

Deep Reinforcement Learning and Control

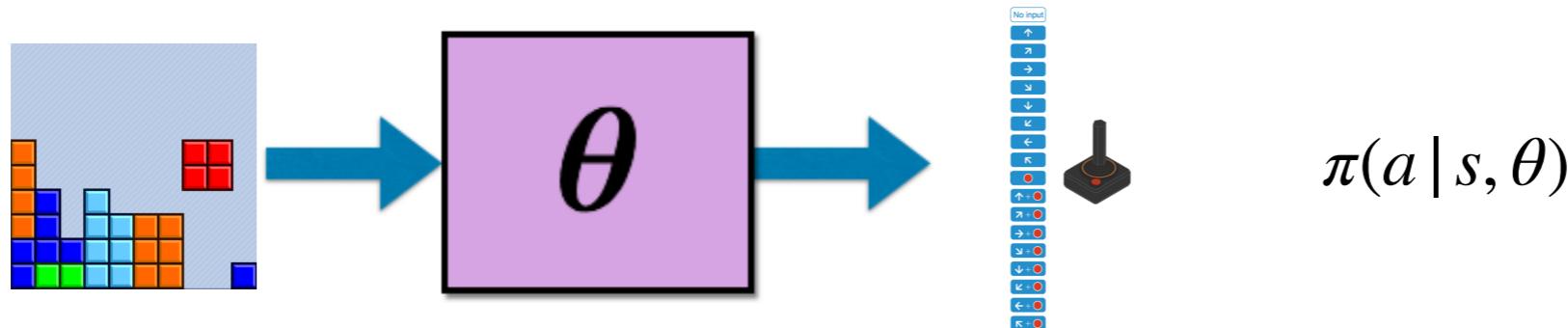
# Evolutionary Methods for Policy Search

Fall 2025, CMU 10-703

Instructors: Katerina Fragkiadaki & Aviral Kumar



# Reinforcement Learning



Given an initial state distribution  $\mu_0(s_0)$ , estimate parameters  $\theta$  of a policy  $\pi_\theta$  so that, the trajectories  $\tau$  sampled from this policy have maximum returns, i.e., sum of rewards  $R(\tau)$ .

$$\max_{\theta} . \quad U(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau) | \pi_\theta, \mu_0(s_0)]$$

$\tau$  : trajectory, a sequence of state, action, rewards, a game fragment or a full game:

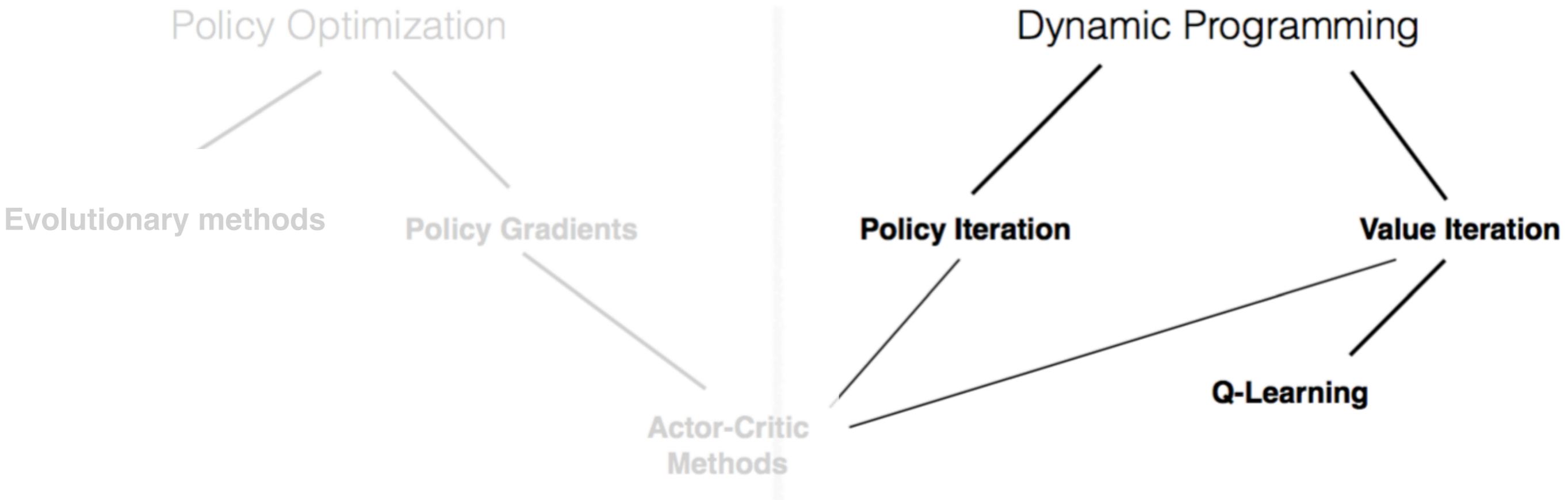
$$\tau : s_0, a_0, r_0, s_1, a_1, r_1, s_2, a_2, r_2, \dots, s_T, a_T, r_T$$

$R(\tau)$  : reward of a trajectory: (discounted) sum of the rewards of the individual state/actions

$$R(\tau) = \sum_{t=1}^T r_t$$

# Policy Optimization and RL

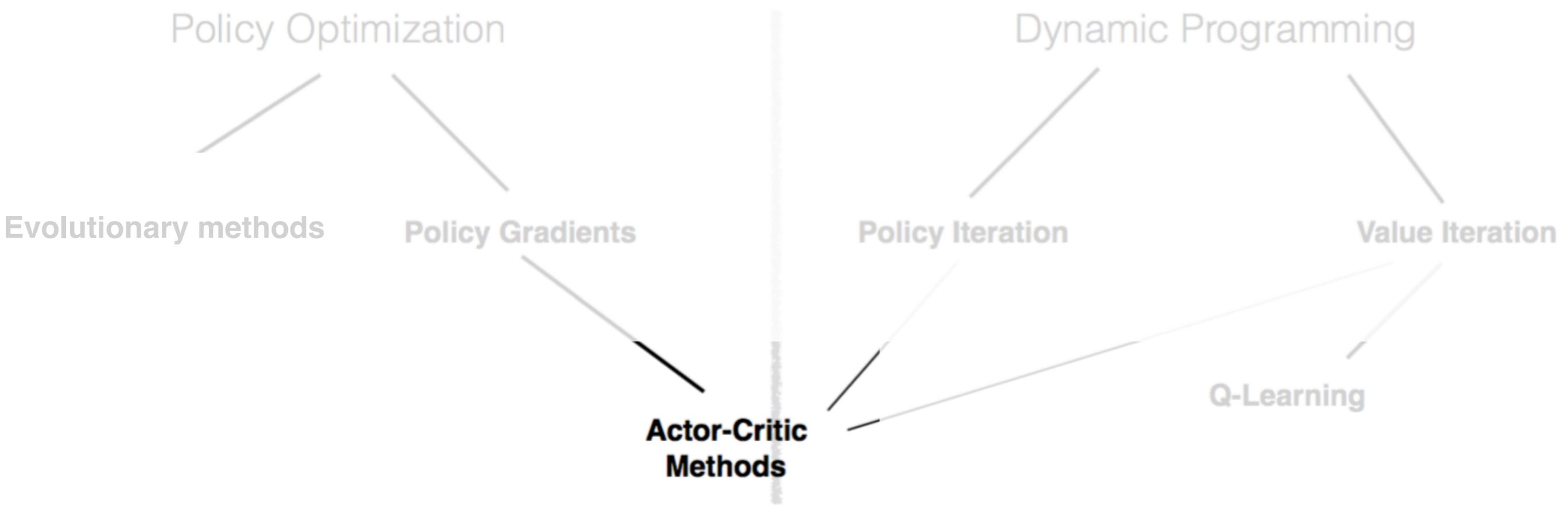
$$\max_{\theta} . \quad U(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau) | \pi_\theta, \mu_0(s_0)]$$



Dynamic programming methods search over values for states without explicitly estimating policies. They use recursive updates that capture how one state results from another.

# Black-box Policy Optimization

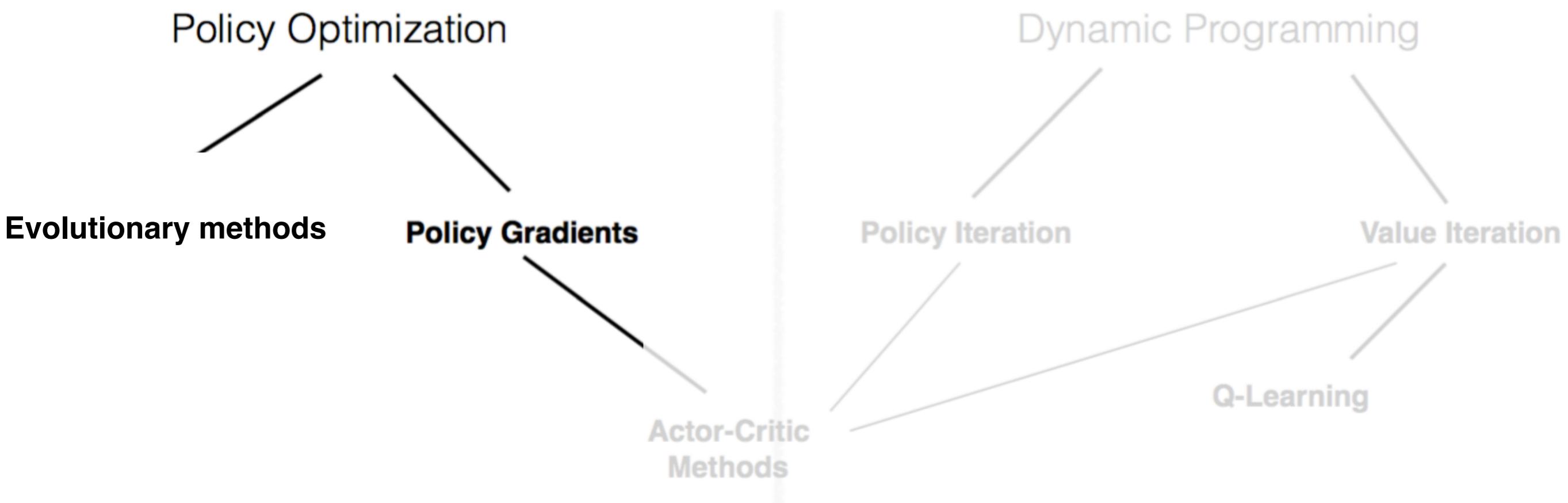
$$\max_{\theta} . \quad U(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau) | \pi_\theta, \mu_0(s_0)]$$



Actor critic method both search over policies and over state values, and use the state values, as opposed to the raw rewards to guide policy search. Estimate of value function are usually way less noisy than raw rewards

# Black-box Policy Optimization

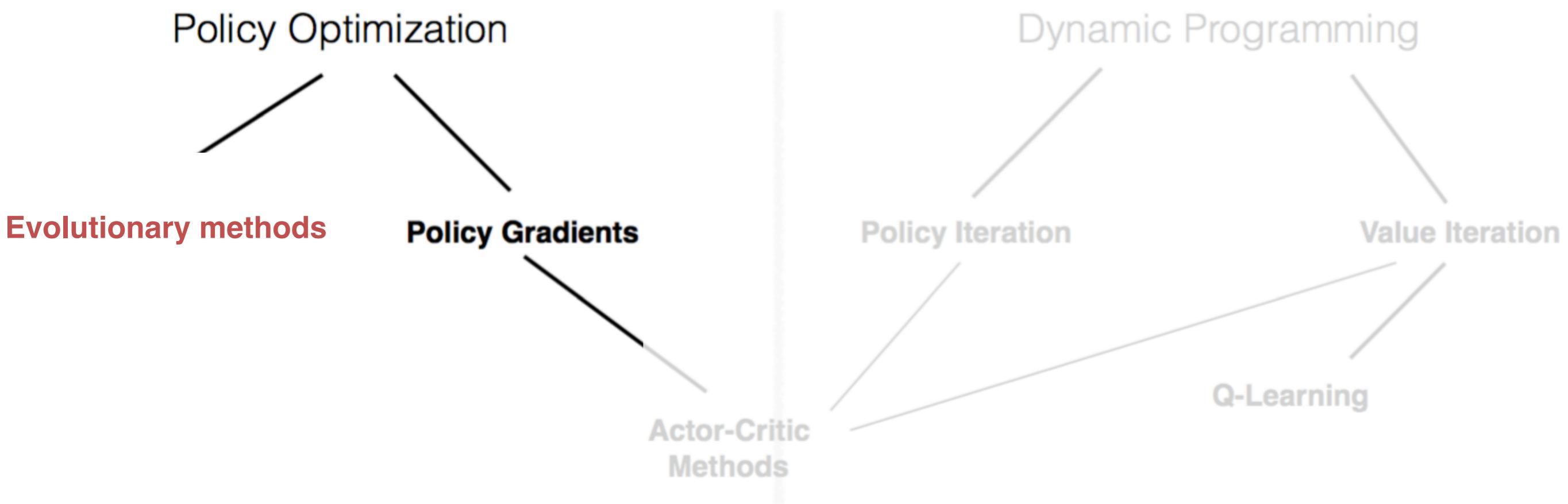
$$\max_{\theta} . \quad U(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau) | \pi_\theta, \mu_0(s_0)]$$



Policy optimization just searches over policy parameters without computing values for states: it just looks at the raw rewards obtained as we search over the parameters. It does not exploit the structure of the states

# Black-box Policy Optimization

$$\max_{\theta} . \quad U(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau) | \pi_\theta, \mu_0(s_0)]$$



Policy optimization searches over policy parameters without computing values for states. It does not exploit any recursions, or reward decompositions over timesteps.

# Black-box policy optimization

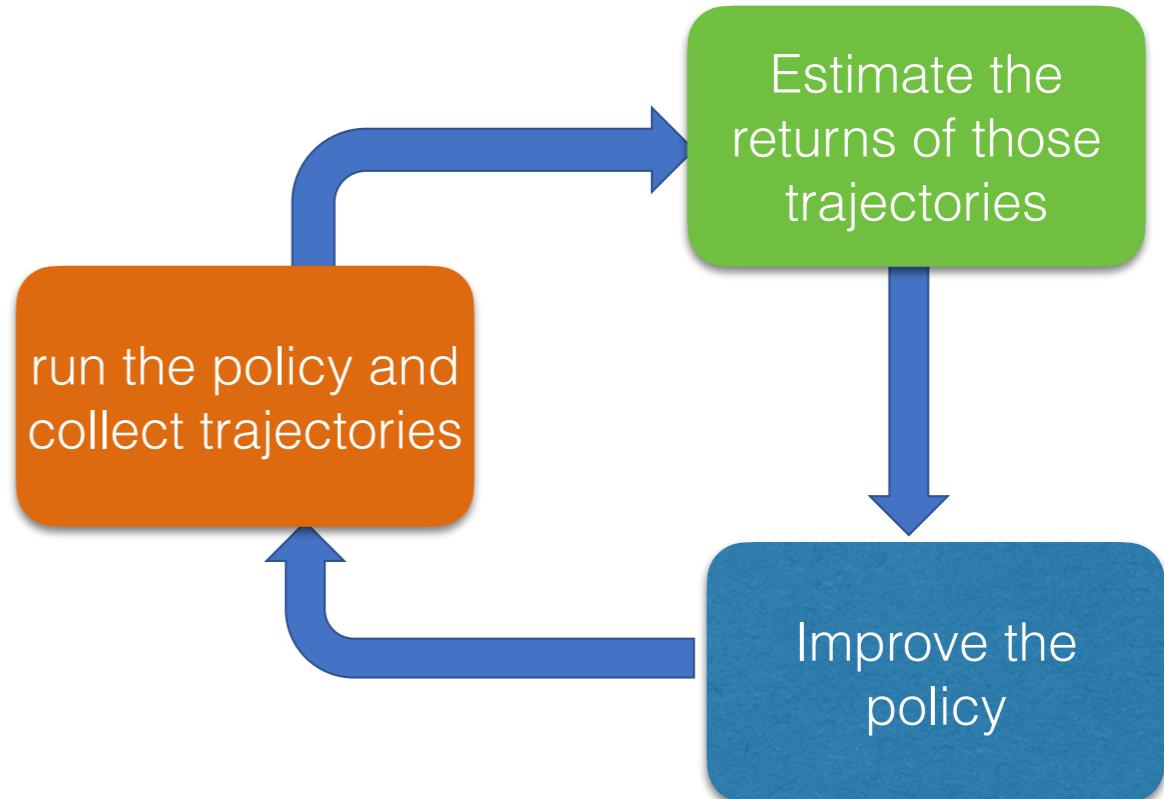
Initialize the policy parameters  $\theta$  randomly.

1. Perturb policy parameters,

2. Run the resulting policy, collect trajectories and evaluate their returns.

3. Promote the policy parameters that resulted in trajectories that gave the largest return improvement.

4. GOTO 1.



- No gradient information
- no information regarding the structure of the reward, that it is additive over states or that states are interconnected in a particular way.

# Evolutionary methods for policy search

$$\max_{\theta} . \quad U(\theta) = \mathbb{E}_{\tau \sim \pi_\theta} [R(\tau) | \pi_\theta, \mu_0(s_0)]$$

*General algorithm:*

*Initialize a population of parameter vectors (genotypes)*

1. *Make random perturbations (mutations) to each parameter vector*
2. *Evaluate the perturbed parameter vector (fitness)*
3. *Keep the perturbed vector if the result improves (selection)*
4. *GOTO 1*

*Simple and biologically plausible...*

# Gaussian Density

Perhaps the most common probability density

$$\mathcal{N}(y|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right)$$

$\sigma^2$  is the variance of the density and  $\mu$  is the mean.

# Multivariate Gaussian Density

Perhaps the most common probability density

$$\mathcal{N}(\mathbf{y} \mid \boldsymbol{\mu}, \boldsymbol{\Sigma}) = \frac{1}{\sqrt{(2\pi)^k |\boldsymbol{\Sigma}|}} \exp \left( -\frac{1}{2} (\mathbf{y} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{y} - \boldsymbol{\mu}) \right)$$

# Cross-entropy method (CEM)

- Policy parameters are sampled from a multivariate Gaussian distribution with a diagonal covariance matrix.
- We will update the mean and variances of the parameter elements towards samples that have highest fitness scores.

**Input:** parameter space  $\Theta$ , number of parameter vectors  $n$ , proportion  $\rho \leq 1$ , noise  $\eta$

**Initialize:** Set the parameter  $\mu = \bar{0}$  and  $\sigma^2 = 100I$  ( $I$  is the identity matrix)

**for**  $k = 1, 2, \dots$  **do**

    Generate a random sample of  $n$  parameter vectors  $\{\theta_i\}_{i=1}^n \sim \mathcal{N}(\mu, \sigma^2 I)$

    For each  $\theta_i$ , play  $L$  games and calculate the average number of rows removed (score) by the controller

    Select  $\lfloor \rho n \rfloor$  parameters with the highest score  $\theta'_1, \dots, \theta'_{\lfloor \rho n \rfloor}$

    Update  $\mu$  and  $\sigma$ :  $\mu(j) = \frac{1}{\lfloor \rho n \rfloor} \sum_{i=1}^{\lfloor \rho n \rfloor} \theta'_i(j)$  and  $\sigma^2(j) = \frac{1}{\lfloor \rho n \rfloor} \sum_{i=1}^{\lfloor \rho n \rfloor} [\theta'_i(j) - \mu(j)]^2 + \eta$

Works embarrassingly well in low-dimensions, e.g., to search for a linear policy over the 22 Bertsekas features for Tetris.

# Cross-entropy method

Work embarrassingly well in low-dimensions

Method	Mean Score	Reference
<b>Nonreinforcement learning</b>		
Hand-coded	631,167	Dellacherie (Fahey, 2003)
Genetic algorithm	586,103	(Böhm et al., 2004)
<b>Reinforcement learning</b>		
Relational reinforcement learning+kernel-based regression	≈50	Ramon and Driessens (2004)
Policy iteration	3183	Bertsekas and Tsitsiklis (1996)
Least squares policy iteration	<3000	Lagoudakis, Parr, and Littman (2002)
Linear programming + Bootstrap	4274	Farias and van Roy (2006)
Natural policy gradient	≈6800	Kakade (2001)
CE+RL	21,252	
CE+RL, constant noise	72,705	
CE+RL, decreasing noise	348,895	

István Szita and András Lörincz. “Learning Tetris using the noisy cross-entropy method”. In: *Neural computation* 18.12 (2006), pp. 2936–2941

$$\mu \in \mathbb{R}^{22}$$

Much later:

---

Approximate Dynamic Programming Finally  
Performs Well in the Game of Tetris

[NIPS 2013]

**Victor Gabilon**  
INRIA Lille - Nord Europe,  
Team SequeL, FRANCE  
*victor.gabilon@inria.fr*

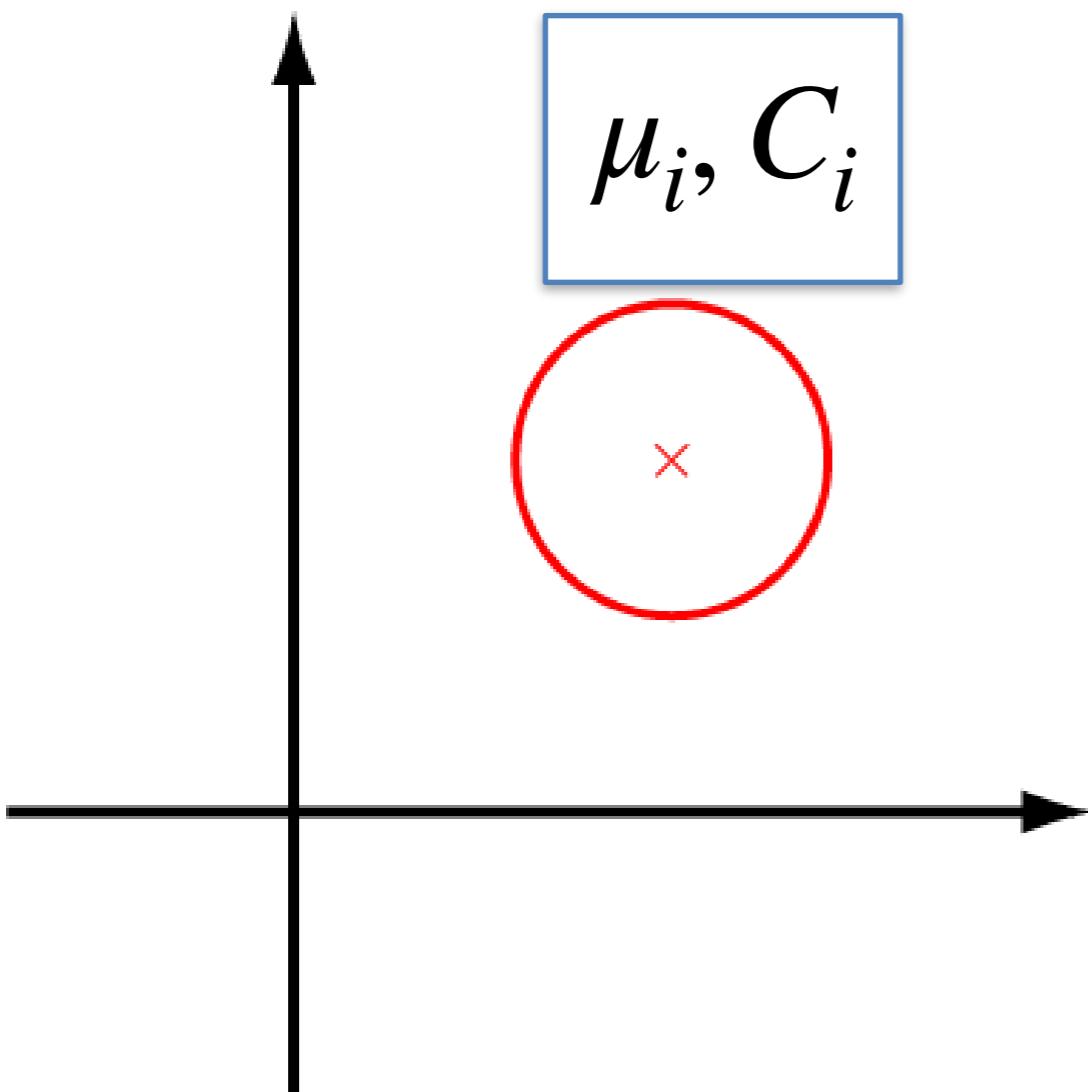
**Mohammad Ghavamzadeh\***  
INRIA Lille - Team SequeL  
& Adobe Research  
*mohammad.ghavamzadeh@inria.fr*

**Bruno Scherrer**  
INRIA Nancy - Grand Est,  
Team Maia, FRANCE  
*bruno.scherrer@inria.fr*

# Covariance Matrix Adaptation (CMA-ES)

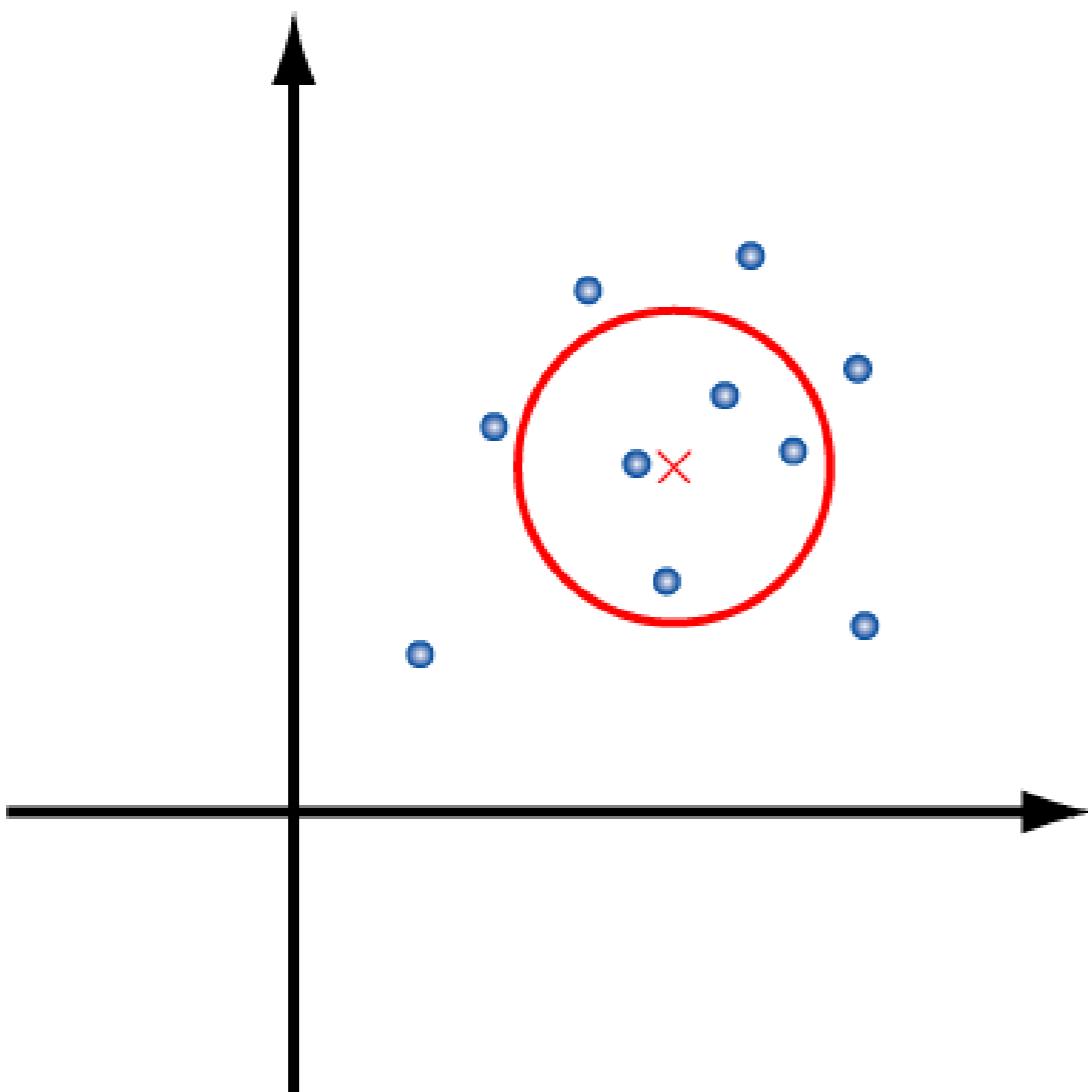
- Policy parameters are sampled from a multivariate Gaussian distribution with a full covariance matrix.
- We will update the mean and the covariance matrix to model samples that have highest fitness scores.

# Covariance Matrix Adaptation



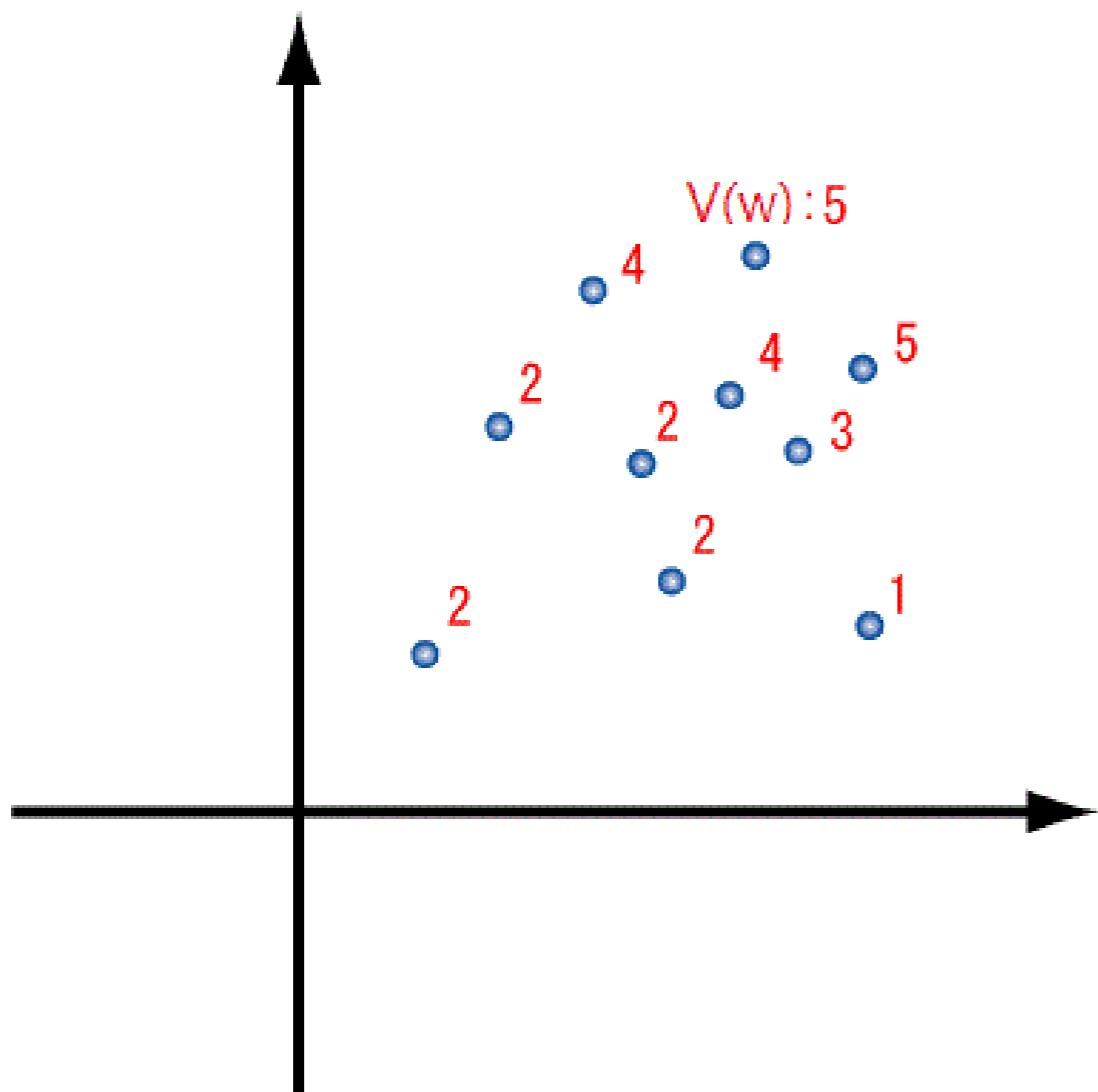
- Sample
- Select elites
- Update mean
- Update covariance
- iterate

# Covariance Matrix Adaptation



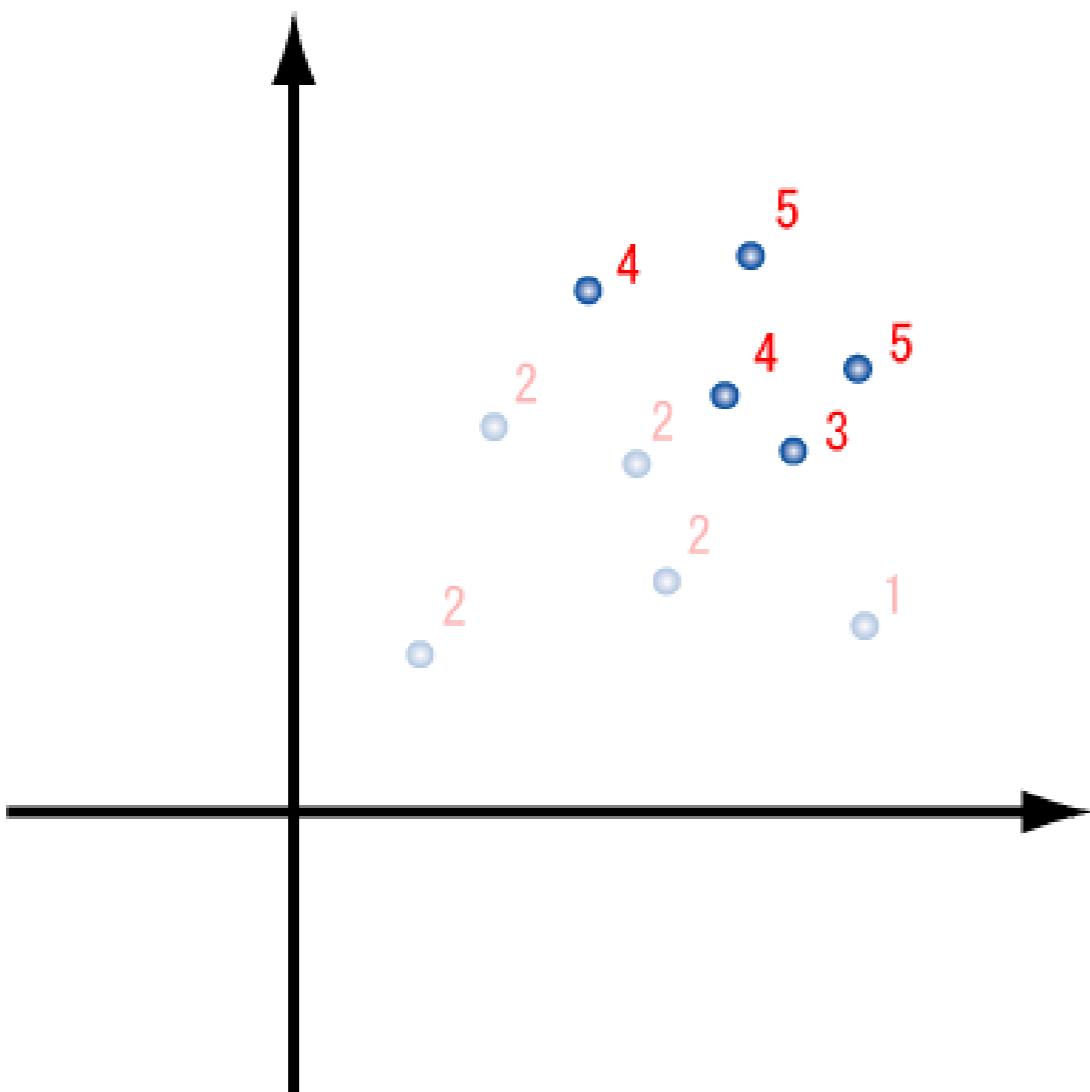
- Sample
- Select elites
- Update mean
- Update covariance
- iterate

# Covariance Matrix Adaptation



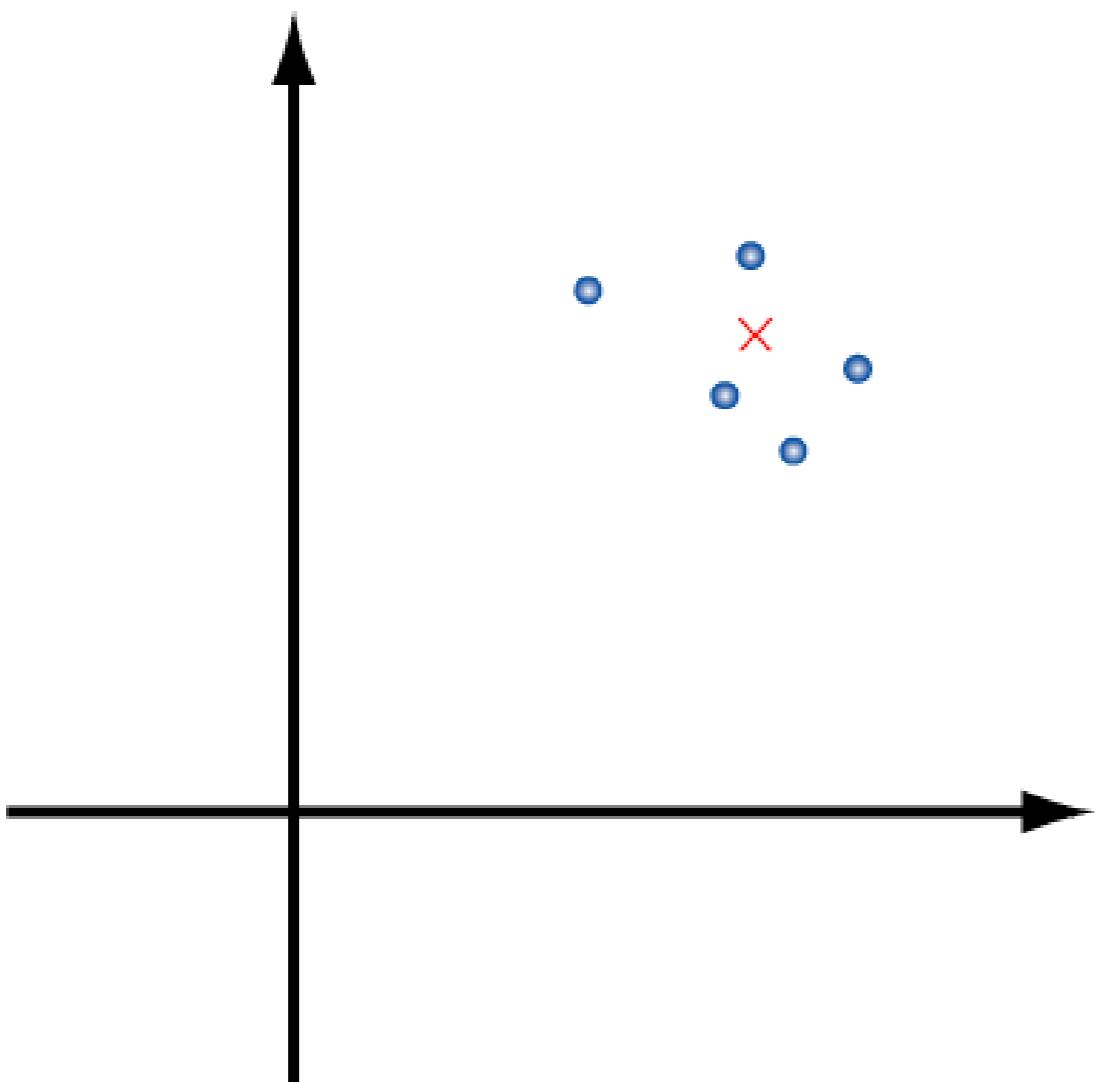
- Sample
- Select elites
- Update mean
- Update covariance
- iterate

# Covariance Matrix Adaptation



- Sample
- Select elites
- Update mean
- Update covariance
- iterate

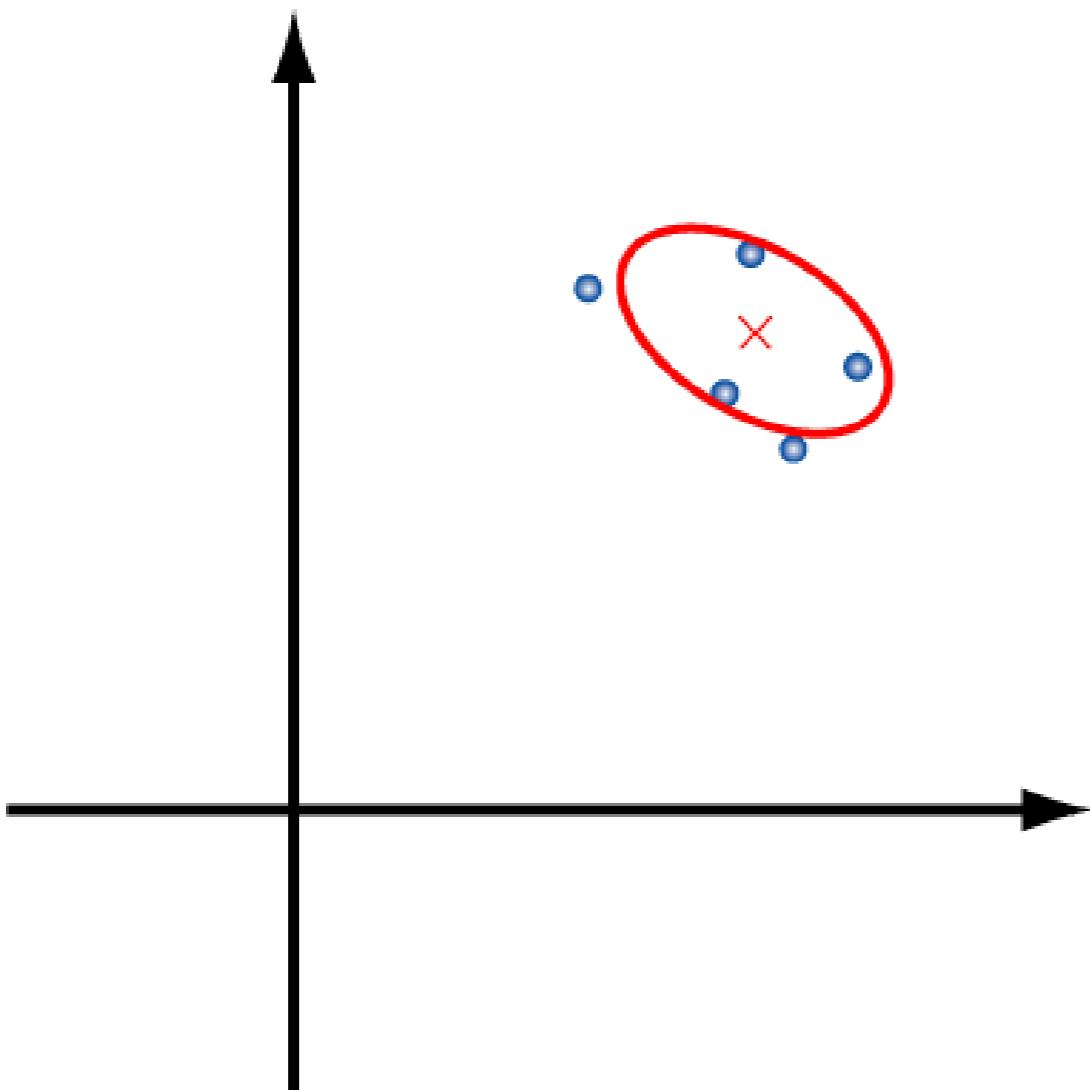
# Covariance Matrix Adaptation



- Sample
- Select elites
- **Update mean**
- Update covariance
- iterate

$$\mu_{t+1} = \mu_t + \alpha \sum_{i=1}^{n_{elit}} w_i (\theta_i^{elit,t} - \mu_t)$$
$$\mu_{t+1} = \sum_{i=1}^{n_{elit}} w_i \theta_i^{elit,t}$$

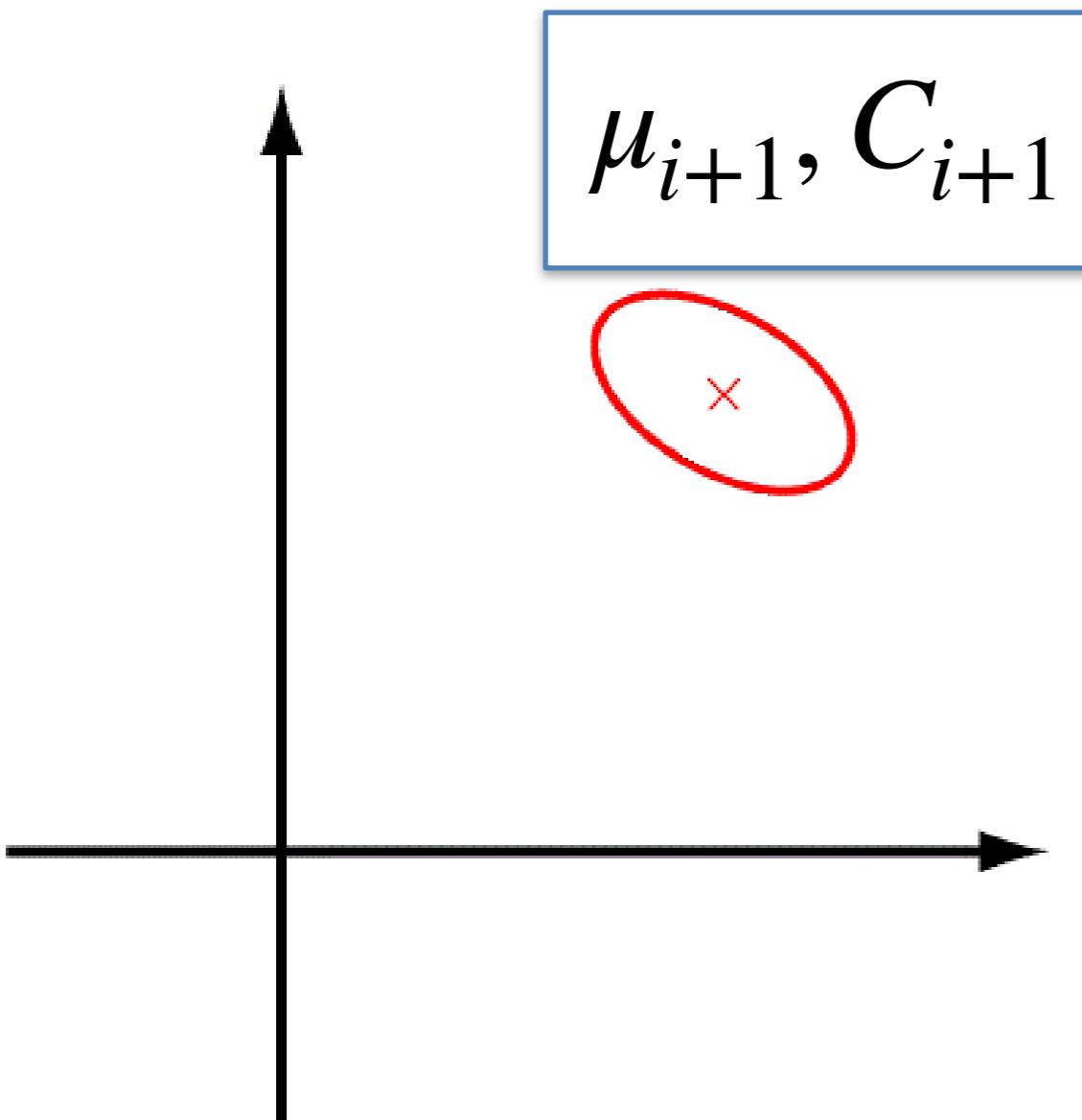
# Covariance Matrix Adaptation



- Sample
- Select elites
- Update mean
- **Update covariance**
- iterate

$$\Sigma_{t+1} = \text{Cov}(\theta_1^{elit,t}, \theta_2^{elit,t}, \dots) + \epsilon I$$

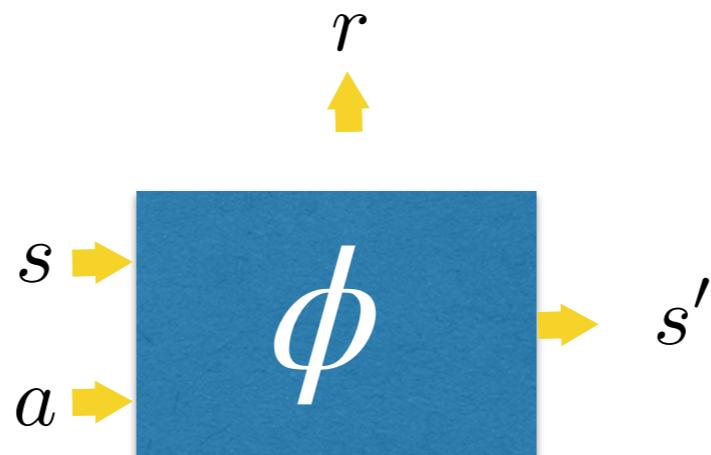
# Covariance Matrix Adaptation



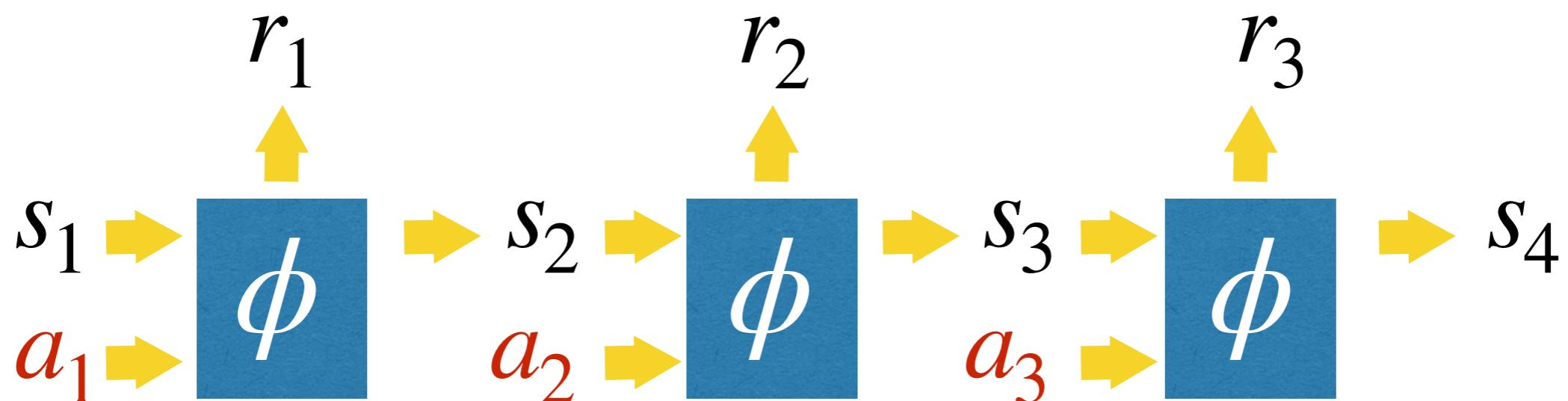
- Sample
- Select elites
- Update mean
- Update covariance
- **iterate**

# Evolution for Action Selection

- A learned model:

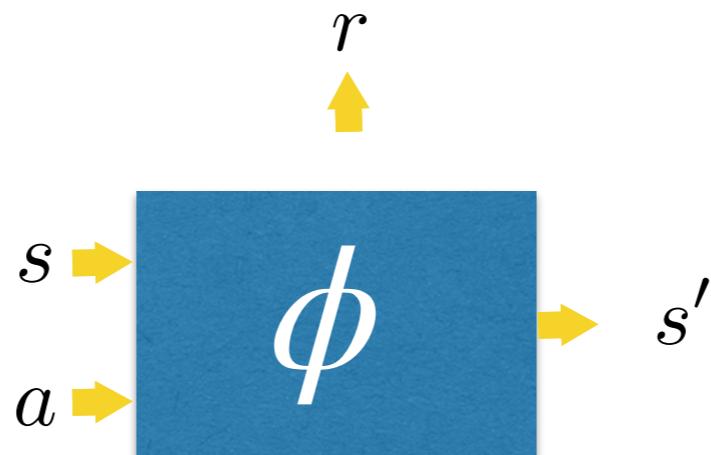


- Objective: maximizing sum of rewards



# Evolution for Action Selection

- A learned model:



- Objective: reaching a desired state:  $\|s_4 - s_*\|$



# Natural Evolutionary Strategies (NES)

- In CEM and CMA-ES, we have been selecting the best (elite) parameter offsprings to update the parameter distribution.
- NES considers **every** offspring.

# Natural Evolutionary Strategies

- Consider the parameters of our policy  $\theta \in \mathbb{R}^d$  follow a Gaussian distribution with mean  $\mu \in \mathbb{R}^d$  and diagonal and fixed covariance matrix  $\sigma^2 \mathbf{I}_d$ :  $\theta \sim P_\mu(\theta)$ .
- Goal:  $\max_{\mu} . \mathbb{E}_{\theta \sim P_\mu(\theta)} F(\theta)$ , where fitness score  $F(\theta) = \mathbb{E}_{\tau \sim \pi_\theta, s_0 \sim \mu_0(s)} R(\tau)$ .

# Natural Evolutionary Strategies

- Consider the parameters of our policy  $\theta \in \mathbb{R}^d$  follow a Gaussian distribution with mean  $\mu \in \mathbb{R}^d$  and diagonal and fixed covariance matrix  $\sigma^2 \mathbf{I}_d$ :  $\theta \sim P_\mu(\theta)$ .
- Goal:  $\max_{\mu} . \mathbb{E}_{\theta \sim P_\mu(\theta)} F(\theta)$ , where fitness score  $F(\theta) = \mathbb{E}_{\tau \sim \pi_\theta, s_0 \sim \mu_0(s)} R(\tau)$ .

$$\nabla_\mu \mathbb{E}_{\theta \sim P_\mu(\theta)} [F(\theta)] = \nabla_\mu \int P_\mu(\theta) F(\theta) d\theta$$

# Natural Evolutionary Strategies

- Consider the parameters of our policy  $\theta \in \mathbb{R}^d$  follow a Gaussian distribution with mean  $\mu \in \mathbb{R}^d$  and diagonal and fixed covariance matrix  $\sigma^2 \mathbf{I}_d$ :  $\theta \sim P_\mu(\theta)$ .
- Goal:  $\max_{\mu} . \mathbb{E}_{\theta \sim P_\mu(\theta)} F(\theta)$ , where fitness score  $F(\theta) = \mathbb{E}_{\tau \sim \pi_\theta, s_0 \sim \mu_0(s)} R(\tau)$ .

$$\begin{aligned}\nabla_\mu \mathbb{E}_{\theta \sim P_\mu(\theta)} [F(\theta)] &= \nabla_\mu \int P_\mu(\theta) F(\theta) d\theta \\ &= \int \nabla_\mu P_\mu(\theta) F(\theta) d\theta\end{aligned}$$

# Natural Evolutionary Strategies

- Consider the parameters of our policy  $\theta \in \mathbb{R}^d$  follow a Gaussian distribution with mean  $\mu \in \mathbb{R}^d$  and diagonal and fixed covariance matrix  $\sigma^2 \mathbf{I}_d$ :  $\theta \sim P_\mu(\theta)$ .
- Goal:  $\max_{\mu} . \mathbb{E}_{\theta \sim P_\mu(\theta)} F(\theta)$ , where fitness score  $F(\theta) = \mathbb{E}_{\tau \sim \pi_\theta, s_0 \sim \mu_0(s)} R(\tau)$ .

$$\begin{aligned}\nabla_\mu \mathbb{E}_{\theta \sim P_\mu(\theta)} [F(\theta)] &= \nabla_\mu \int P_\mu(\theta) F(\theta) d\theta \\ &= \int \nabla_\mu P_\mu(\theta) F(\theta) d\theta \\ &= \int P_\mu(\theta) \frac{\nabla_\mu P_\mu(\theta)}{P_\mu(\theta)} F(\theta) d\theta\end{aligned}$$

# Natural Evolutionary Strategies

- Consider the parameters of our policy  $\theta \in \mathbb{R}^d$  follow a Gaussian distribution with mean  $\mu \in \mathbb{R}^d$  and diagonal and fixed covariance matrix  $\sigma^2 \mathbf{I}_d$ :  $\theta \sim P_\mu(\theta)$ .
- Goal:  $\max_{\mu} . \mathbb{E}_{\theta \sim P_\mu(\theta)} F(\theta)$ , where fitness score  $F(\theta) = \mathbb{E}_{\tau \sim \pi_\theta, s_0 \sim \mu_0(s)} R(\tau)$ .

$$\begin{aligned}\nabla_\mu \mathbb{E}_{\theta \sim P_\mu(\theta)} [F(\theta)] &= \nabla_\mu \int P_\mu(\theta) F(\theta) d\theta \\ &= \int \nabla_\mu P_\mu(\theta) F(\theta) d\theta \\ &= \int P_\mu(\theta) \frac{\nabla_\mu P_\mu(\theta)}{P_\mu(\theta)} F(\theta) d\theta \\ &= \int P_\mu(\theta) \nabla_\mu \log P_\mu(\theta) F(\theta) d\theta\end{aligned}$$

# Natural Evolutionary Strategies

- Consider the parameters of our policy  $\theta \in \mathbb{R}^d$  follow a Gaussian distribution with mean  $\mu \in \mathbb{R}^d$  and diagonal and fixed covariance matrix  $\sigma^2 \mathbf{I}_d$ :  $\theta \sim P_\mu(\theta)$ .
- Goal:  $\max_{\mu} . \mathbb{E}_{\theta \sim P_\mu(\theta)} F(\theta)$ , where fitness score  $F(\theta) = \mathbb{E}_{\tau \sim \pi_\theta, s_0 \sim \mu_0(s)} R(\tau)$ .

$$\begin{aligned}\nabla_\mu \mathbb{E}_{\theta \sim P_\mu(\theta)} [F(\theta)] &= \nabla_\mu \int P_\mu(\theta) F(\theta) d\theta \\ &= \int \nabla_\mu P_\mu(\theta) F(\theta) d\theta \\ &= \int P_\mu(\theta) \frac{\nabla_\mu P_\mu(\theta)}{P_\mu(\theta)} F(\theta) d\theta \\ &= \int P_\mu(\theta) \nabla_\mu \log P_\mu(\theta) F(\theta) d\theta \\ &= \mathbb{E}_{\theta \sim P_\mu(\theta)} [\nabla_\mu \log P_\mu(\theta) F(\theta)]\end{aligned}$$

We approximate expectations by sampling!

# Natural Evolutionary Strategies

- Consider the parameters of our policy  $\theta \in \mathbb{R}^d$  follow a Gaussian distribution with mean  $\mu \in \mathbb{R}^d$  and diagonal and fixed covariance matrix  $\sigma^2 \mathbf{I}_d$ :  $\theta \sim P_\mu(\theta)$ .
- Goal:  $\max_{\mu} . \mathbb{E}_{\theta \sim P_\mu(\theta)} F(\theta)$ , where fitness score  $F(\theta) = \mathbb{E}_{\tau \sim \pi_\theta, s_0 \sim \mu_0(s)} R(\tau)$ .

$$\begin{aligned}\nabla_\mu \mathbb{E}_{\theta \sim P_\mu(\theta)} [F(\theta)] &= \nabla_\mu \int P_\mu(\theta) F(\theta) d\theta \\ &= \int \nabla_\mu P_\mu(\theta) F(\theta) d\theta \\ &= \int P_\mu(\theta) \frac{\nabla_\mu P_\mu(\theta)}{P_\mu(\theta)} F(\theta) d\theta \\ &= \int P_\mu(\theta) \nabla_\mu \log P_\mu(\theta) F(\theta) d\theta \\ &= \mathbb{E}_{\theta \sim P_\mu(\theta)} \left[ \nabla_\mu \log P_\mu(\theta) F(\theta) \right] \\ &\approx \frac{1}{N} \sum_{i=1}^N \left[ \nabla_\mu \log P_\mu(\theta) |_{\theta=\theta_i} F(\theta_i) \right] \quad \theta_i \sim P_\mu(\theta)\end{aligned}$$

# Policy Gradients

$$\begin{aligned} \max_{\theta} . \quad U(\theta) &= \mathbb{E}_{\tau \sim P_{\theta}(\tau)} [R(\tau)] \\ \nabla_{\theta} U(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim P_{\theta}(\tau)} [R(\tau)] \\ &= \nabla_{\theta} \sum_{\tau} P_{\theta}(\tau) R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P_{\theta}(\tau) R(\tau) \\ &= \sum_{\tau} P_{\theta}(\tau) \frac{\nabla_{\theta} P_{\theta}(\tau)}{P_{\theta}(\tau)} R(\tau) \\ &= \sum_{\tau} P_{\theta}(\tau) \nabla_{\theta} \log P_{\theta}(\tau) R(\tau) \\ &= \mathbb{E}_{\tau \sim P_{\theta}(\tau)} [\nabla_{\theta} \log P_{\theta}(\tau) R(\tau)] \end{aligned}$$

Sample estimate:

$$\nabla_{\theta} U(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log P_{\theta}(\tau^{(i)}) R(\tau^{(i)})$$

# Evolutionary Methods

$$\max_{\mu} . \ U(\mu) = \mathbb{E}_{\theta \sim P_\mu(\theta)} [F(\theta)]$$

$$\begin{aligned}\nabla_\mu U(\mu) &= \nabla_\mu \mathbb{E}_{\theta \sim P_\mu(\theta)} [F(\theta)] \\ &= \nabla_\mu \int P_\mu(\theta) F(\theta) d\theta \\ &= \int \nabla_\mu P_\mu(\theta) F(\theta) d\theta \\ &= \int P_\mu(\theta) \frac{\nabla_\mu P_\mu(\theta)}{P_\mu(\theta)} F(\theta) d\theta \\ &= \int P_\mu(\theta) \nabla_\mu \log P_\mu(\theta) F(\theta) d\theta \\ &= \mathbb{E}_{\theta \sim P_\mu(\theta)} [\nabla_\mu \log P_\mu(\theta) F(\theta)]\end{aligned}$$

Sample estimate:

$$\nabla_\mu U(\mu) \approx \frac{1}{N} \sum_{i=1}^N \nabla_\mu \log P_\mu(\theta^{(i)}) F(\theta^{(i)})$$

# Policy gradients VS Evolutionary methods

Considers distribution over actions

$$\max_{\theta} . \quad U(\theta) = \mathbb{E}_{\tau \sim P_{\theta}(\tau)} [R(\tau)]$$

$$\begin{aligned}\nabla_{\theta} U(\theta) &= \nabla_{\theta} \mathbb{E}_{\tau \sim P_{\theta}(\tau)} [R(\tau)] \\ &= \nabla_{\theta} \sum_{\tau} P_{\theta}(\tau) R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P_{\theta}(\tau) R(\tau) \\ &= \sum_{\tau} P_{\theta}(\tau) \frac{\nabla_{\theta} P_{\theta}(\tau)}{P_{\theta}(\tau)} R(\tau) \\ &= \sum_{\tau} P_{\theta}(\tau) \nabla_{\theta} \log P_{\theta}(\tau) R(\tau) \\ &= \mathbb{E}_{\tau \sim P_{\theta}(\tau)} [\nabla_{\theta} \log P_{\theta}(\tau) R(\tau)]\end{aligned}$$

Sample estimate:

$$\nabla_{\theta} U(\theta) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\theta} \log P_{\theta}(\tau^{(i)}) R(\tau^{(i)})$$

Considers distribution over policy parameters

$$\max_{\mu} . \quad U(\mu) = \mathbb{E}_{\theta \sim P_{\mu}(\theta)} [F(\theta)]$$

$$\begin{aligned}\nabla_{\mu} U(\mu) &= \nabla_{\mu} \mathbb{E}_{\theta \sim P_{\mu}(\theta)} [F(\theta)] \\ &= \nabla_{\mu} \int P_{\mu}(\theta) F(\theta) d\theta \\ &= \int \nabla_{\mu} P_{\mu}(\theta) F(\theta) d\theta \\ &= \int P_{\mu}(\theta) \frac{\nabla_{\mu} P_{\mu}(\theta)}{P_{\mu}(\theta)} F(\theta) d\theta \\ &= \int P_{\mu}(\theta) \nabla_{\mu} \log P_{\mu}(\theta) F(\theta) d\theta \\ &= \mathbb{E}_{\theta \sim P_{\mu}(\theta)} [\nabla_{\mu} \log P_{\mu}(\theta) F(\theta)]\end{aligned}$$

Sample estimate:

$$\nabla_{\mu} U(\mu) \approx \frac{1}{N} \sum_{i=1}^N \nabla_{\mu} \log P_{\mu}(\theta^{(i)}) F(\theta^{(i)})$$

# Policy gradients VS Evolutionary methods

- We are sampling in both cases...
  - PG: sampling in action space
  - ES: sampling in parameter space

# Reminder: Sampling from a multivariate Gaussian

We want to generate samples  $\theta_i \sim P_\mu(\theta)$  where  $P$  is the Gaussian distribution with mean  $\mu \in \mathbb{R}^d$  and diagonal and fixed covariance matrix  $\sigma^2 \mathbf{I}_d$ .

Imagine we have access to random vectors  $\epsilon_i \in \mathbb{R}^d$ ,  $\epsilon_i \sim \mathcal{N}(0, \mathbf{I}_d)$

$$\begin{aligned}\theta_1 &= \mu + \sigma \epsilon_1 \\ \theta_2 &= \mu + \sigma \epsilon_2\end{aligned}$$

The samples have the desired mean and variance

# Reparametrization Trick

- A sample from a normal distribution  $z \sim N(\mu, \Sigma)$  can be rewritten as follows:

$$z = \mu + \Sigma^{\frac{1}{2}}\epsilon \text{ where } \epsilon \sim N(0, I).$$

- We just need a standard normal sampler to sample from an arbitrary normal distribution.

# Cumulative Distribution Function (CDF)

Given a PDF  $p(x)$ ,

- CDF  $F_X(x): \mathbb{R} \rightarrow [0,1]$ ,

$$F_X(x) = \Pr(X \leq x) = \int_{-\infty}^x p(t)dt$$

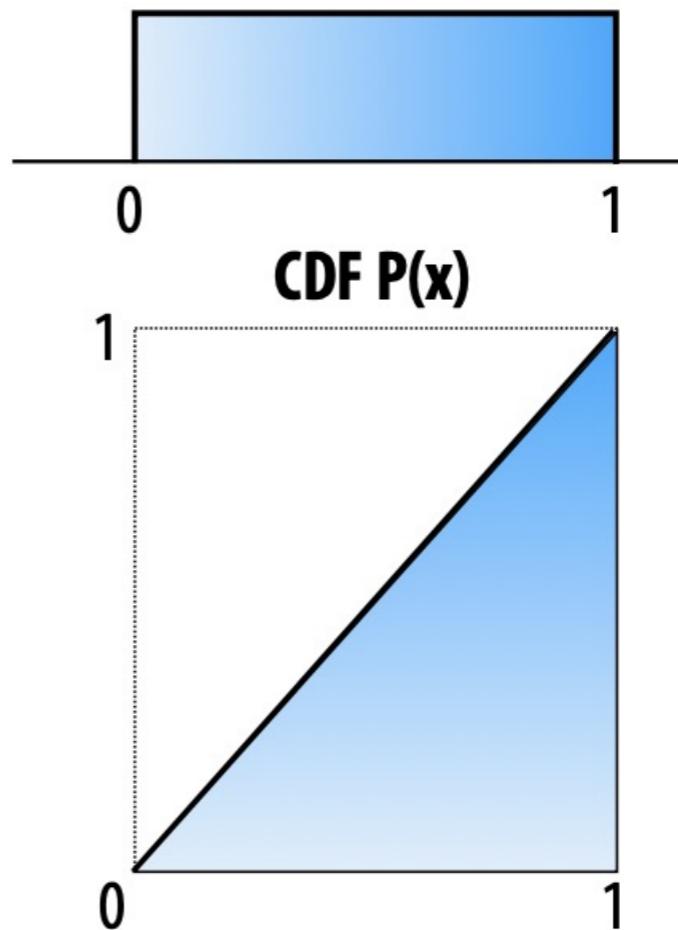
- $\Pr(a < X \leq b)$

$$= \int_a^b p(t)dt = F_X(b) - F_X(a)$$

- $F_X(1) = 1$

**Uniform distribution:  $p(x) = c$**

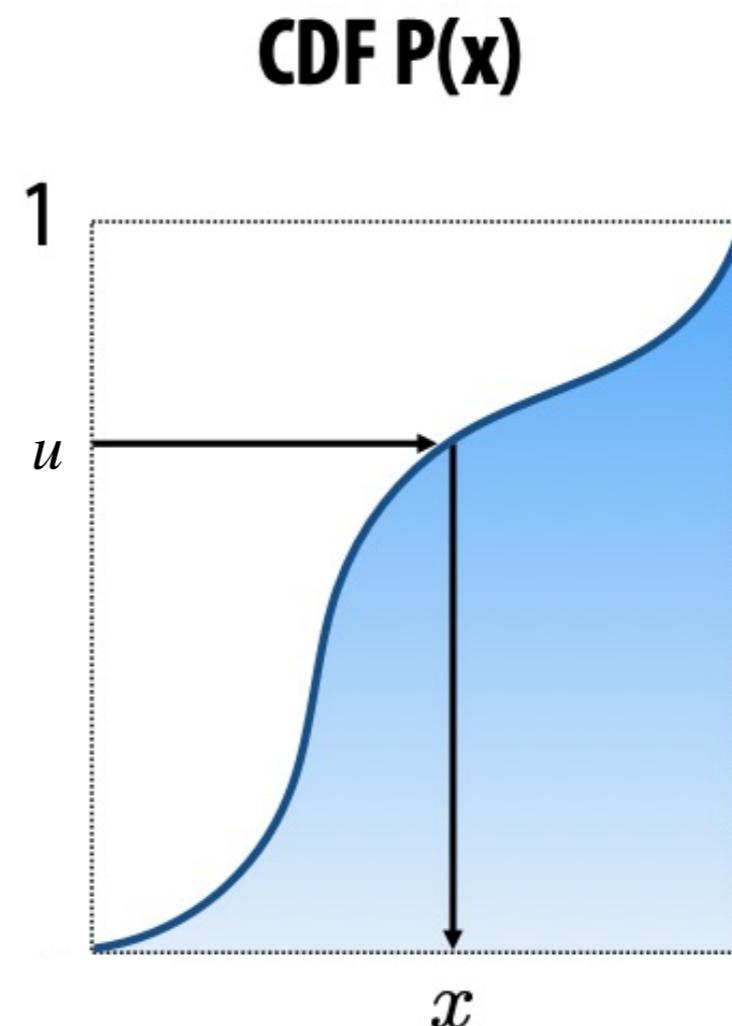
(for random variable  $X$  defined on  $[0,1]$  domain)



# Inverse Transform Sampling

To randomly sample based on the given PDF,

- Compute CDF  $F_X(x)$
- Draw  $u \sim \mathcal{U}(0, 1)$ .
- Take  $x = F_X^{-1}(u)$ .



# Reparametrization Trick: Sampling from a general Gaussian density

$$x_i \sim \mathcal{N}(0,1)$$

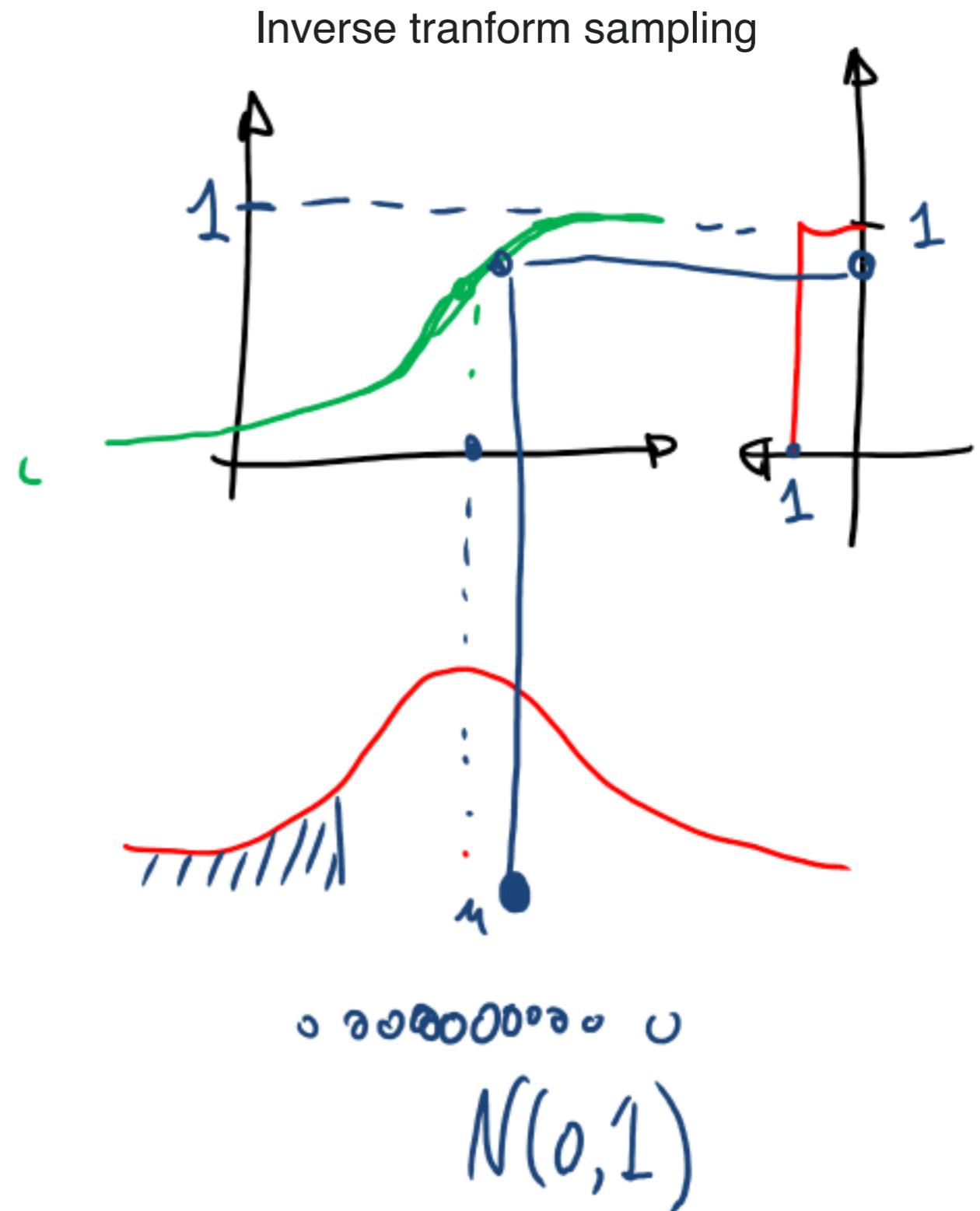
$$x_i \sim \mathcal{N}(\mu, \sigma^2)$$

$$\sim \mu + \sigma \mathcal{N}(0,1)$$

$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}_i \sim \mathcal{N} \left( \begin{bmatrix} \mu_1 \\ \mu_2 \end{bmatrix}, \begin{bmatrix} \Sigma_{11} & \Sigma_{12} \\ \Sigma_{21} & \Sigma_{22} \end{bmatrix} \right)$$

$$\mathbf{x} \sim \mathcal{N}(\mu, \Sigma) \quad x \sim \mu + L\mathcal{N}(0, I)$$

Cholesky decomposition  $\Sigma = \underline{LL}^T$



# Reparametrization Trick: Sampling from a general Gaussian density

- A sample from a normal distribution  $z \sim N(\mu, \Sigma)$  can be rewritten as follows:

$$z = \mu + \Sigma^{\frac{1}{2}}\epsilon \text{ where } \epsilon \sim N(0, I).$$

- We just need a standard normal sampler to sample from an arbitrary normal distribution.

# A concrete example

- Suppose  $\theta \sim P_\mu(\theta)$  is a Gaussian distribution with mean  $\mu$ , and covariance matrix  $\sigma^2 \mathbf{I}_d$ .

$$\log P_\mu(\theta) = -\frac{\|\theta - \mu\|^2}{2\sigma^2} + \text{const}$$

$$\nabla_\mu \log P_\mu(\theta) = \frac{\theta - \mu}{\sigma^2}$$

# A concrete example

- Suppose  $\theta \sim P_\mu(\theta)$  is a Gaussian distribution with mean  $\mu$ , and covariance matrix  $\sigma^2 \mathbf{I}_d$ .

$$\log P_\mu(\theta) = -\frac{\|\theta - \mu\|^2}{2\sigma^2} + \text{const}$$

$$\nabla_\mu \log P_\mu(\theta) = \frac{\theta - \mu}{\sigma^2}$$

- We draw two parameter samples  $\theta_1, \theta_2$ . For each sampled parameter vector, we run the policy and obtain a set of trajectories, or a single trajectory. Then:

$$\mathbb{E}_{\theta \sim P_\mu(\theta)} \left[ \nabla_\mu \log P_\mu(\theta) F(\theta) \right] \approx \frac{1}{N} \sum_{i=1}^N \left[ \nabla_\mu \log P_\mu(\theta) \Big|_{\theta=\theta_i} F(\theta_i) \right]$$

# A concrete example

- Suppose  $\theta \sim P_\mu(\theta)$  is a Gaussian distribution with mean  $\mu$ , and covariance matrix  $\sigma^2 \mathbf{I}_d$ .

$$\log P_\mu(\theta) = -\frac{\|\theta - \mu\|^2}{2\sigma^2} + \text{const}$$

$$\nabla_\mu \log P_\mu(\theta) = \frac{\theta - \mu}{\sigma^2}$$

- We draw two parameter samples  $\theta_1, \theta_2$ . For each sampled parameter vector, we run the policy and obtain a set of trajectories, or a single trajectory. Then:

$$\mathbb{E}_{\theta \sim P_\mu(\theta)} \left[ \nabla_\mu \log P_\mu(\theta) F(\theta) \right] \approx \frac{1}{N} \sum_{i=1}^N \left[ \nabla_\mu \log P_\mu(\theta) \Big|_{\theta=\theta_i} F(\theta_i) \right]$$

$$\approx \frac{1}{2} \left[ F(\theta_1) \frac{\theta_1 - \mu}{\sigma^2} + F(\theta_2) \frac{\theta_2 - \mu}{\sigma^2} \right], \text{ where } F(\theta_1) = R(\tau_1), F(\theta_2) = R(\tau_2)$$

- **Q:** Can we simplify this expression more?

# A concrete example

- Suppose  $\theta \sim P_\mu(\theta)$  is a Gaussian distribution with mean  $\mu$ , and covariance matrix  $\sigma^2 \mathbf{I}_d$ .

$$\log P_\mu(\theta) = -\frac{\|\theta - \mu\|^2}{2\sigma^2} + \text{const}$$

$$\nabla_\mu \log P_\mu(\theta) = \frac{\theta - \mu}{\sigma^2}$$

- We draw two parameter samples  $\theta_1, \theta_2$ . For each sampled parameter vector, we run the policy and obtain a set of trajectories, or a single trajectory. Then:

$$\mathbb{E}_{\theta \sim P_\mu(\theta)} \left[ \nabla_\mu \log P_\mu(\theta) F(\theta) \right] \approx \frac{1}{N} \sum_{i=1}^N \left[ \nabla_\mu \log P_\mu(\theta) \Big|_{\theta=\theta_i} F(\theta_i) \right]$$

$$\boxed{\begin{aligned}\theta_1 &= \mu + \sigma \epsilon_1 \\ \theta_2 &= \mu + \sigma \epsilon_2\end{aligned}}$$

$$\approx \frac{1}{2} \left[ F(\theta_1) \frac{\theta_1 - \mu}{\sigma^2} + F(\theta_2) \frac{\theta_2 - \mu}{\sigma^2} \right], \text{ where } F(\theta_1) = R(\tau_1), F(\theta_2) = R(\tau_2)$$

# A concrete example

- Suppose  $\theta \sim P_\mu(\theta)$  is a Gaussian distribution with mean  $\mu$ , and covariance matrix  $\sigma^2 \mathbf{I}_d$ .

$$\log P_\mu(\theta) = -\frac{\|\theta - \mu\|^2}{2\sigma^2} + \text{const}$$

$$\nabla_\mu \log P_\mu(\theta) = \frac{\theta - \mu}{\sigma^2}$$

- We draw two parameter samples  $\theta_1, \theta_2$ . For each sampled parameter vector, we run the policy and obtain a set of trajectories, or a single trajectory. Then:

$$\mathbb{E}_{\theta \sim P_\mu(\theta)} \left[ \nabla_\mu \log P_\mu(\theta) F(\theta) \right] \approx \frac{1}{N} \sum_{i=1}^N \left[ \nabla_\mu \log P_\mu(\theta) \Big|_{\theta=\theta_i} F(\theta_i) \right]$$

$$\boxed{\begin{aligned}\theta_1 &= \mu + \sigma \epsilon_1 \\ \theta_2 &= \mu + \sigma \epsilon_2\end{aligned}}$$

$$\approx \frac{1}{2} \left[ F(\theta_1) \frac{\theta_1 - \mu}{\sigma^2} + F(\theta_2) \frac{\theta_2 - \mu}{\sigma^2} \right], \text{ where } F(\theta_1) = R(\tau_1), F(\theta_2) = R(\tau_2)$$

$$= \frac{1}{2\sigma} [F(\theta_1)\epsilon_1 + F(\theta_2)\epsilon_2]$$

# Natural Evolutionary Strategies (NES)

- In CEM and CMA-ES, we have been selecting the best (elite) parameter offsprings.
- NES considers **every** offspring.

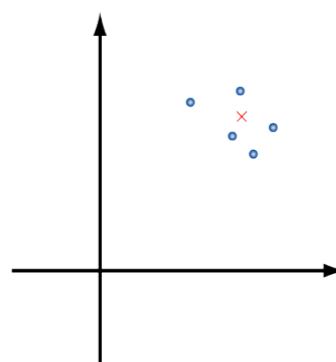
## Algorithm 1: Evolutionary Strategies

1. **Input:** Learning rate  $\alpha$ , noise standard deviation  $\sigma$ , initial policy parameters  $\theta_0$
2. **for**  $t = 0, 1, 2, \dots$  **do**
3.   Sample  $\epsilon_1, \dots, \epsilon_n \sim \mathcal{N}(0, \mathbf{I}_d)$
4.   Compute returns  $F_i = F(\mu_t + \sigma\epsilon_i)$  for  $i = 1, 2, \dots, n$
5.   Set  $\mu_{t+1} \leftarrow \mu_t + \alpha \frac{1}{n\sigma} \sum_{i=1}^n F_i \epsilon_i$
6. **end for**

# Compare the two update rules for the mean

NES:  $\mu_{t+1} = \mu_t + \frac{\alpha}{n\sigma} \sum_{i=1}^n F(\theta_i) \epsilon_i$  All offsprings participate

CMA-ES:  $\mu_{t+1} = \mu_t + \alpha \sum_{i=1}^{n_{elit}} w_i (\theta_i^{elit,t} - \mu_t)$



- Sample
- Select elites
- **Update mean**
- Update covariance
- iterate

# Question: Can evolutionary methods scale?

- Evolutionary methods work well on relatively low-dimensional problems (small number of parameter dimensions).
- Can they be used to optimize policies represented as neural networks of thousands parameters?

---

# Evolution Strategies as a Scalable Alternative to Reinforcement Learning

---

**Tim Salimans**

**Jonathan Ho**

**Xi Chen**

OpenAI

**Szymon Sidor**

**Ilya Sutskever**

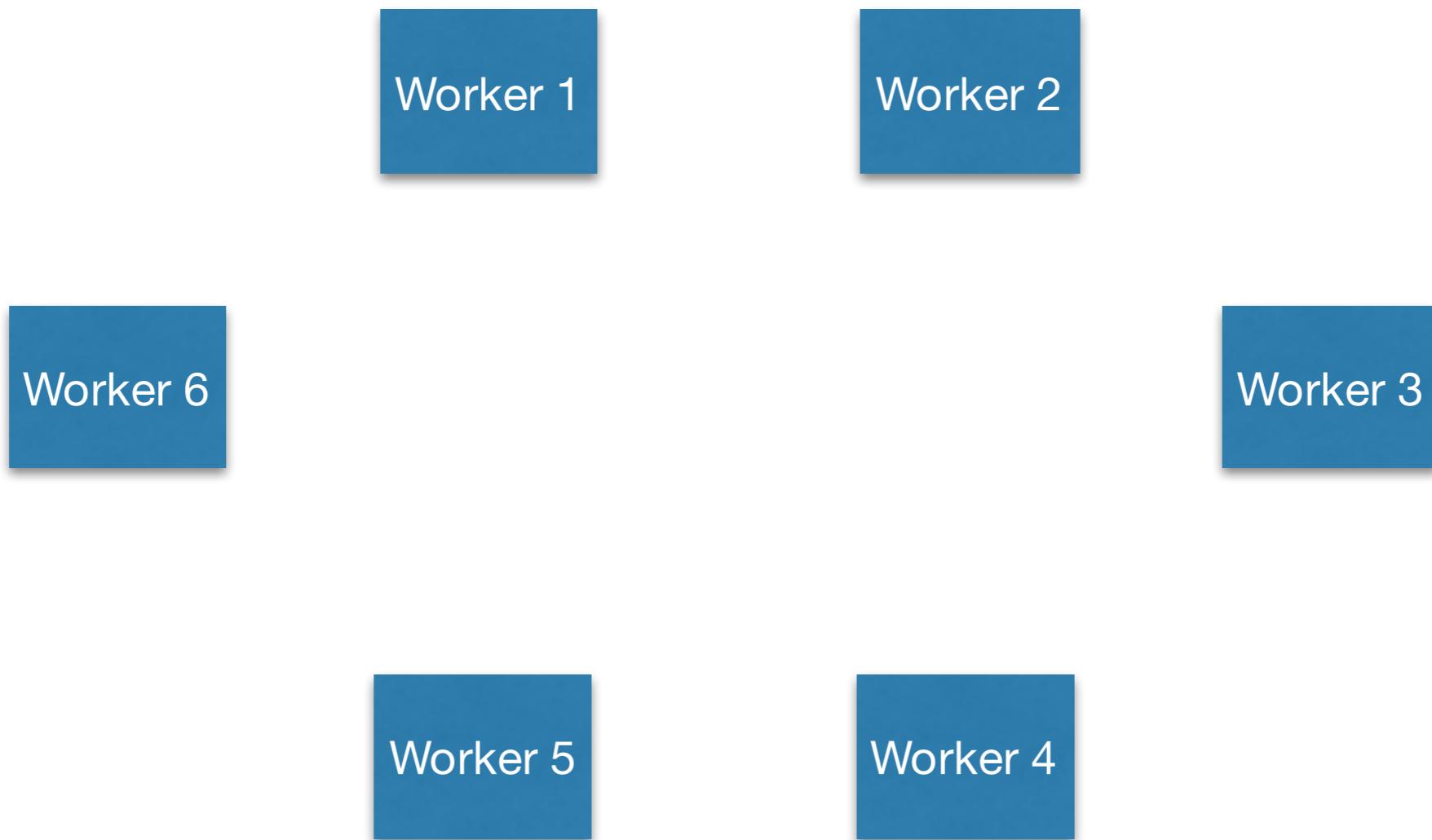
## Algorithm 1: Evolutionary Strategies

1. **Input:** Learning rate  $\alpha$ , noise standard deviation  $\sigma$ , initial policy parameters  $\theta_0$
2. **for**  $t = 0, 1, 2, \dots$  **do**
3.   Sample  $\epsilon_1, \dots, \epsilon_n \sim \mathcal{N}(0, \mathbf{I}_d)$
4.   Compute returns  $F_i = F(\mu_t + \sigma\epsilon_i)$  for  $i = 1, 2, \dots, n$
5.   Set  $\mu_{t+1} \leftarrow \mu_t + \alpha \frac{1}{n\sigma} \sum_{i=1}^n F_i \epsilon_i$
6. **end for**

Main contribution: Parallelization with a need for tiny only cross-worker communication

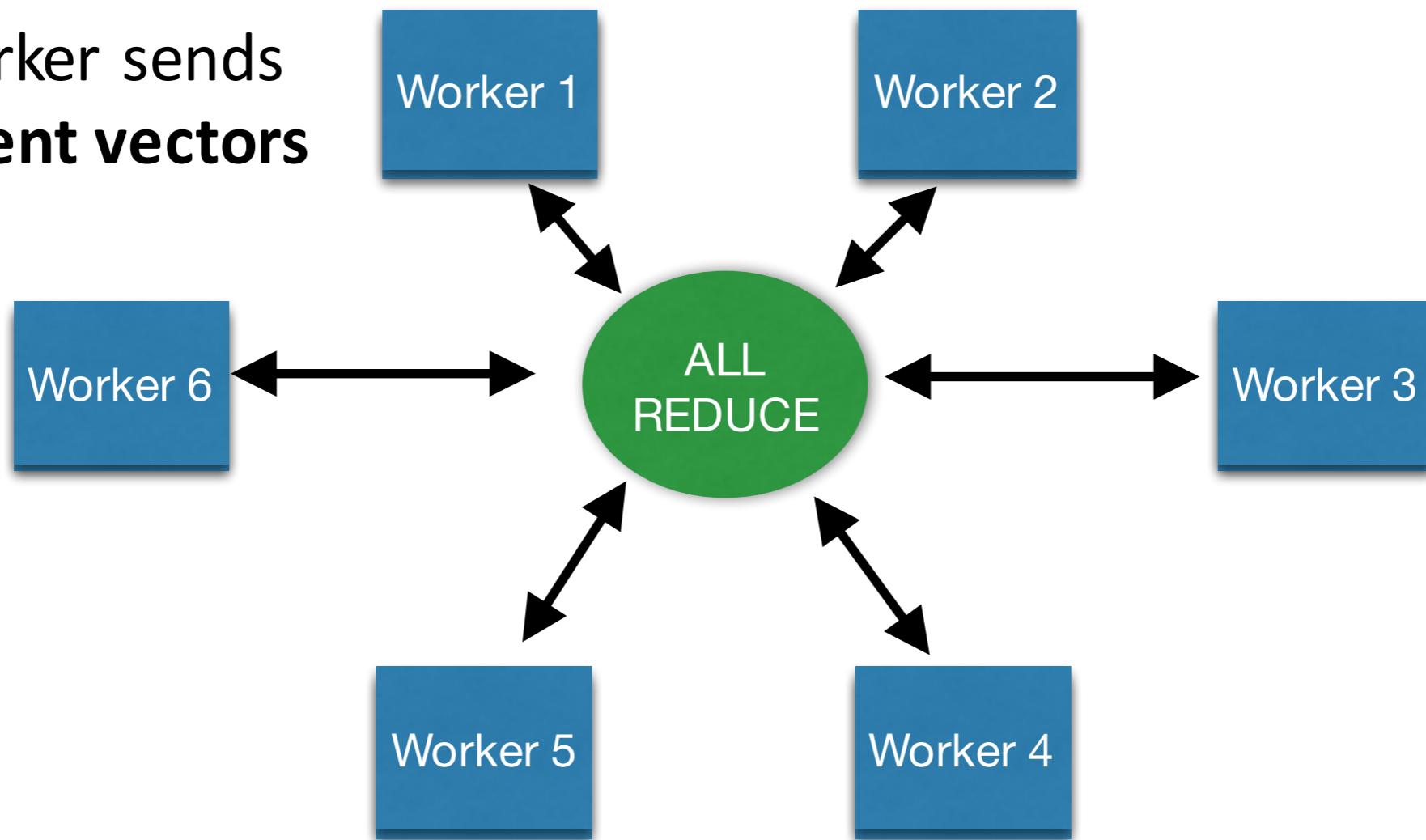
# Distributed stochastic gradient descent

Every worker get a set of examples, and computes a gradient vector.  
The master averages those gradient vectors.

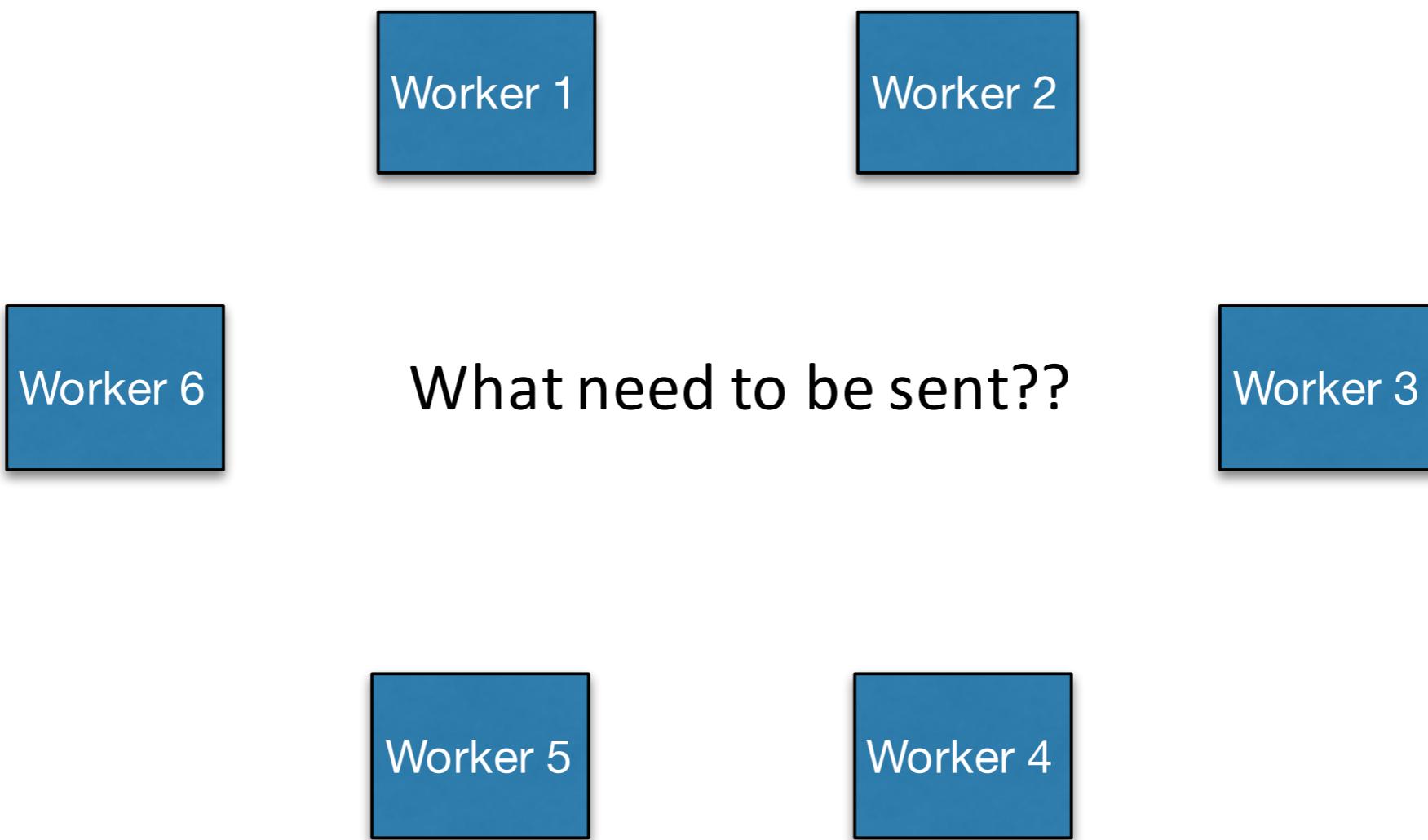


# Distributed stochastic gradient descent

Each worker sends  
**big gradient vectors**



# Distributed Evolution



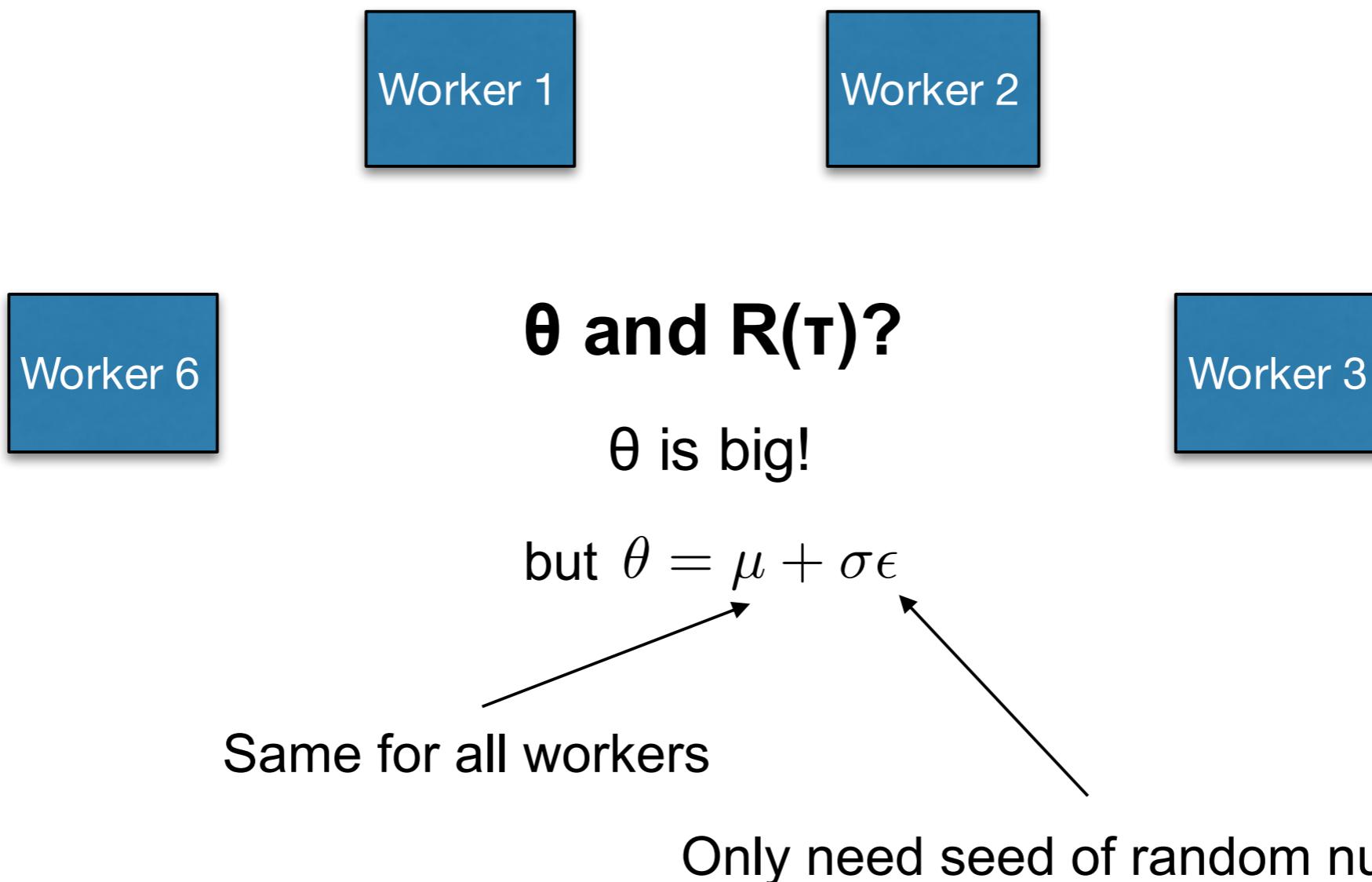
# Distributed Evolution



## Algorithm 1: Evolutionary Strategies

1. **Input:** Learning rate  $\alpha$ , noise standard deviation  $\sigma$ , initial policy parameters  $\theta_0$
2. **for**  $t = 0, 1, 2, \dots$  **do**
3.   Sample  $\epsilon_1, \dots, \epsilon_n \sim \mathcal{N}(0, \mathbf{I}_d)$
4.   Compute returns  $F_i = F(\mu_t + \sigma\epsilon_i)$  for  $i = 1, 2, \dots, n$
5.   Set  $\mu_{t+1} \leftarrow \mu_t + \alpha \frac{1}{n\sigma} \sum_{i=1}^n F_i \epsilon_i$
6. **end for**

# Distributed Evolution



# Distributed Evolution

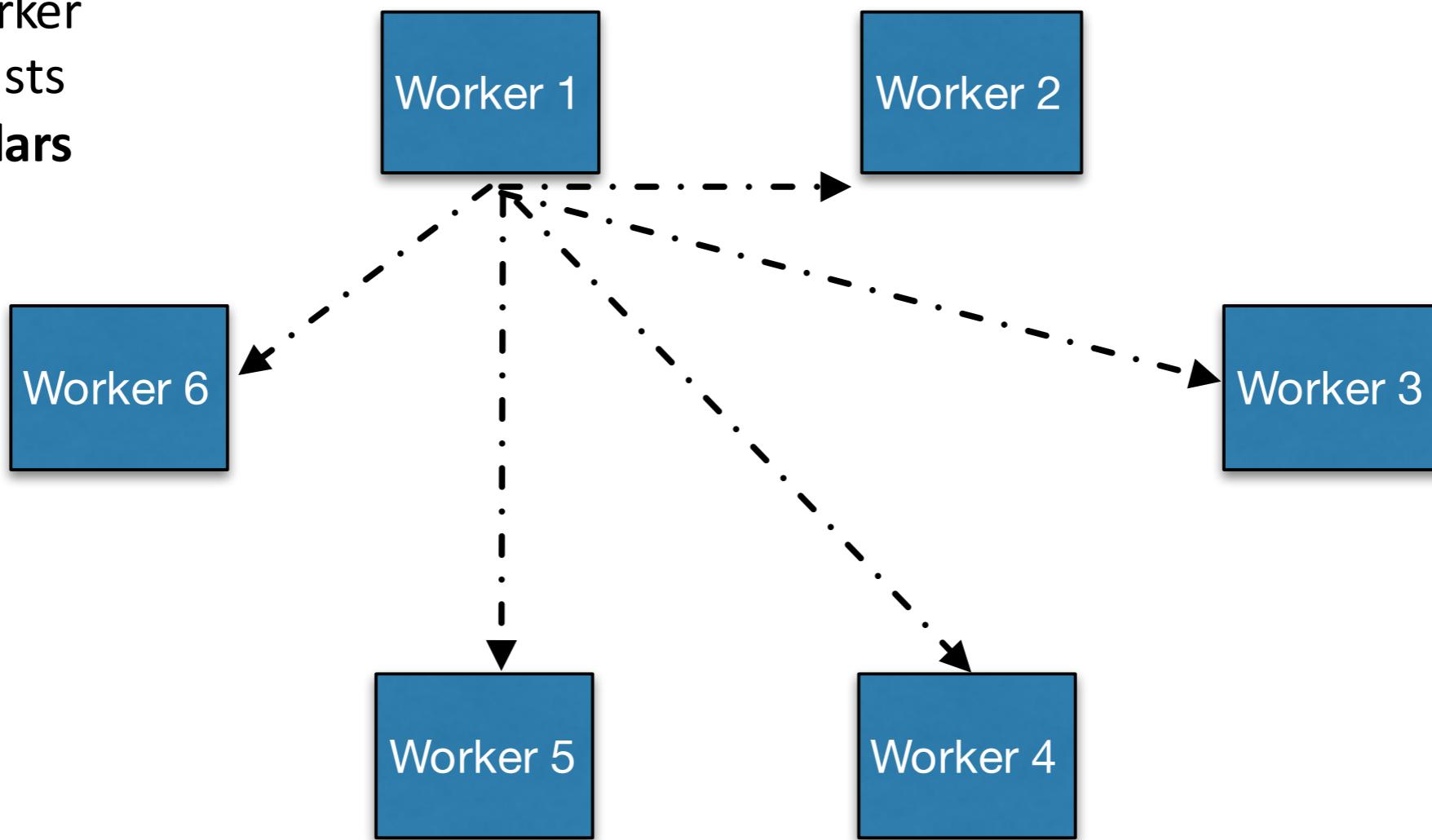
## Algorithm 2 Parallelized Evolution Strategies

- 1: **Input:** Learning rate  $\alpha$ , noise standard deviation  $\sigma$ , initial policy parameters  $\theta_0$
- 2: **Initialize:**  $n$  workers with known random seeds, and initial parameters  $\theta_0$
- 3: **for**  $t = 0, 1, 2, \dots$  **do**
- 4:   **for** each worker  $i = 1, \dots, n$  **do**
- 5:     Sample  $\epsilon_i \sim \mathcal{N}(0, I)$
- 6:     Compute returns  $F_i = F(\mu_t - \sigma\epsilon_i)$
- 7:   **end for**
- 8:   Send all scalar returns  $F_i$  from each worker to every other worker
- 9:   **for** each worker  $i = 1, \dots, n$  **do**
- 10:     Reconstruct all perturbations  $\epsilon_j$  for  $j = 1, \dots, n$
- 11:     Set  $\mu_{t+1} \leftarrow \mu_t + \alpha \frac{1}{n\sigma} \sum_{j=1}^n F_j \epsilon_j$
- 12:   **end for**
- 13: **end for**

[Salimans, Ho, Chen, Sutskever, 2017]

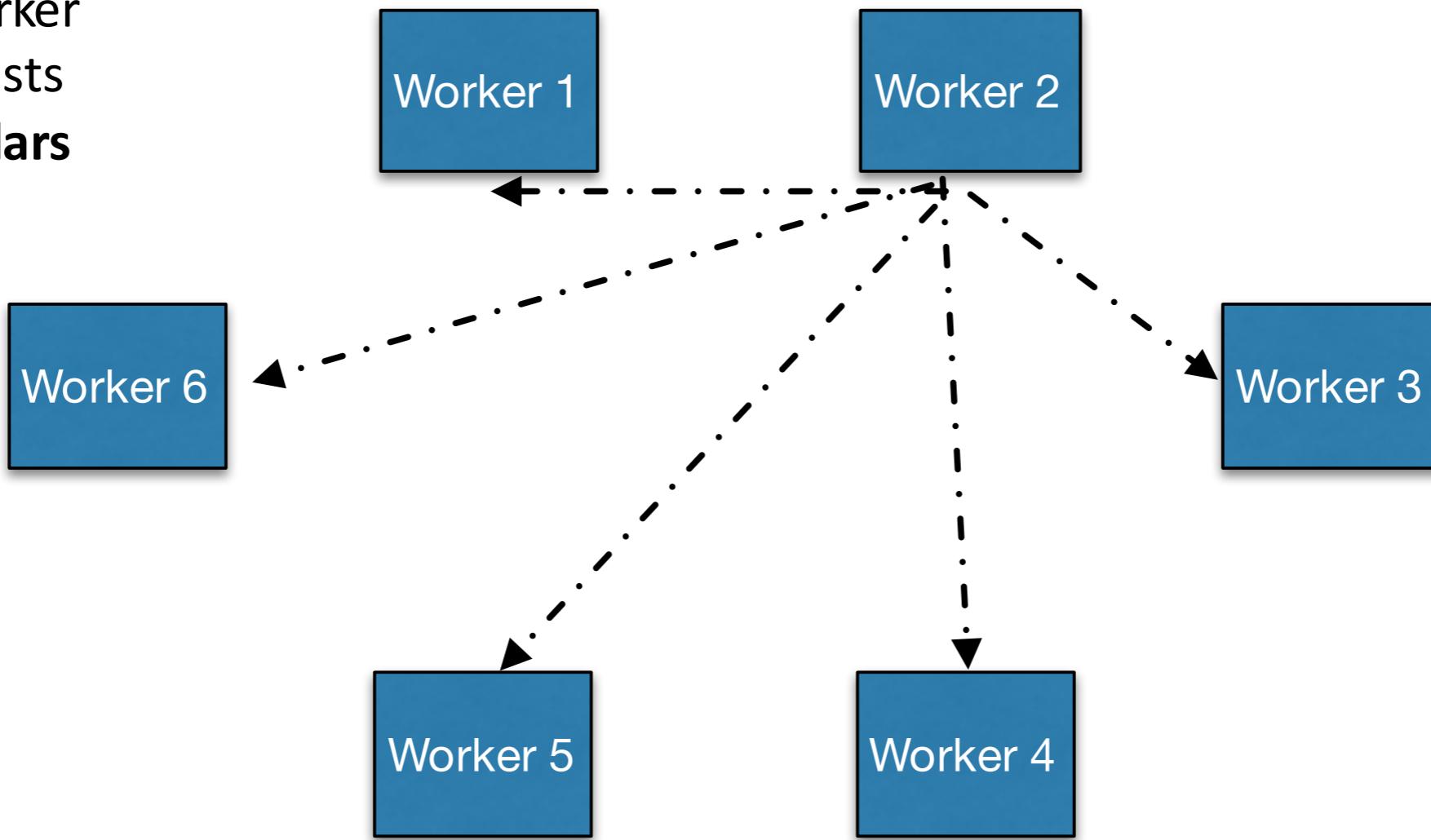
# Distributed Evolution

Each worker  
broadcasts  
**tiny scalars**



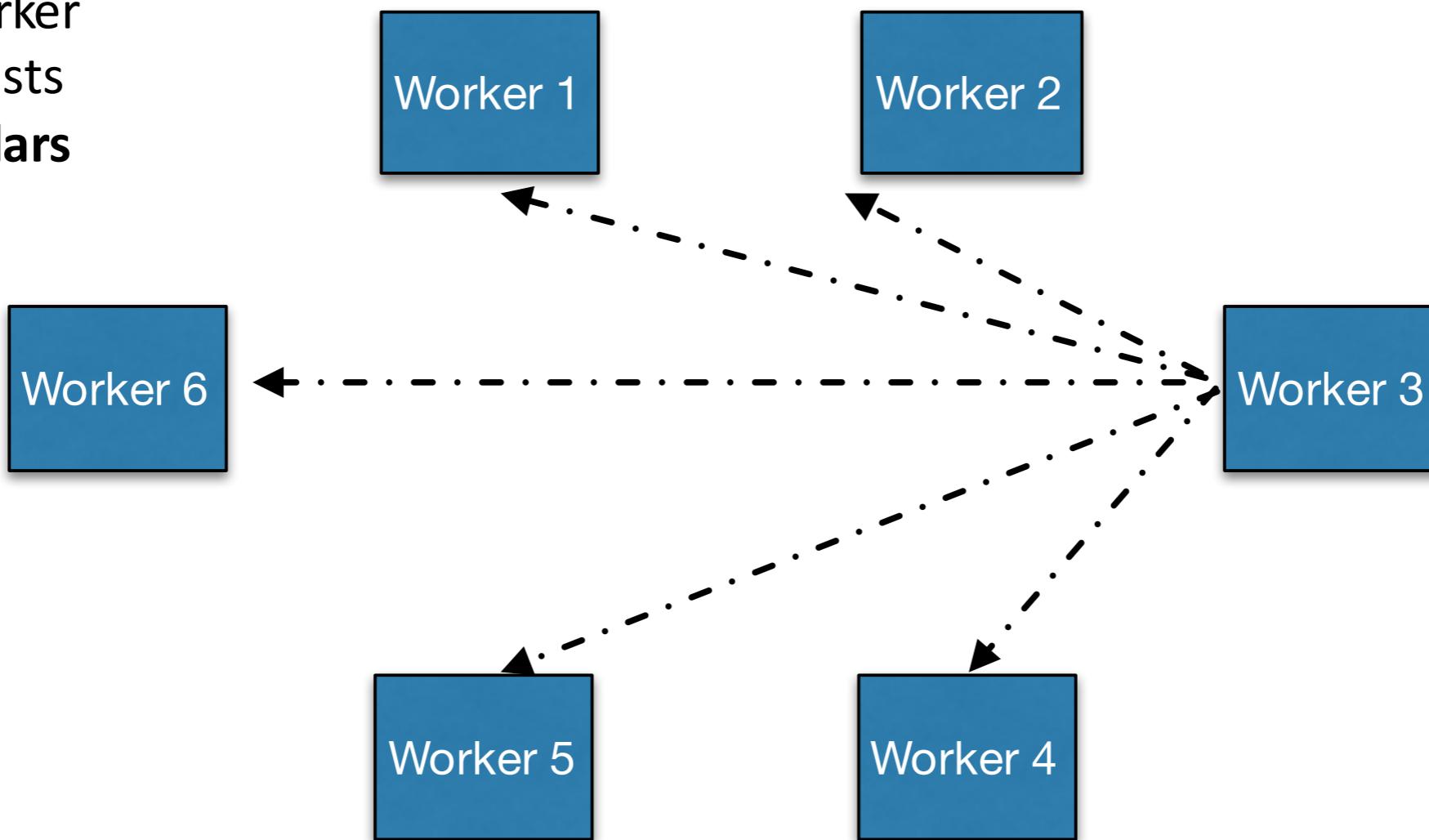
# Distributed Evolution

Each worker  
broadcasts  
**tiny scalars**

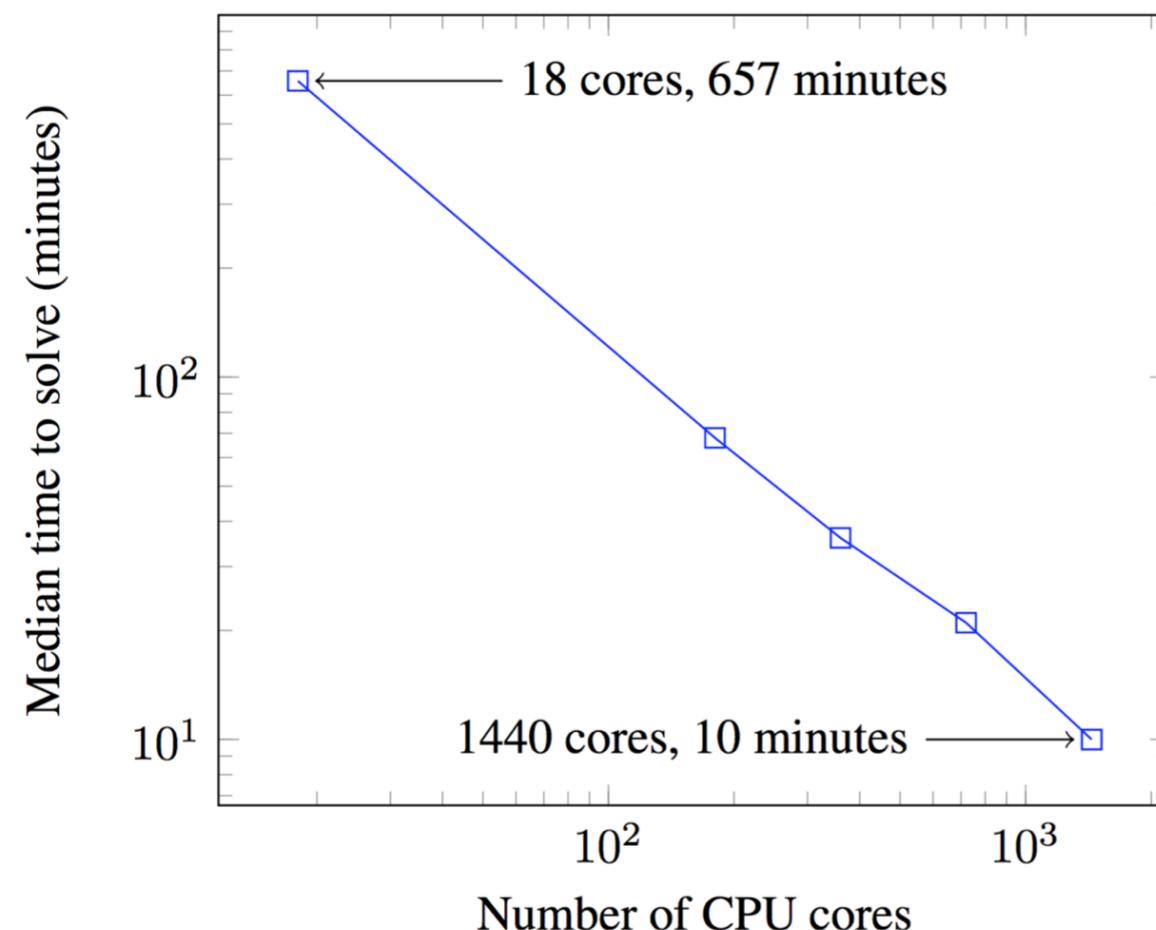


# Distributed Evolution

Each worker  
broadcasts  
**tiny scalars**



# Distributed Evolution Scales Very Well :-)



*Figure 1.* Time to reach a score of 6000 on 3D Humanoid with different number of CPU cores. Experiments are repeated 7 times and median time is reported.

# Limitations

Q: Can Evolutionary Search get stuck in local optima?

A: Yes, sure it can and it does.

Q: how can we get ES to be less stuck in local optima?

# ES can get stuck in local optima

- How can we help search methods search better?

It turns out, searching for policies:

- 1.in multiple related tasks
- 2.and in multiple related environments

can help ES search methods do better. We do not have any formal guarantees, but empirically it turns out this is the case.