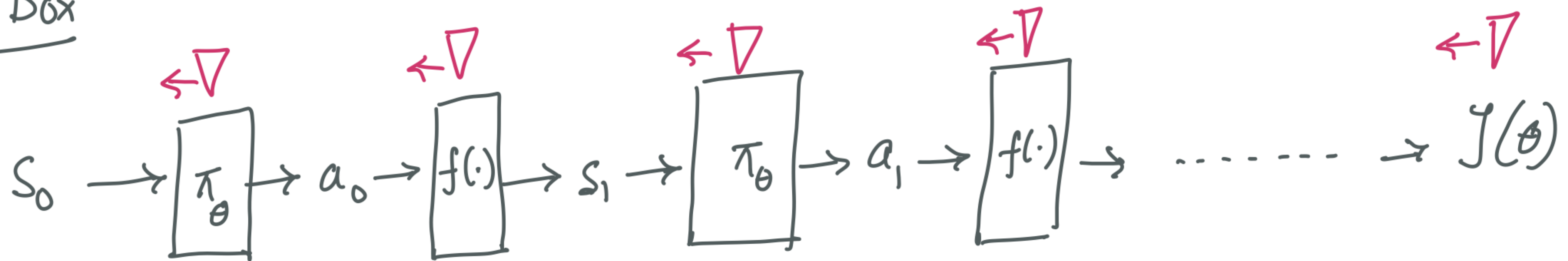


# Black-box vs White-box vs Gray-box

Black Box



White Box





How can we take  
gradients if we don't  
know the dynamics?



# The Likelihood Ratio Trick!





# REINFORCE

---

**Algorithm 20:** The REINFORCE algorithm.

---

Start with an arbitrary initial policy  $\pi_\theta$

**while** *not converged* **do**

Run simulator with  $\pi_\theta$  to collect  $\{\zeta^{(i)}\}_{i=1}^N$

Compute estimated gradient

$$\tilde{\nabla}_\theta J = \frac{1}{N} \sum_{i=1}^N \left[ \left( \sum_{t=0}^{T-1} \nabla_\theta \log \pi_\theta \left( a_t^{(i)} | s_t^{(i)} \right) \right) R(\zeta^{(i)}) \right]$$

Update parameters  $\theta \leftarrow \theta + \alpha \tilde{\nabla}_\theta J$

**return**  $\pi_\theta$

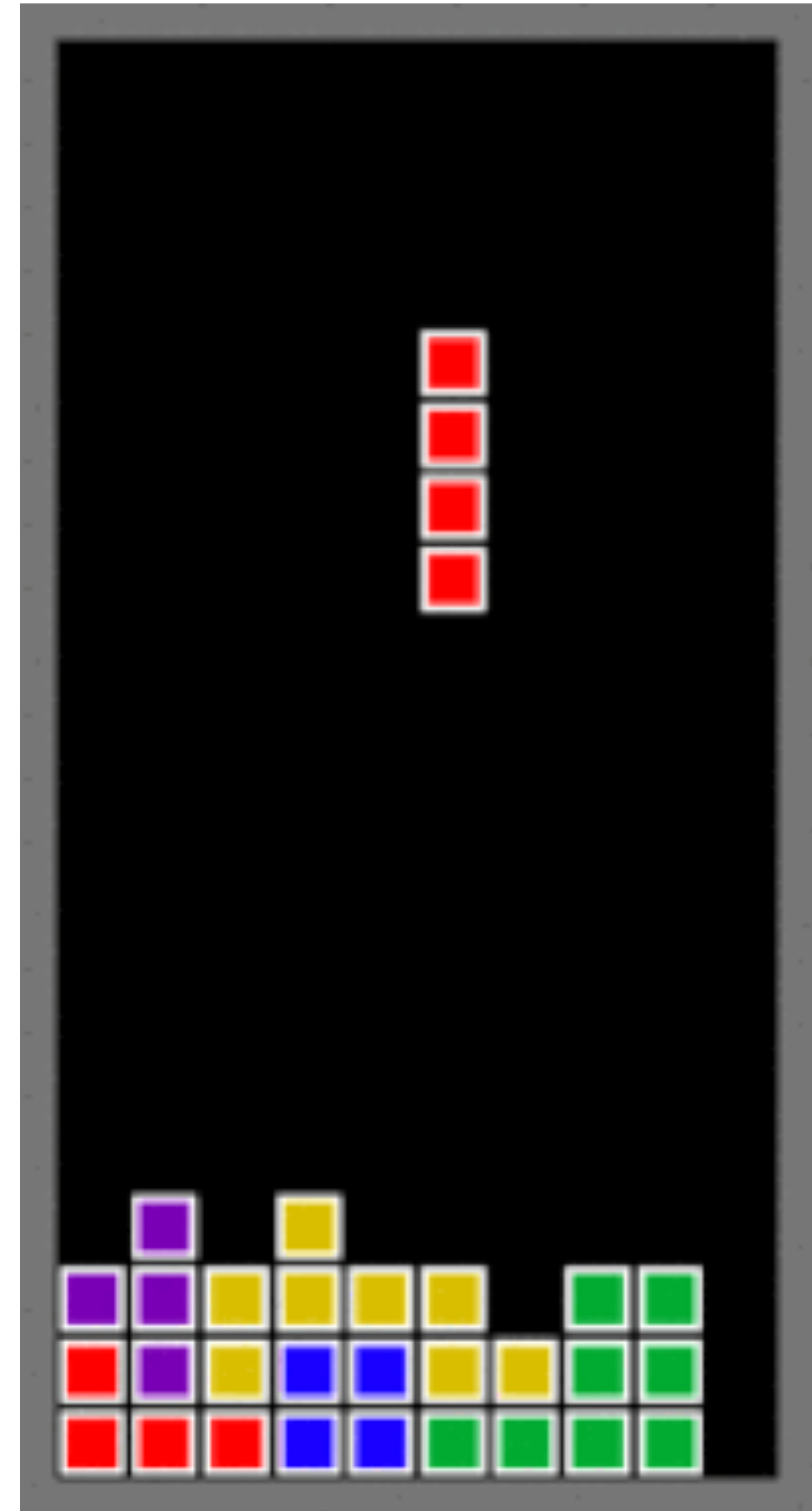
---

# Tetris Policy

$$\pi_{\theta}(a|s) = \frac{\exp(\theta^{\top} f(s, a))}{\sum_{a'} \exp(\theta^{\top} f(s, a'))}$$

$f_1(s, a) = \#$  number of holes

$f_2(s, a) = \#$  max height



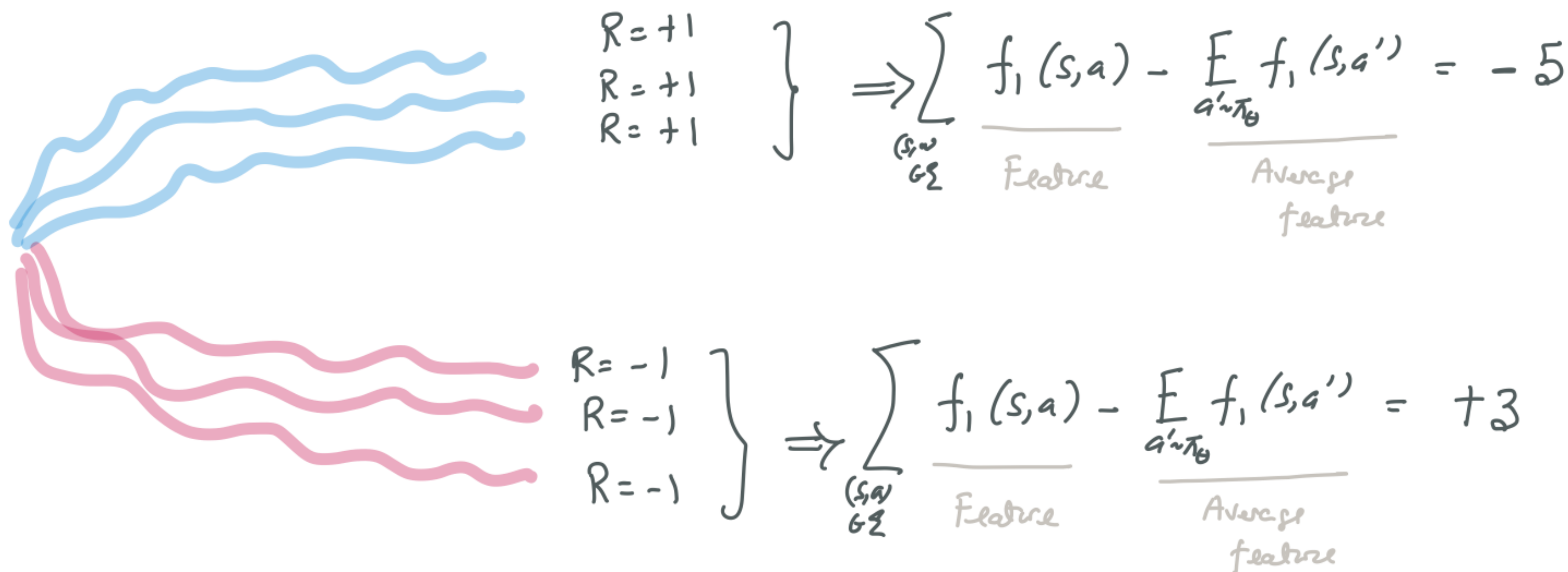


# Chugging through the gradient ..

$$\begin{aligned}\nabla_{\theta} \log \pi_{\theta}(a|s) &= \nabla_{\theta} \left[ \theta^{\top} f(s, a) - \log \sum_{a'} \exp \left( \theta^{\top} f(s, a') \right) \right] \\ &= f(s, a) - \frac{\sum_{a'} f(s, a') \exp \left( \theta^{\top} f(s, a') \right)}{\sum_{a'} \exp \left( \theta^{\top} f(s, a') \right)} \\ &= f(s, a) - \sum_{a'} f(s, a') \pi_{\theta} (a' | s) \\ &= f(s, a) - E_{\pi_{\theta}(a'|s)} [f(s, a')]\end{aligned}$$

# Understanding the REINFORCE update

LET  $f_1(s,a) = \# \text{ holes.}$



$$\begin{aligned} \Theta_1 &= \Theta_0 + \sum_{\substack{(s,a) \\ \in \mathcal{Z}}} \left( \nabla_{\Theta} \log \pi_{\Theta}(a|s) \right) R(\mathcal{Z}) = \Theta_0 + \alpha \left( -5 \times (+1) + 3 \times (-1) \right) \\ &= \Theta_0 - \alpha 8 \quad (\text{Bump down this feature}) \end{aligned}$$

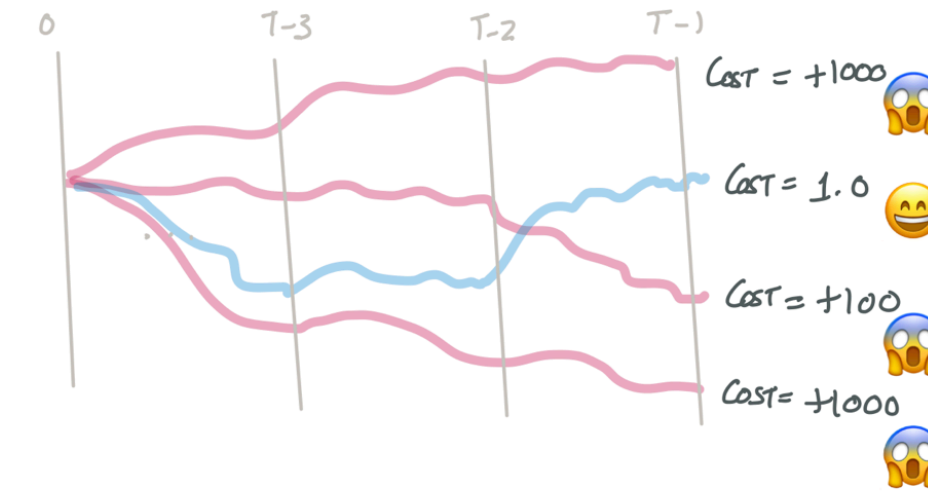


# tl;dr

## The Goal of Policy Optimization

$$\pi_{\theta}(s) = \arg \min_a \theta^T f(s, a)$$

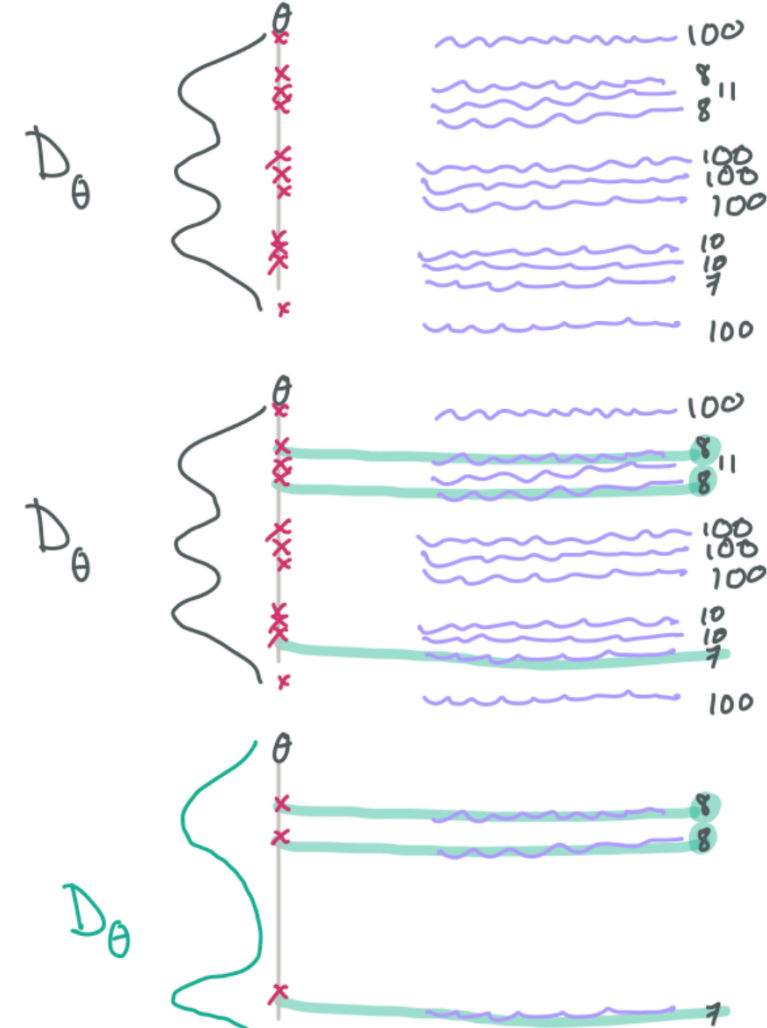
$$\min_{\theta} J(\theta) = \sum_{t=0}^{T-1} \mathbb{E}_{\pi_{\theta}} c(s_t, a_t)$$



Black-box

Gray-box

## The Cross Entropy Algorithm



EVALUATE EACH  $\theta_i$

- EXECUTE POLICY MULTIPLE TIMES

FIND TOP 'B' PERCENTS (e.g. 25%)

FIT A NEW DISTRIBUTION

$D_{\theta}$

## REINFORCE

**Algorithm 20:** The REINFORCE algorithm.

Start with an arbitrary initial policy  $\pi_{\theta}$

**while not converged do**

Run simulator with  $\pi_{\theta}$  to collect  $\{\zeta^{(i)}\}_{i=1}^N$

Compute estimated gradient

$$\tilde{\nabla}_{\theta} J = \frac{1}{N} \sum_{i=1}^N \left[ \left( \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi_{\theta} (a_t^{(i)} | s_t^{(i)}) \right) R(\zeta^{(i)}) \right]$$

Update parameters  $\theta \leftarrow \theta + \alpha \tilde{\nabla}_{\theta} J$

**return**  $\pi_{\theta}$