# End-to-End Training of Deep Visuomotor Policies

Sergey Levine*, Chelsea Finn*, Trevor Darrell, Pieter Abbeel.
JMLR 17, 2016

(335 cites)

# Visuomotor?

visuomotor 🔊
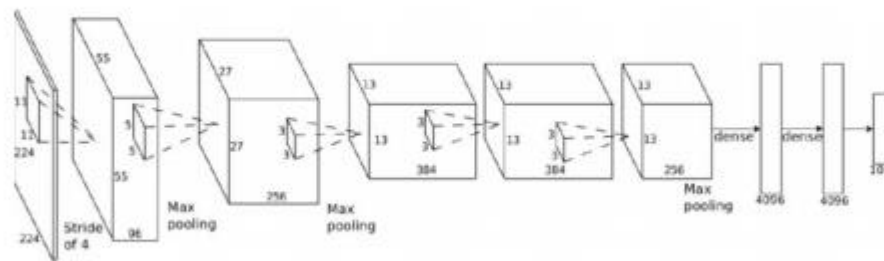
Also found in: Dictionary.

vis·u·o·mo·tor (viz'yū-ō-mō'tŏr),

Denoting the ability to synchronize visual information with physical movement, for example, driving a car or playing a video game of skill.
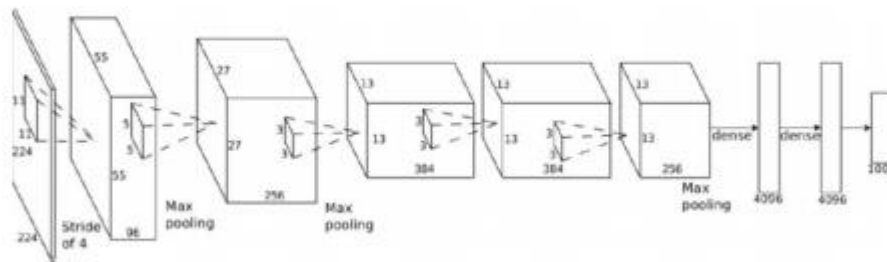
"CITE" 🔗 Farlex Partner Medical Dictionary © Farlex 2012

시각(카메라) 정보를 통해 신체(모터)를 움직이는 것이라고 봐도 무방할 듯

tiger
tiger
tiger cat
jaguar
lynx

perception
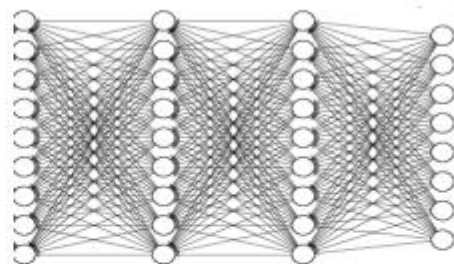
Action
(run away)

tiger
tiger
tiger cat
jaguar
lynx

action

# sensorimotor loop



Action
(run away)

no direct supervision
actions have consequences

## general-purpose neural network policy



$$\theta = \arg\min_\theta E_{\pi_\theta}\left[\sum_{t=1}^{T} c(\mathbf{x}_t, \mathbf{u}_t)\right]$$

$$\pi_\theta(\mathbf{u}_t|\mathbf{o}_t) - \text{control policy}$$

$\mathbf{o}_t$ − observation (may or may not be equal to $\mathbf{x}_t$)

# general-purpose neural network policy



$$\theta = \arg\min_\theta E_{\pi_\theta}\left[\sum_{t=1}^{T} c(\mathbf{x}_t, \mathbf{u}_t)\right]$$

$\pi_\theta(\mathbf{u}_t|\mathbf{o}_t)$ – control policy

$\mathbf{o}_t$ – observation (may or may not be equal to $\mathbf{x}_t$)
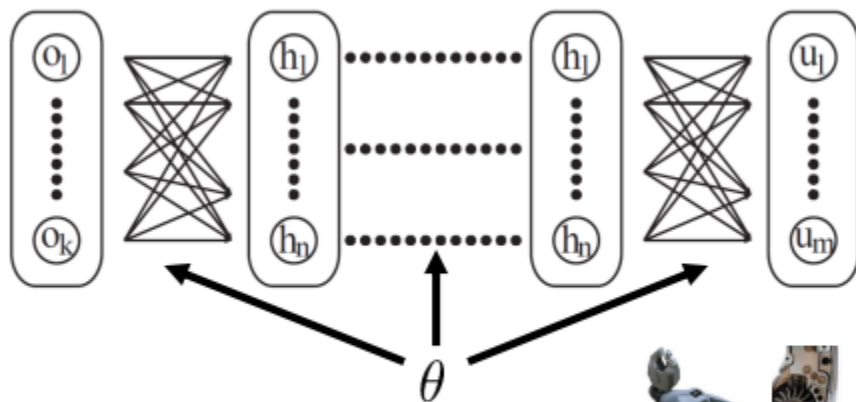
## general-purpose neural network policy



$$\theta = \arg\min_\theta E_{\pi_\theta}\left[\sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t)\right]$$

$\pi_\theta(\mathbf{u}_t|\mathbf{o}_t)$ – control policy

$\mathbf{o}_t$ – observation (may or may not be equal to $\mathbf{x}_t$)

## general-purpose neural network policy

$$\theta = \arg\min_\theta E_{\pi_\theta}\left[\sum_{t=1}^{T} c(\mathbf{x}_t, \mathbf{u}_t)\right]$$

$\pi_\theta(\mathbf{u}_t|\mathbf{o}_t)$ – control policy

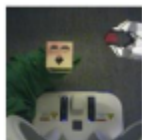$\mathbf{o}_t$ – observation (may or may not be equal to $\mathbf{x}_t$)
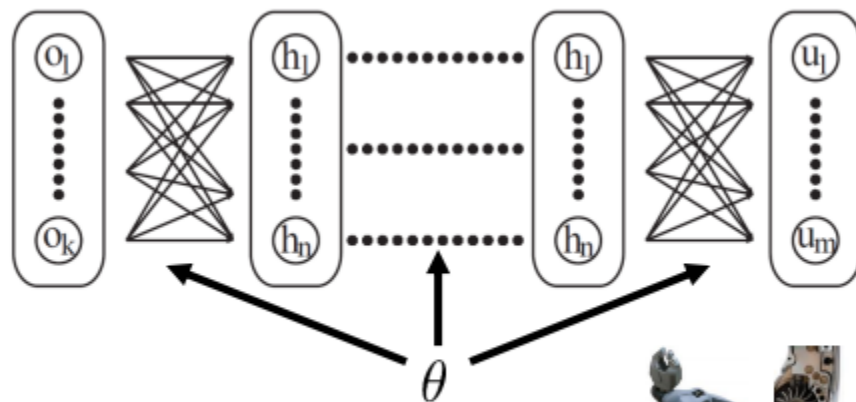
general-purpose neural network policy



$$\theta = \arg\min_\theta E_{\pi_\theta}\left[\sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t)\right]$$

$\pi_\theta(\mathbf{u}_t|\mathbf{o}_t)$ – control policy

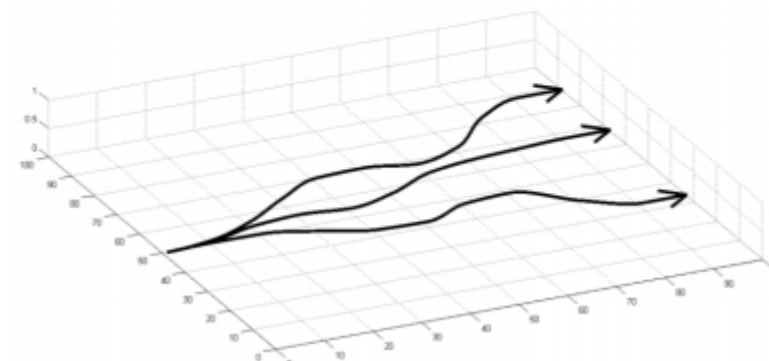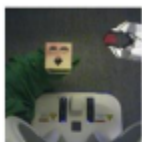$\mathbf{o}_t$ – observation (may or may not be equal to $\mathbf{x}_t$)
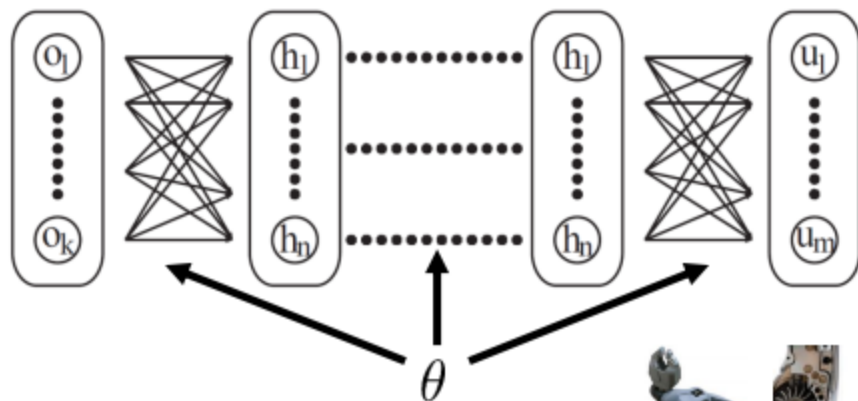
## general-purpose neural network policy



$$\theta = \arg\min_\theta E_{\pi_\theta}\left[\sum_{t=1}^{T} c(\mathbf{x}_t, \mathbf{u}_t)\right]$$

$\pi_\theta(\mathbf{u}_t|\mathbf{o}_t)$ – control policy

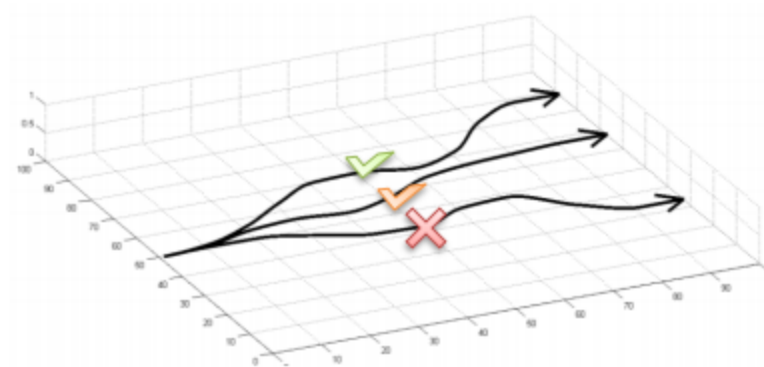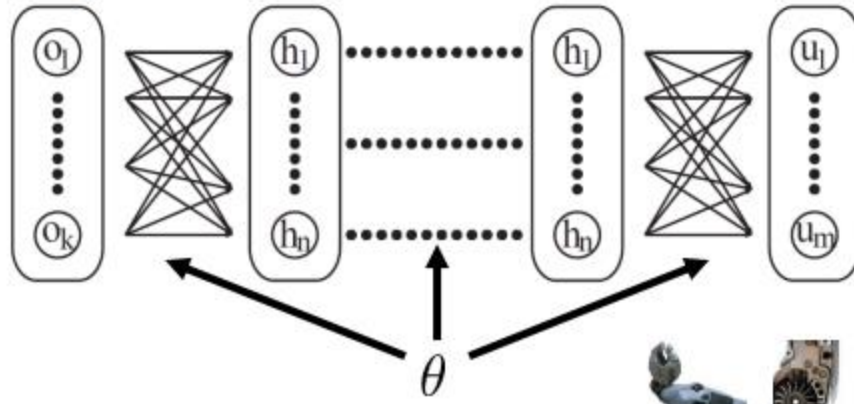$\mathbf{o}_t$ – observation (may or may not be equal to $\mathbf{x}_t$)

Mnih et al. '13     Schulman et al. '14 & '15

general-purpose neural network policy



$$\theta = \arg\min_\theta E_{\pi_\theta} \left[ \sum_{t=1}^{T} c(\mathbf{x}_t, \mathbf{u}_t) \right]$$

$\pi_\theta(\mathbf{u}_t | \mathbf{o}_t)$ – control policy

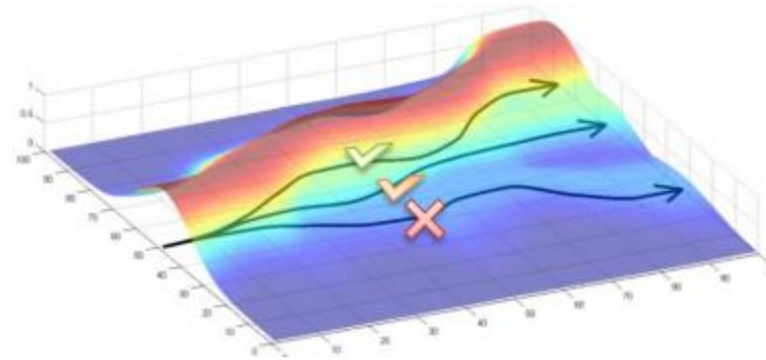$\mathbf{o}_t$ – observation (may or may not be equal to $\mathbf{x}_t$)

Mnih et al. '13    Schulman et al. '14 & '15

policy search (RL)    complex dynamics    complex policy    HARD
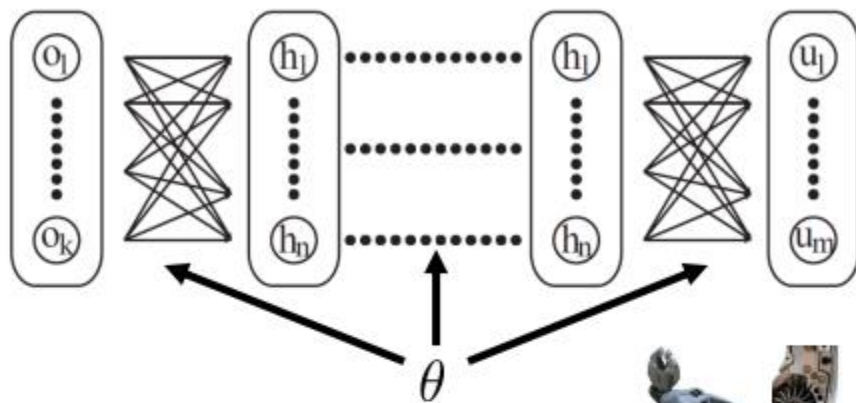
## general-purpose neural network policy

$$\theta = \arg\min_\theta E_{\pi_\theta}\left[\sum_{t=1}^T c(\mathbf{x}_t, \mathbf{u}_t)\right]$$

$\pi_\theta(\mathbf{u}_t|\mathbf{o}_t)$ – control policy

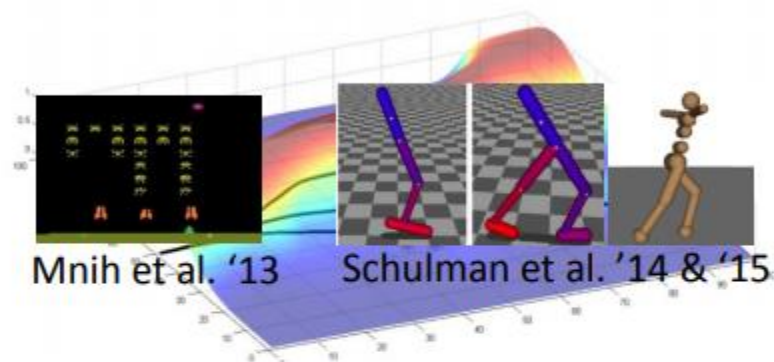$\mathbf{o}_t$ – observation (may or may not be equal to $\mathbf{x}_t$)

Mnih et al. '13    Schulman et al. '14 & '15

| | | | |
|---|---|---|---|
| policy search (RL) | complex dynamics | complex policy | HARD |
| supervised learning | ~~complex dynamics~~ | complex policy | EASY |

## general-purpose neural network policy



$$\theta = \arg\min_\theta E_{\pi_\theta}\left[\sum_{t=1}^{T} c(\mathbf{x}_t, \mathbf{u}_t)\right]$$

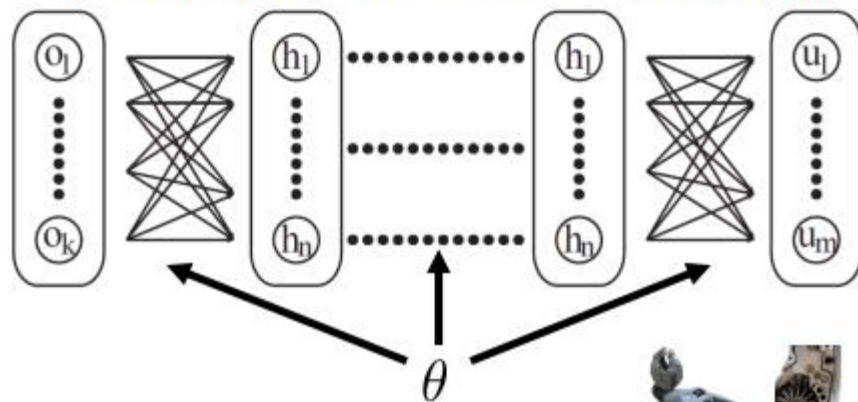$\pi_\theta(\mathbf{u}_t|\mathbf{o}_t)$ – control policy

$\mathbf{o}_t$ – observation (may or may not be equal to $\mathbf{x}_t$)

Mnih et al. '13     Schulman et al. '14 & '15
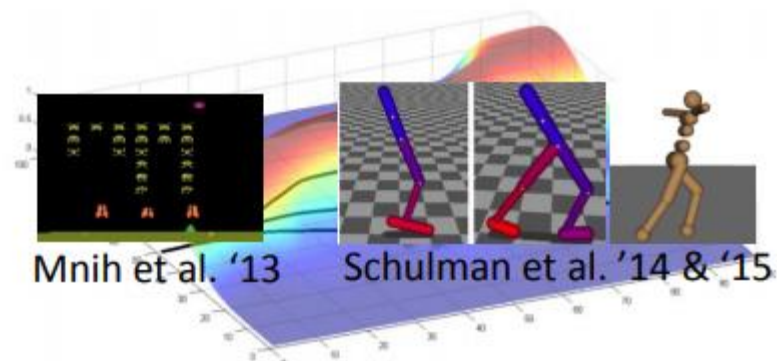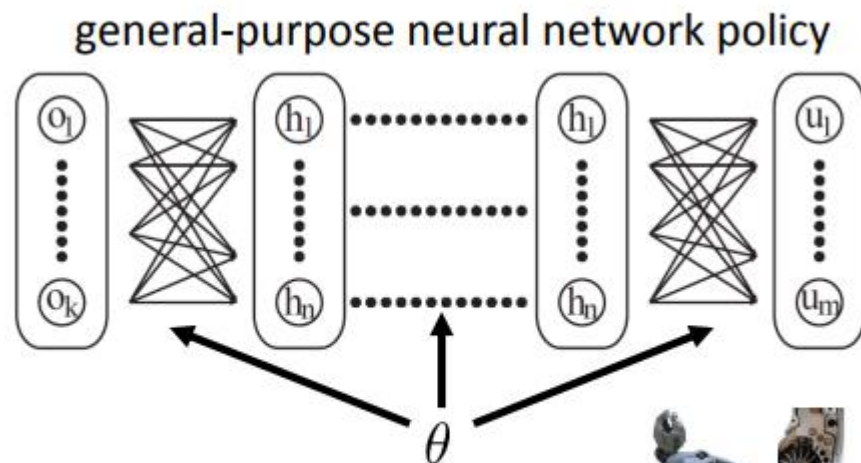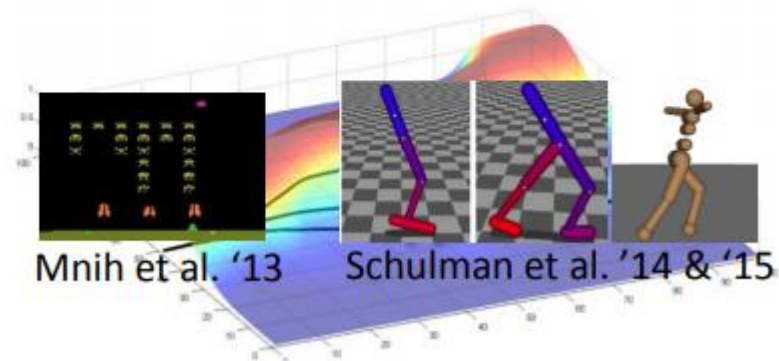
| | | | |
|---|---|---|---|
| policy search (RL) | complex dynamics | complex policy | HARD |
| supervised learning | ~~complex dynamics~~ | complex policy | EASY |
| optimal control | complex dynamics | ~~complex policy~~ | EASY |

1. break up the task:
   separately solve N
   different task instances

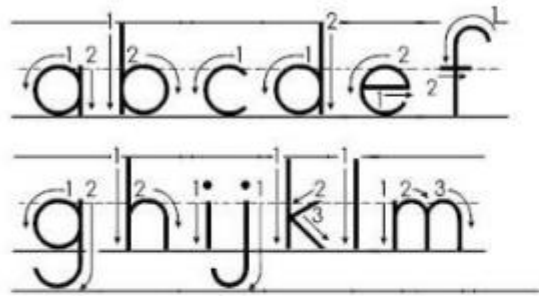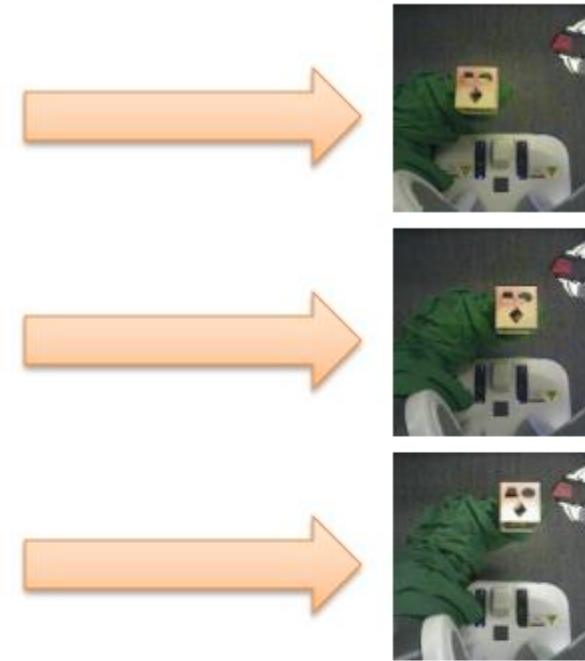1. break up the task: separately solve N different task instances

2. use supervised learning

trajectory-centric RL (fully observed)

state to action

supervised learning

observation to action

# Guided Policy Search



trajectory-centric RL       supervised learning

expectation under
current policy

$$\min_{\theta} E_{\pi_\theta}[c(\tau)]$$



$$\min_{\theta, p(\tau)} E_p[c(\tau)]$$

trajectory distribution(s)

$$s.t. \ \pi_\theta(\mathbf{u}_t | \mathbf{o}(\mathbf{x}_t)) = p(\mathbf{u}_t | \mathbf{x}_t) \ \forall t, \mathbf{x}_t, \mathbf{u}_t$$

$$\min_{\theta} E_{\pi_\theta}[c(\tau)]$$

$$\min_{\theta, p(\tau)} E_p[c(\tau)]$$

$$s.t. \quad \pi_\theta(\mathbf{u}_t | \mathbf{o}(\mathbf{x}_t)) = p(\mathbf{u}_t | \mathbf{x}_t) \quad \forall t, \mathbf{x}_t, \mathbf{u}_t$$

$$\min_{\theta} \overbrace{E_{\pi_\theta}[c(\tau)]}^{\text{expectation under current policy}}$$

$$\min_{\theta, p(\tau)} E_p[c(\tau)]$$

$$s.t. \quad \pi_\theta(\mathbf{u}_t|\mathbf{o}(\mathbf{x}_t)) = p(\mathbf{u}_t|\mathbf{x}_t) \quad \forall t, \mathbf{x}_t, \mathbf{u}_t$$
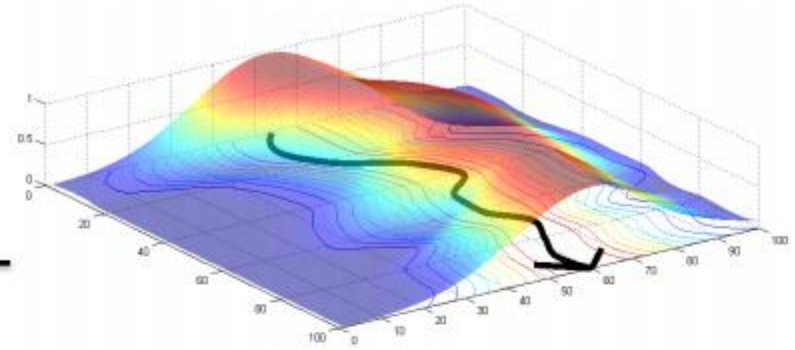
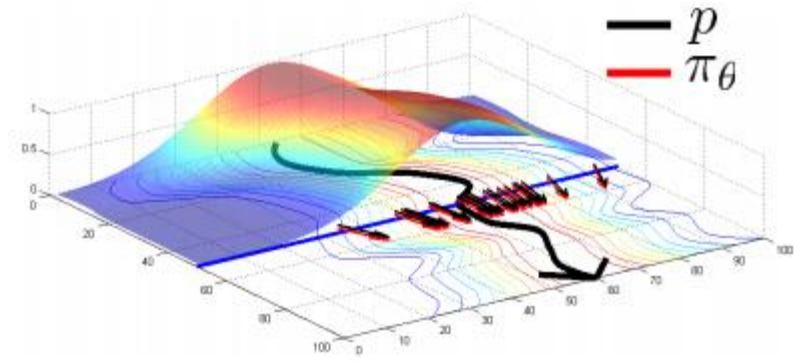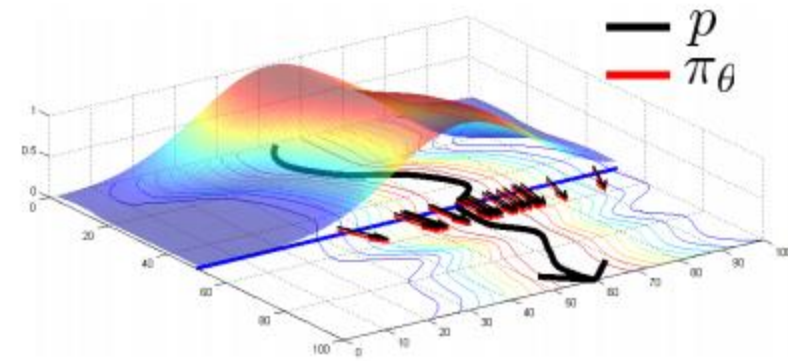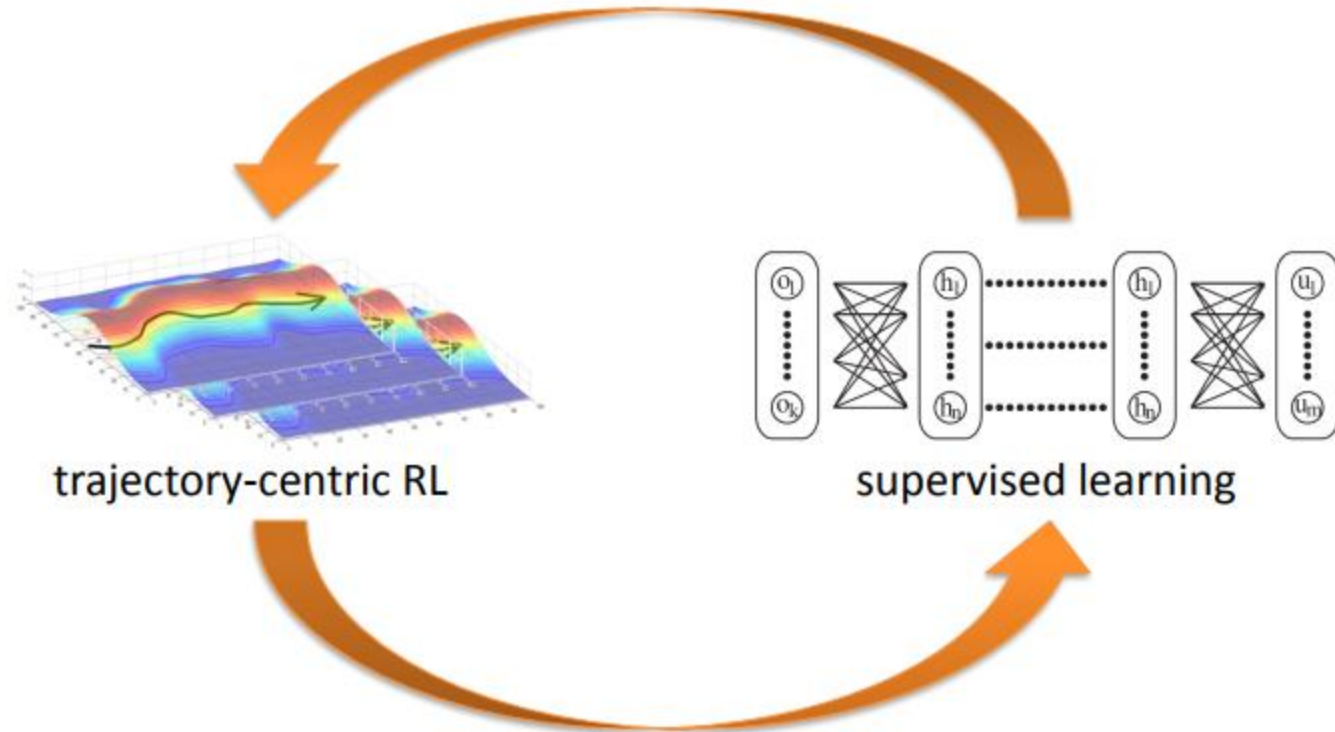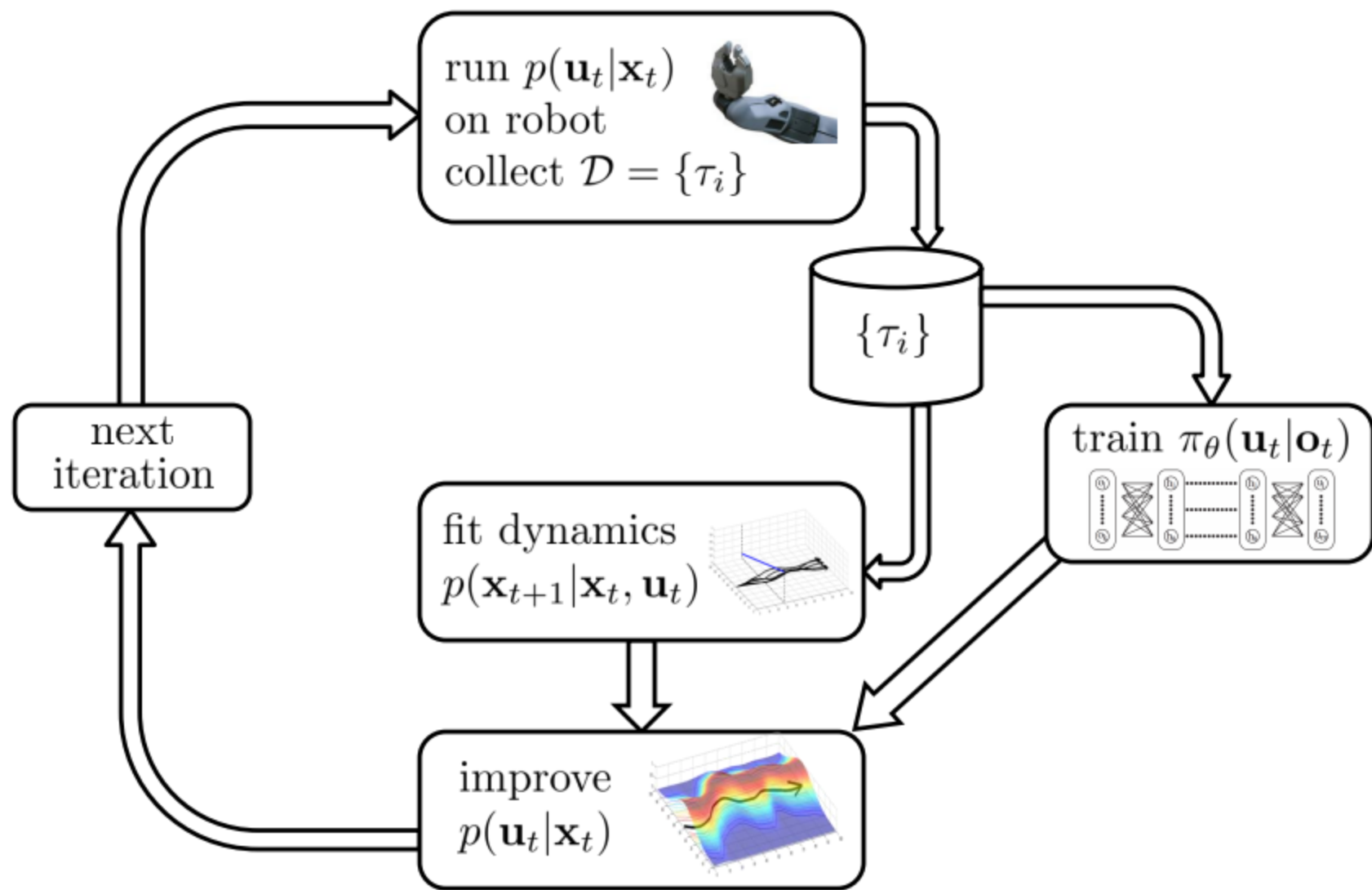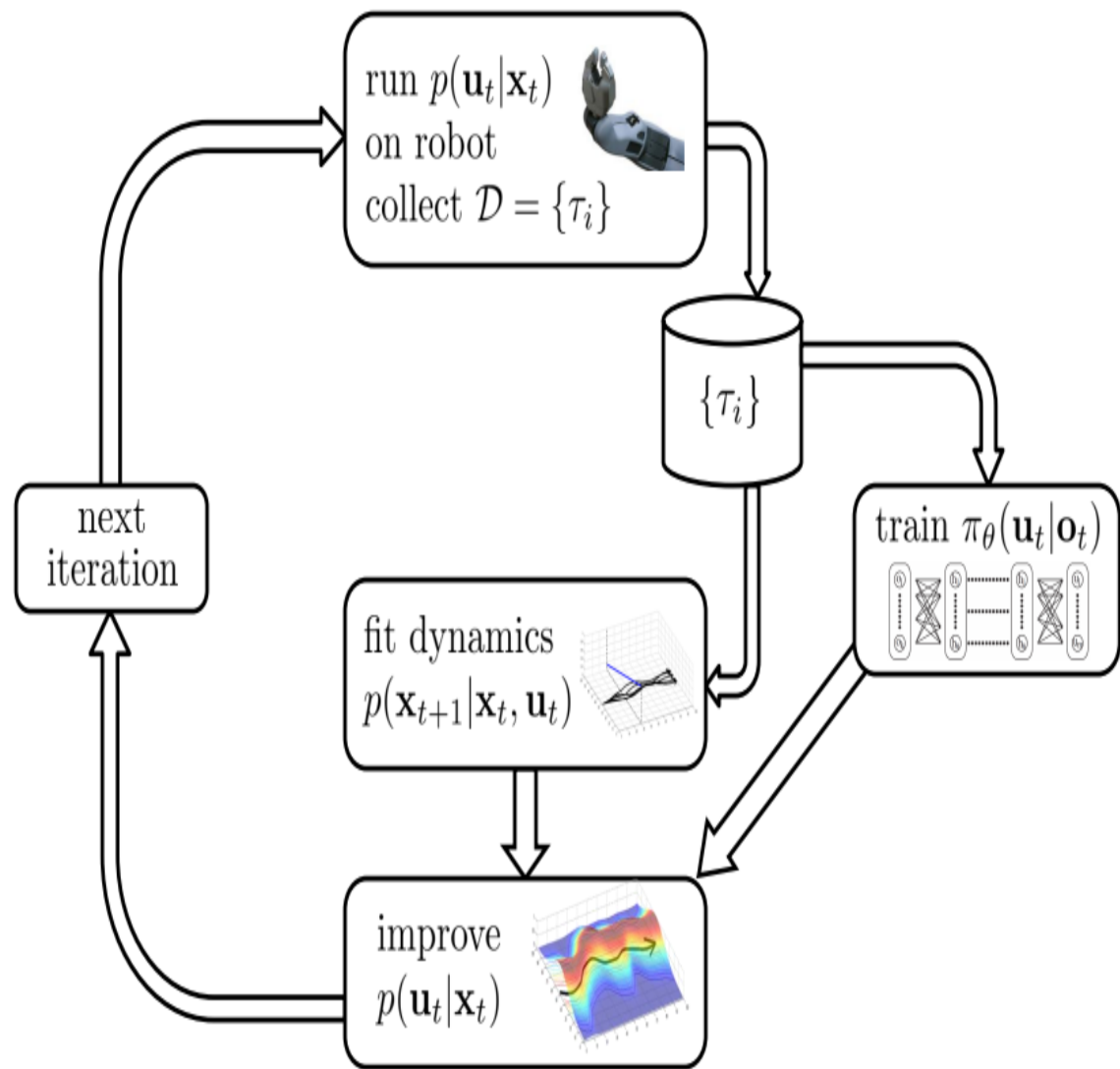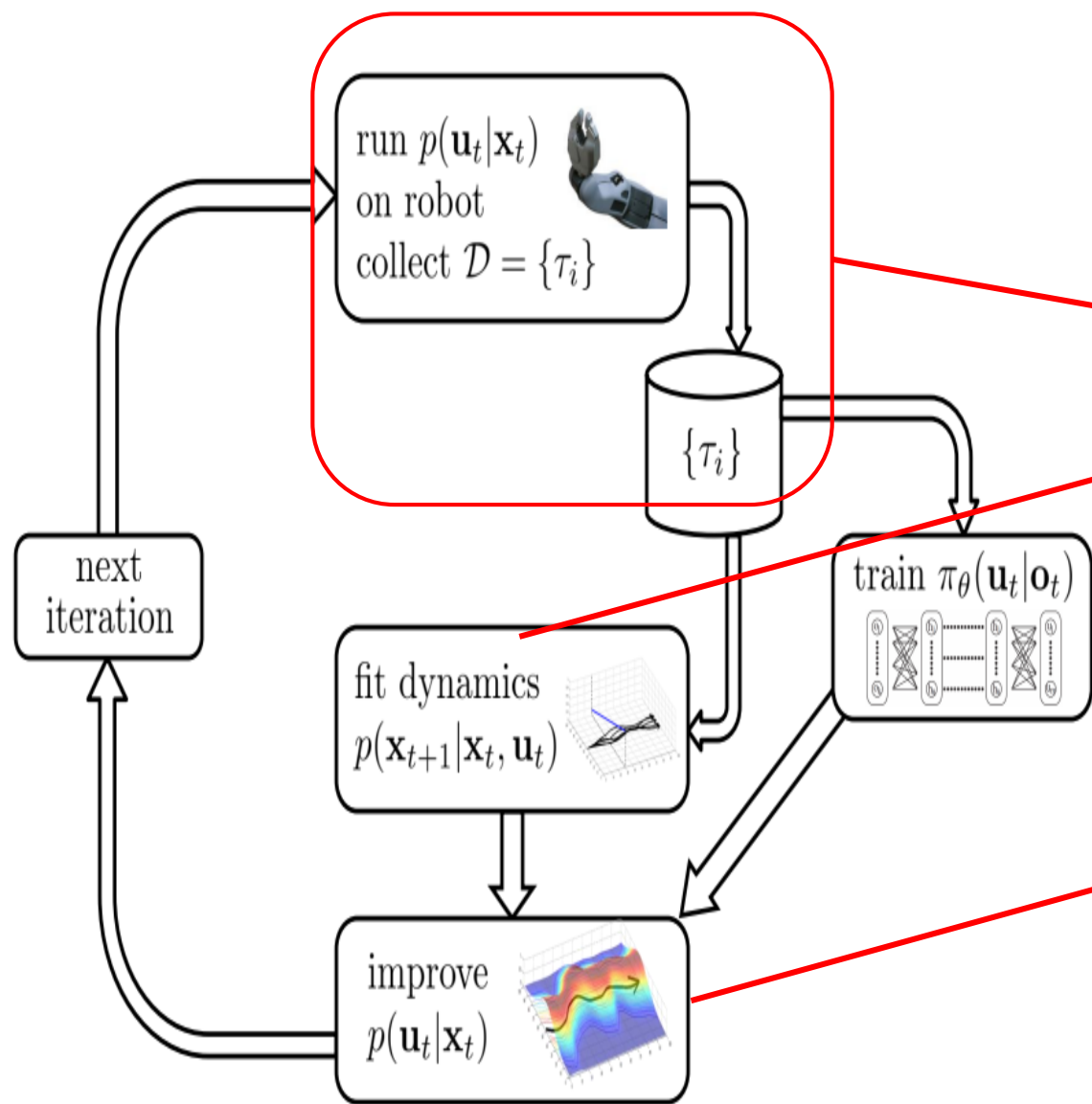solve using Bregman ADMM (BADMM), a type of dual decomposition method

trajectory-centric RL

supervised learning

2017-09-06 jyyang

**Algorithm 1** Guided Policy Search

1: Generate DDP solutions $\pi_{\mathcal{G}_1}, \ldots, \pi_{\mathcal{G}_n}$
2: Sample $\zeta_1, \ldots, \zeta_m$ from $q(\zeta) = \frac{1}{n} \sum_i \pi_{\mathcal{G}_i}(\zeta)$
3: Initialize $\theta^\star \leftarrow \arg\max_\theta \sum_i \log \pi_{\theta^\star}(\zeta_i)$
4: Build initial sample set $\mathcal{S}$ from $\pi_{\mathcal{G}_1}, \ldots, \pi_{\mathcal{G}_n}, \pi_{\theta^\star}$
5: **for** iteration $k = 1$ to $K$ **do**
6:     Choose current sample set $\mathcal{S}_k \subset \mathcal{S}$
7:     Optimize $\theta_k \leftarrow \arg\max_\theta \Phi_{\mathcal{S}_k}(\theta)$
8:     Append samples from $\pi_{\theta_k}$ to $\mathcal{S}_k$ and $\mathcal{S}$
9:     Optionally generate adaptive guiding samples
10:     Estimate the values of $\pi_{\theta_k}$ and $\pi_{\theta^\star}$ using $\mathcal{S}_k$
11:     **if** $\pi_{\theta_k}$ is better than $\pi_{\theta^\star}$ **then**
12:         Set $\theta^\star \leftarrow \theta_k$
13:         Decrease $w_r$
14:     **else**
15:         Increase $w_r$
16:         Optionally, resample from $\pi_{\theta^\star}$
17:     **end if**
18: **end for**
19: Return the best policy $\pi_{\theta^\star}$
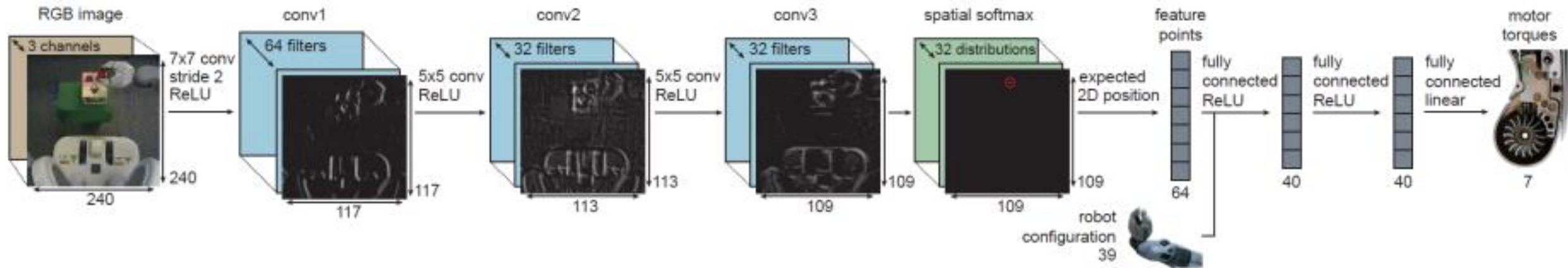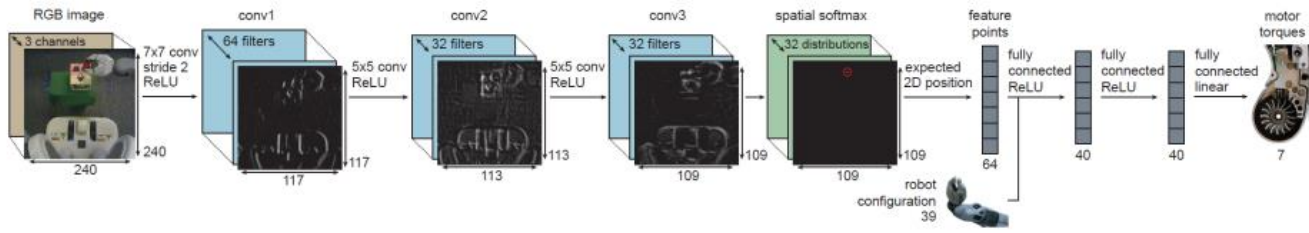
2017-09-06 jyyang

**Algorithm 1** Guided Policy Search

1: Generate DDP solutions $\pi_{\mathcal{G}_1}, \ldots, \pi_{\mathcal{G}_n}$
2: Sample $\zeta_1, \ldots, \zeta_m$ from $q(\zeta) = \frac{1}{n} \sum_i \pi_{\mathcal{G}_i}(\zeta)$
3: Initialize $\theta^\star \leftarrow \arg\max_\theta \sum_i \log \pi_{\theta^\star}(\zeta_i)$
4: Build initial sample set $\mathcal{S}$ from $\pi_{\mathcal{G}_1}, \ldots, \pi_{\mathcal{G}_n}, \pi_{\theta^\star}$
5: **for** iteration $k = 1$ to $K$ **do**
6:    Choose current sample set $\mathcal{S}_k \subset \mathcal{S}$
7:    Optimize $\theta_k \leftarrow \arg\max_\theta \Phi_{\mathcal{S}_k}(\theta)$
8:    Append samples from $\pi_{\theta_k}$ to $\mathcal{S}_k$ and $\mathcal{S}$
9:    Optionally generate adaptive guiding samples
10:    Estimate the values of $\pi_{\theta_k}$ and $\pi_{\theta^\star}$ using $\mathcal{S}_k$
11:    **if** $\pi_{\theta_k}$ is better than $\pi_{\theta^\star}$ **then**
12:       Set $\theta^\star \leftarrow \theta_k$
13:       Decrease $w_r$
14:    **else**
15:       Increase $w_r$
16:       Optionally, resample from $\pi_{\theta^\star}$
17:    **end if**
18: **end for**
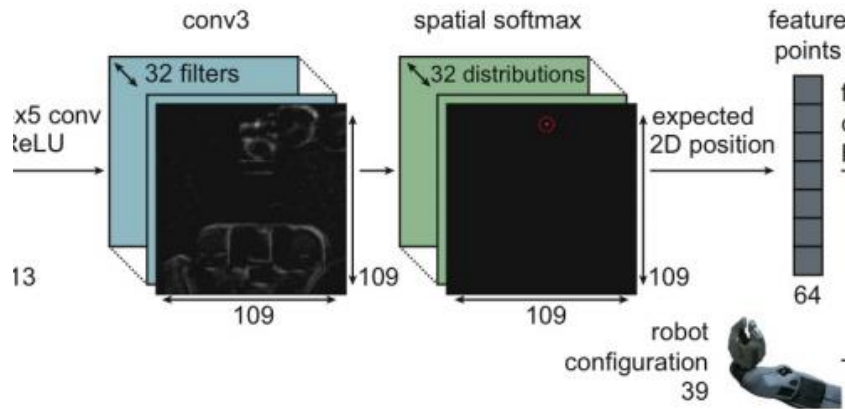19: Return the best policy $\pi_{\theta^\star}$

2017-09-06 jyyang

# Network architecture

# spatial softmax ~~라고 읽고 spatial soft=argmax라 부른다~~

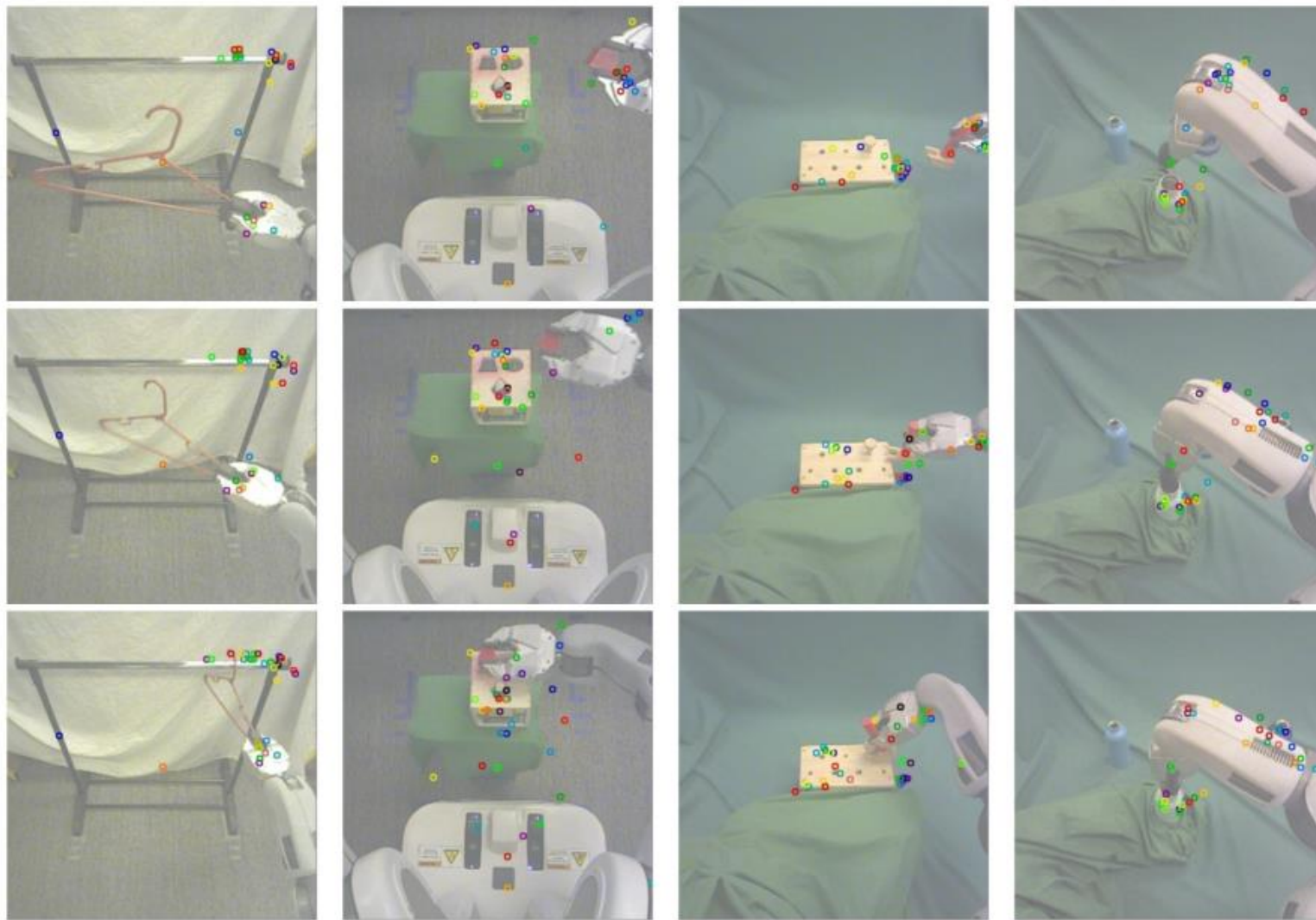1. 채널별로 softmax를 구함
2. 채널별로 구한 softmax 결과와 x, y 좌표 값을 곱한 후 더한다.



$$s_{cij} = e^{a_{cij}} / \sum_{i'j'} e^{a_{ci'j'}}$$

$$f_{cx} = \sum_{ij} s_{cij} x_{ij} \text{ and } f_{cy} = \sum_{ij} s_{cij} y_{ij}$$

$$\mathrm{softargmax}(x) = \sum_i \frac{e^{\beta x_i}}{\sum_j e^{\beta x_j}} i$$

(a) hanger   (b) cube   (c) hammer   (d) bottle

Figure 10: Feature points tracked by the policy during task execution for each of the four tasks. Each feature point is displayed in a different random color, with consistent coloring across images. The policy finds features on the target object and the robot gripper and arm. In the bottle cap task, note that the policy correctly ignores the distractor bottle in the background, even though it was not present during training.

Position 2

real time                              autonomous execution

2017-09-06  jyyang

# Training

- SL + RL
- Cost function