# Policy Learning
# Using Weak Supervision

Jingkang Wang*[1,2], Hongyi Guo*[3], Zhaowei Zhu*[4], Yang Liu[4]

# Deep Learning in Sequential Decision Making
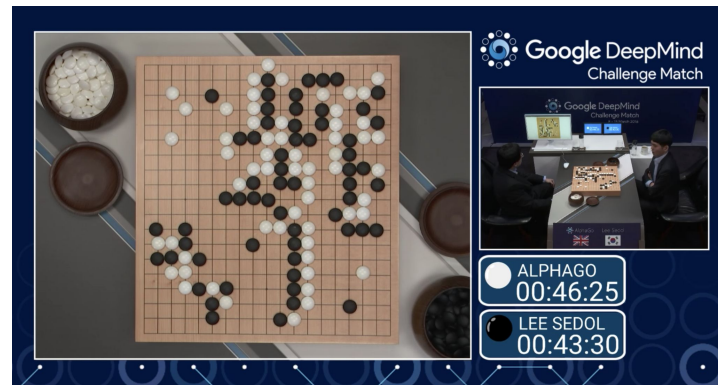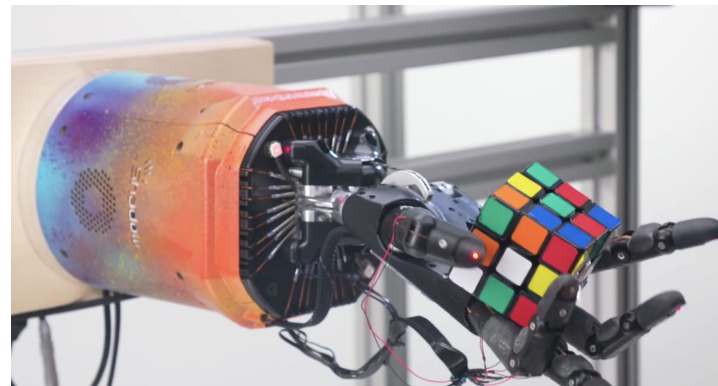

Atari2600 Games [Mnih et al., 2015]


AlphaGo [Silver et al., 2017]


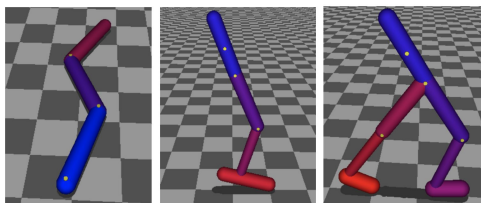Self-Driving [Amini et al., 2020]


Self-Driving [OpenAI, 2019]

2

# Markov Decision Process (MDP)



[Mnih et al., 2013]

[Schulman et al., 2015]

[AlphaGo versus Lee Sedol]

MDP: $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{P}, \gamma \rangle$

Action $a_t \in \mathcal{A}$

Agent

$\pi : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$

Environment

# Markov Decision Process (MDP)


[Mnih et al., 2013]


[Schulman et al., 2015]


[AlphaGo versus Lee Sedol]

MDP: $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{P}, \gamma \rangle$

Action $a_t \in \mathcal{A}$

Agent

$\pi : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$

Environment

State $s_{t+1} \in \mathcal{S}$
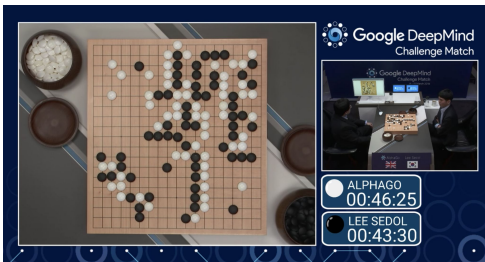Reward $r(s_t, a_t) \in \mathcal{R}$

4

# Markov Decision Process (MDP)

[Mnih et al., 2013]
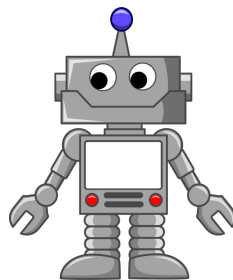
[Schulman et al., 2015]

[AlphaGo versus Lee Sedol]

MDP: $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{R}, \mathcal{P}, \gamma \rangle$
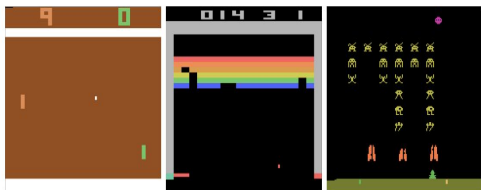
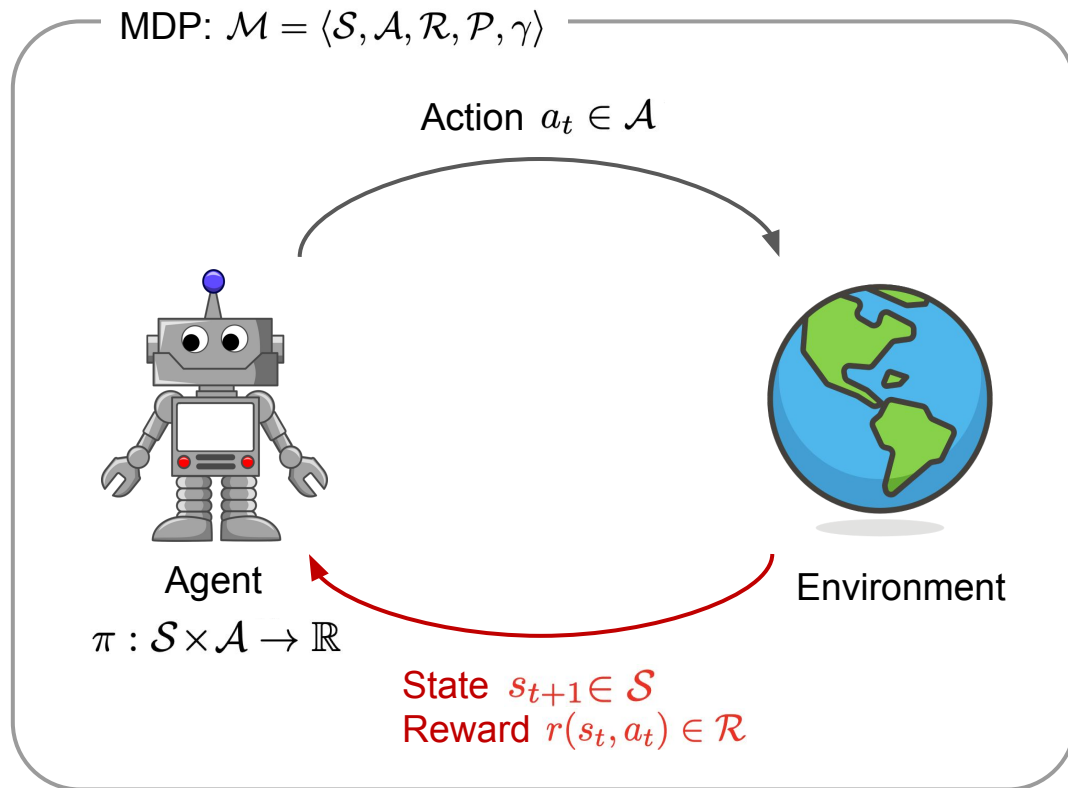Action $a_t \in \mathcal{A}$

**The agent interacts with environment repeatedly**

Agent

$\pi : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$

Environment

State $s_{t+1} \in \mathcal{S}$
Reward $r(s_t, a_t) \in \mathcal{R}$

# Reinforcement Learning

Action $a_t \in \mathcal{A}$

The agent interacts with environment repeatedly

Agent

$\pi : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$

Environment

State $s_t \in \mathcal{S}$
Reward $r(s_t, a_t) \in \mathcal{R}$

Generated trajectory:

$$\tau = \{(s_t, a_t, r_t)\}_{t=0}^{T}$$

**Objective**: maximize the expected reward

$$J(\pi) = \mathbb{E}_{(s_t, a_t, r_t) \sim \tau}[\sum_{t=0}^{T} \gamma^t r_t]$$

**Supervision signal**: *reward*

# Imitation Learning

Action

The agent interacts with expert repeatedly to mimic the policy

Expert $\pi_E$

State

Agent

$$\pi : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$$

Expert demonstrations

$$D_E = \{(s_i, a_i)\}_{i=1}^N$$

**Behavioral Cloning Objective**: maximize the log-likelihood

$$J(\pi) = \mathbb{E}_{(s,a) \sim \mathcal{D}_E} \left[ \log \pi(a|s) \right]$$

**Supervision signal**: *expert action*

**Supervision signal**: *reward + expert action*

**Hybrid objective:** $J(\pi) = \lambda_1 \mathbb{E}_{(s_t, a_t, r_t) \sim \tau} \left[ \sum_{t=0}^{T} \gamma^t r_t \right] + \lambda_2 \mathbb{E}_{(s,a) \sim \mathcal{D}_E} \left[ \log \pi(a|s) \right]$



Action

Action $a_t \in \mathcal{A}$

**The agent interacts with expert repeatedly to mimic the policy**

**The agent interacts with environment repeatedly**

State $s_t \in \mathcal{S}$
Reward $r(s_t, a_t) \in \mathcal{R}$

Expert $\pi_E$

State

Agent
$\pi : \mathcal{S} \times \mathcal{A} \to \mathbb{R}$

Environment

# Policy Learning

Summary:



Teacher

action

state

Behavioral Cloning (BC)

Environment

action

state, reward

Reinforcement Learning (RL)

BC    RL    + Hybrid Learning

# Weak supervision signals are everywhere!

**Weak Supervision:**

- **RL:** The reward may be collected through sensors thus noisy
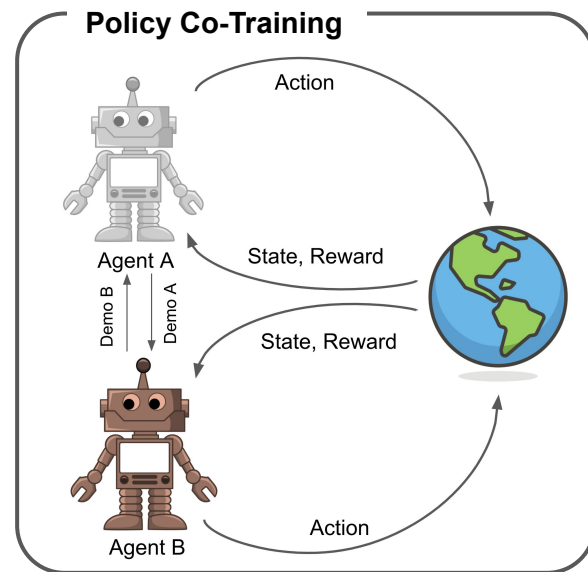
- **IL:** The demonstrations by an expert are often imperfect due to limited resources

# Weakly Supervised Policy Learning

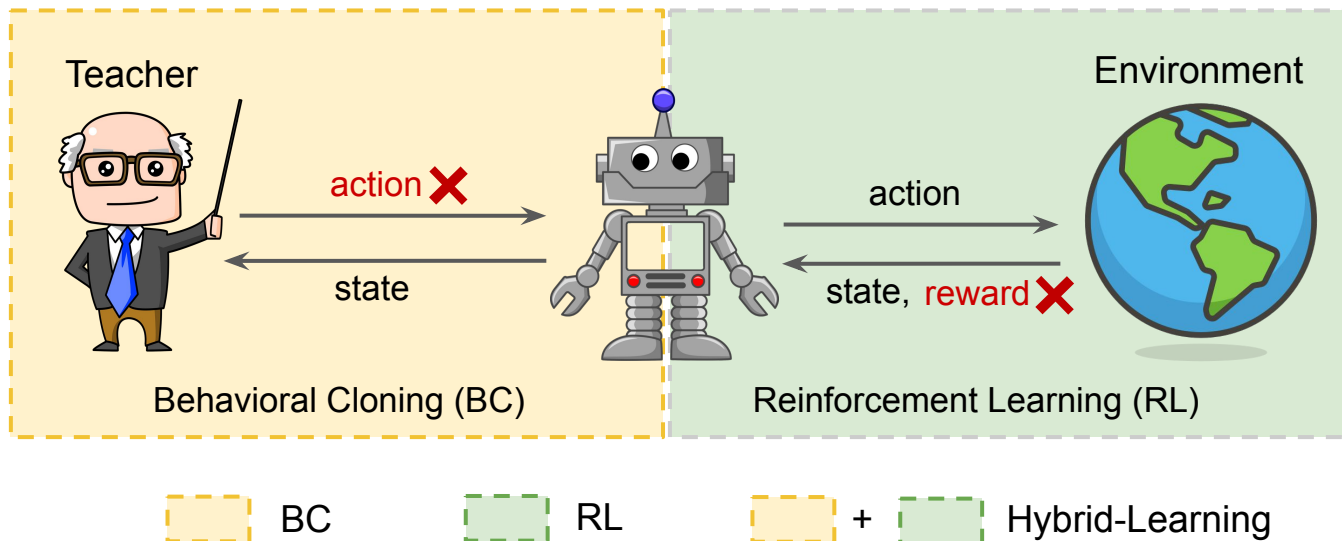Problem: Supervision signals $\tilde{Y}$ (either reward or expert's demonstrations) are *not credible!*



- **RL:** The reward may be collected through sensors thus noisy

- **IL:** The demonstrations by an expert are often imperfect due to limited resources

*Weak Supervision:*

# Weakly Supervised Policy Learning

Weakly Supervised Policy Learning $\{(s_i, a_i), \widetilde{Y}_i\}_{i=1}^{N}$



Teacher

action ✖

state

Behavioral Cloning (BC)

Environment

action

state, reward ✖

Reinforcement Learning (RL)

☐ BC   ☐ RL   ☐ + ☐ Hybrid-Learning

- Objective: $J(\pi) = \mathbb{E}_{(s,a) \sim \tau} \left[ \mathtt{Eva}_\pi \left( (s, a), \tilde{Y} \right) \right]$

# Correlated Agreement (CA)

Solution - CA with weak supervision:   $\text{Eva}_\pi\big((s_i, a_i), \widetilde{Y}_i\big) - \text{Eva}_\pi\big((s_j, a_j), \widetilde{Y}_k\big)$

# Correlated Agreement (CA)

Solution - CA with weak supervision: $\text{Eva}_\pi\left((s_i, a_i), \widetilde{Y}_i\right) - \text{Eva}_\pi\left((s_j, a_j), \widetilde{Y}_k\right)$
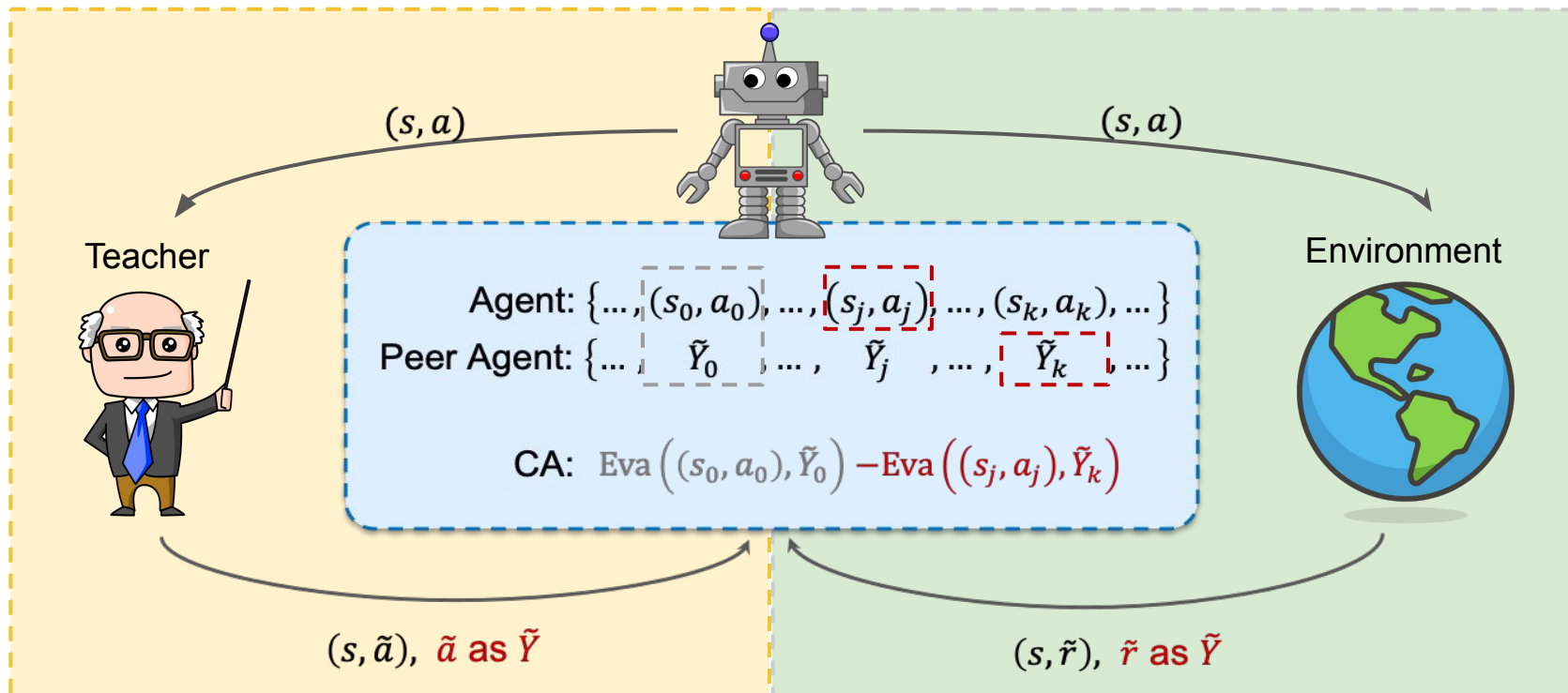


$(s, a)$

$(s, a)$

Teacher

Environment

Prediction: $\{a_1 = a_2 = a_3 = a_4 = 1, \ a_5 = 0\}$

Supervision: $\{\tilde{a}_1 = \tilde{a}_2 = \tilde{a}_3 = \tilde{a}_{4.} = 1, \ \tilde{a}_5 = 0\}$

CA: $1 - 0.75^2 - 0.25^2 = 0.375$

$(s, \tilde{a}), \ \tilde{a}$ as $\widetilde{Y}$

$(s, \tilde{r}), \ \tilde{r}$ as $\widetilde{Y}$

# PeerPL: A Unified Framework for Weakly Supervised PL

Solution - CA with weak supervision:   $\text{Eva}_\pi\big((s_i, a_i), \widetilde{Y}_i\big) - \text{Eva}_\pi\big((s_j, a_j), \widetilde{Y}_k\big)$

# PeerRL

We assume the noisy reward $\tilde{r}$ is generated following a certain function $F : \mathcal{R} \to \tilde{\mathcal{R}}$.

- Discrete with $|\mathcal{R}|$ levels.
- Characterized via an ***unknown*** matrix $\mathbf{C}^{\mathrm{RL}}_{|\mathcal{R}| \times |\mathcal{R}|}$

$$r(s,a) \xrightarrow{\mathbf{C}^{\mathrm{RL}}_{|\mathcal{R}| \times |\mathcal{R}|}} \tilde{r}(s,a)$$

**PeerRL** handles the noisy reward by defining the peer RL reward:

$$\tilde{r}_{\mathrm{peer}}(s,a) = \tilde{r}(s,a) - \xi \cdot \tilde{r}'$$

where $\tilde{r}' \overset{\pi_{\mathrm{sample}}}{\sim} \{\tilde{r}(s,a) | s \in \mathcal{S}, a \in \mathcal{A}\}$ is a reward sampled over all state-action pairs according to a fixed policy $\pi_{\mathrm{sample}}$.

$$\tilde{r}(s,a) \xrightarrow{\mathrm{CA}} \tilde{r}_{\mathrm{peer}}(s,a)$$

Our theory shows that peer RL rewards are robust to noisy rewards *(see Lemma 1 and Theorem 1)*.

16

# Why Peer Reward Works?

- **Hypothesis 1:** PeerRL reduces the bias (while with larger variance like Wang et al., 2020)

noisy reward:
$$\mathbb{E}[\tilde{r}] = \eta \cdot \left( \mathbb{E}[r] + \frac{e_+}{1 - e_- - e_+} r_- + \frac{e_-}{1 - e_- - e_+} r_+ \right)$$

peer reward:
$$\mathbb{E}[\tilde{r}_{\mathrm{peer}}] = \eta \cdot (\mathbb{E}[r] - (1 - p_{\mathrm{peer}}) r_- - p_{\mathrm{peer}} r_+)$$

potentially much larger than $(1 - p_{\mathrm{peer}})$ and $p_{\mathrm{peer}}$ in high noise regime!

- **Hypothesis 2:** PeerRL helps break ties
  - "tie" states indicate that the rewards for different states are the same - unstable and uncertain
  - randomness in discretization model thus breaking ties - more informative for optimization

2-state Markov process (no actions)

$s_1$    $s_2$

$r_1 \sim \mathrm{clamp}[\mathcal{N}(0.6, 1), \min = 0, \max = 1]$

$r_2 \sim \mathrm{clamp}[\mathcal{N}(0.4, 1), \min = 0, \max = 1]$

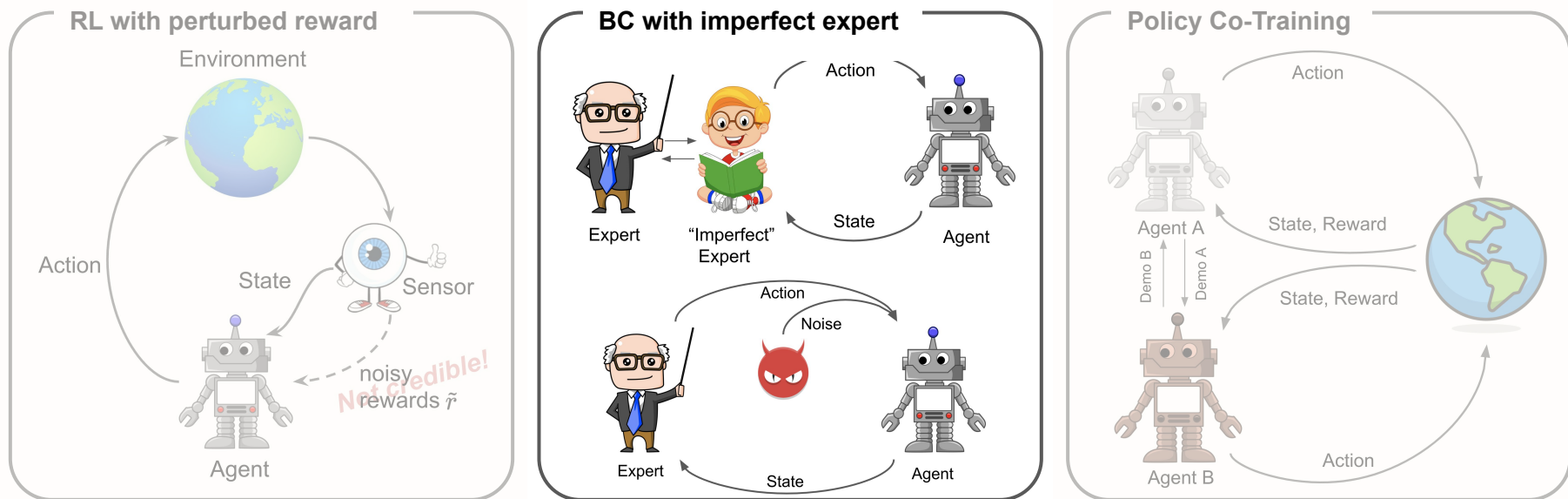| | Correct | Tie | Incorrect |
|---|---|---|---|
| Baseline | 54.6% | 5.6% | 39.8% |
| PeerRL | **58.0%** | **0.3%** | **41.7%** |

**Tie breaking!**

# PeerPL: A Unified Framework for Weakly Supervised PL

Solution - CA with weak supervision: $\mathtt{Eva}_\pi\big((s_i, a_i), \widetilde{Y}_i\big) - \mathtt{Eva}_\pi\big((s_j, a_j), \widetilde{Y}_k\big)$

*PeerBC*

# PeerBC

Available weak demonstrations $\{(s_i, \tilde{a}_i)\}_{i=1}^{N}$ where $\tilde{a}_i \sim \tilde{\pi}_E(\cdot|s_i)$

- The noisy action $\tilde{a}_i$ is independent of the state given the deterministic expert action $\pi_E(s)$
- The noise is characterized by an ***unknown*** confusion matrix $\mathbf{C}_{|\mathcal{A}| \times |\mathcal{A}|}^{BC}$

Again, we use CA with weak supervision to handle the noise

$$a_i \xrightarrow{\mathbf{C}_{|\mathcal{A}| \times |\mathcal{A}|}^{BC}} \tilde{a}_i$$

- Taking cross-entropy loss for example
-

$$\mathsf{Eva}_\pi^{BC}\big((s_i, a_i), \tilde{a}_i\big) \xrightarrow{\text{CA}} J^{BC}(\pi_\theta)$$

$$J^{BC}(\pi_\theta) = \mathbb{E}\Big[\mathsf{Eva}_\pi^{BC}\big((s_i, a_i), \tilde{a}_i\big)\Big] - \xi \cdot \mathbb{E}\Big[\mathsf{Eva}_\pi^{BC}\big((s_j, a_j), \tilde{a}_k\big)\Big]$$

where $\mathsf{Eva}_\pi^{BC}(s, a), \tilde{a}) = -\ell\big(\pi_\theta, (s, \tilde{a})\big) = \log \pi_\theta(\tilde{a}|s)$.

- sufficient amount of weak demonstrations

==*(see Theorem 2)*==

# PeerPL: A Unified Framework for Weakly Supervised PL

Solution - CA with weak supervision: $\mathrm{Eva}_\pi\big((s_i, a_i), \widetilde{Y}_i\big) - \mathrm{Eva}_\pi\big((s_j, a_j), \widetilde{Y}_k\big)$
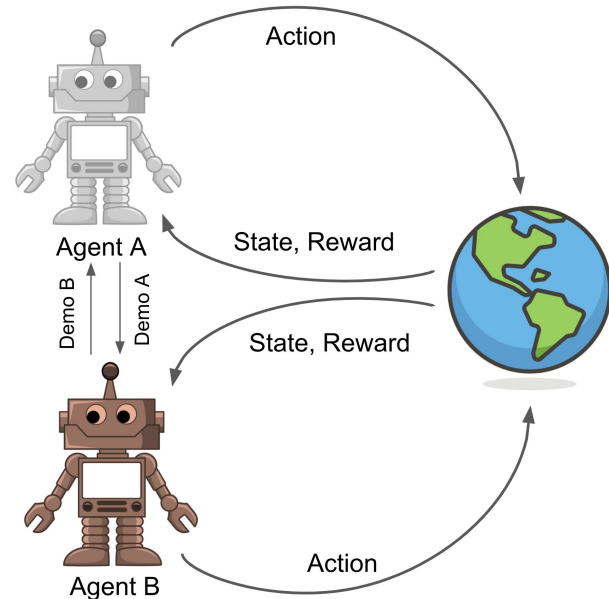
*PeerCT*

# PeerCT

Policy Co-Training (Song et al., 2019) is an instance of hybrid policy learning
- Two agents A,B with policies $\pi^A$ and $\pi^B$ that receive partial observations
- Agents are trained jointly to learn with rewards and noisy demonstrations from each other.
- For instance, consider agent A
  - Besides interacting with environment, A also receives $\{s_i, \pi_B(s_i)\}$ from agent B
  - We consider $\pi^B$ as the noisy version of the optimal policy



Agent A

Action

State, Reward

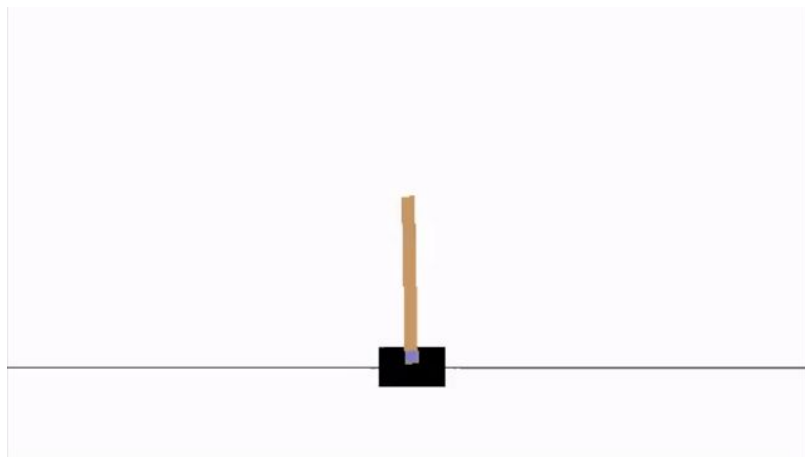Demo B        Demo A

State, Reward

Agent B

Action

Similar to the PeerBC setting, we use CA with weak supervision to handle the noise in imperfect demonstrations

$$J^{\mathrm{CT}}(\pi_\theta) = \mathbb{E}\left[ \mathsf{Eva}_\pi^{\mathrm{RL}}\big((s_i^A, a_i^A), r_i^A\big) + \mathsf{Eva}_\pi^{\mathrm{BC}}\big((s_i^A, a_i^A), a_i'^B\big) \right]$$
$$- \xi \cdot \mathbb{E}\left[ \mathsf{Eva}_\pi^{\mathrm{BC}}\big((s_j^A, a_j^A), a_k'^B\big) \right]$$

*Jialin Song, Ravi Lanka, Yisong Yue, and Masahiro Ono. Co-training for policy learning. In 429 UAI, page 441. AUAI Press, 2019*
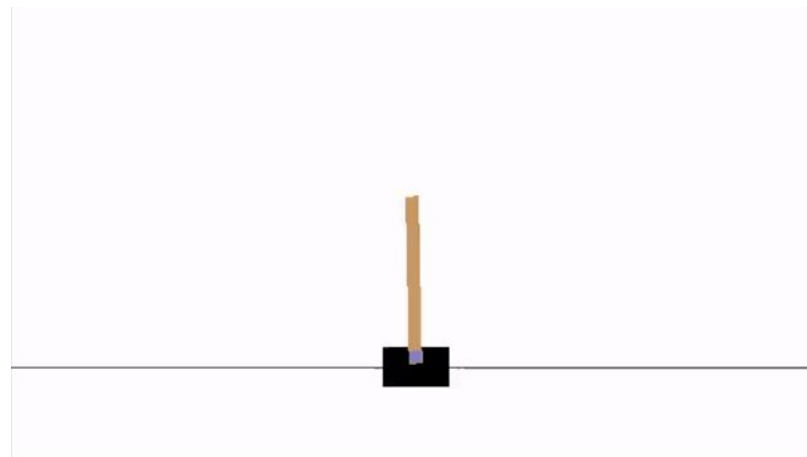
# An example of PeerRL on CartPole

- RL with Noisy Rewards ( $e_- = e_+ = 0.4$ )

- training 10,000 frames using Dueling-DQN
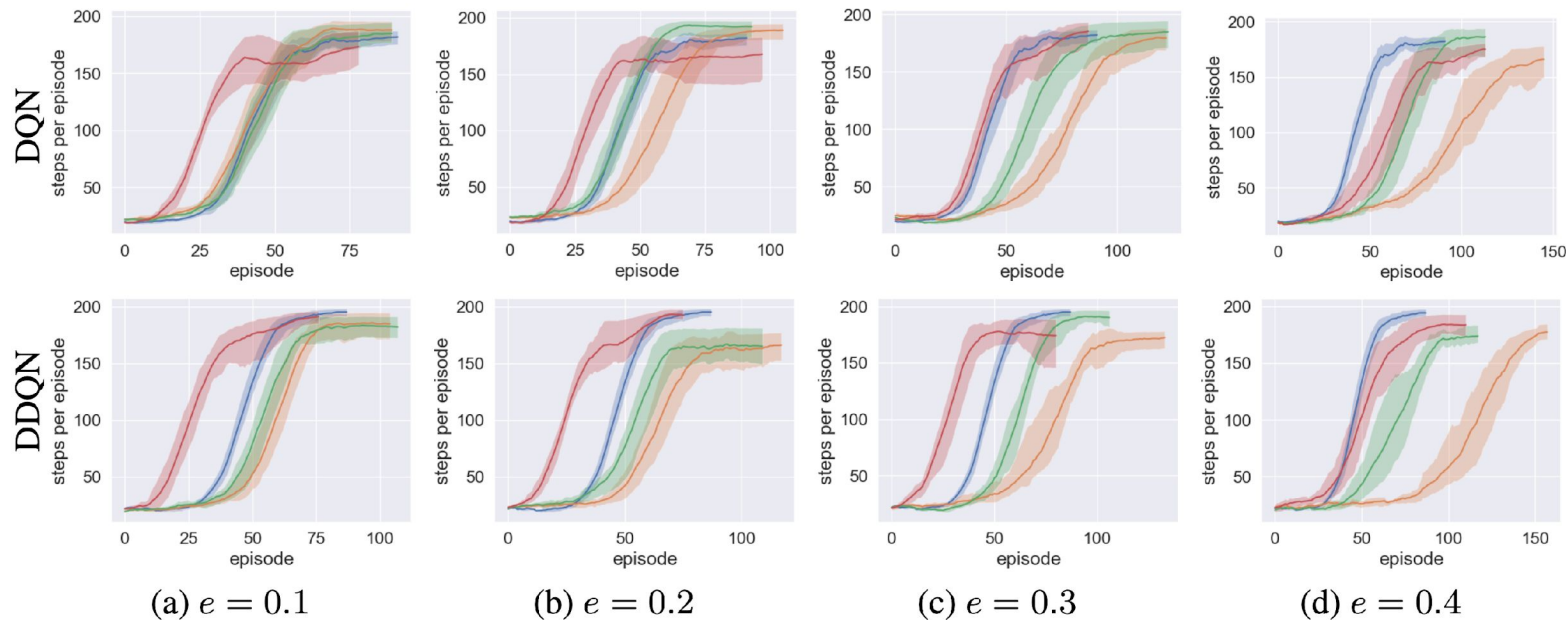


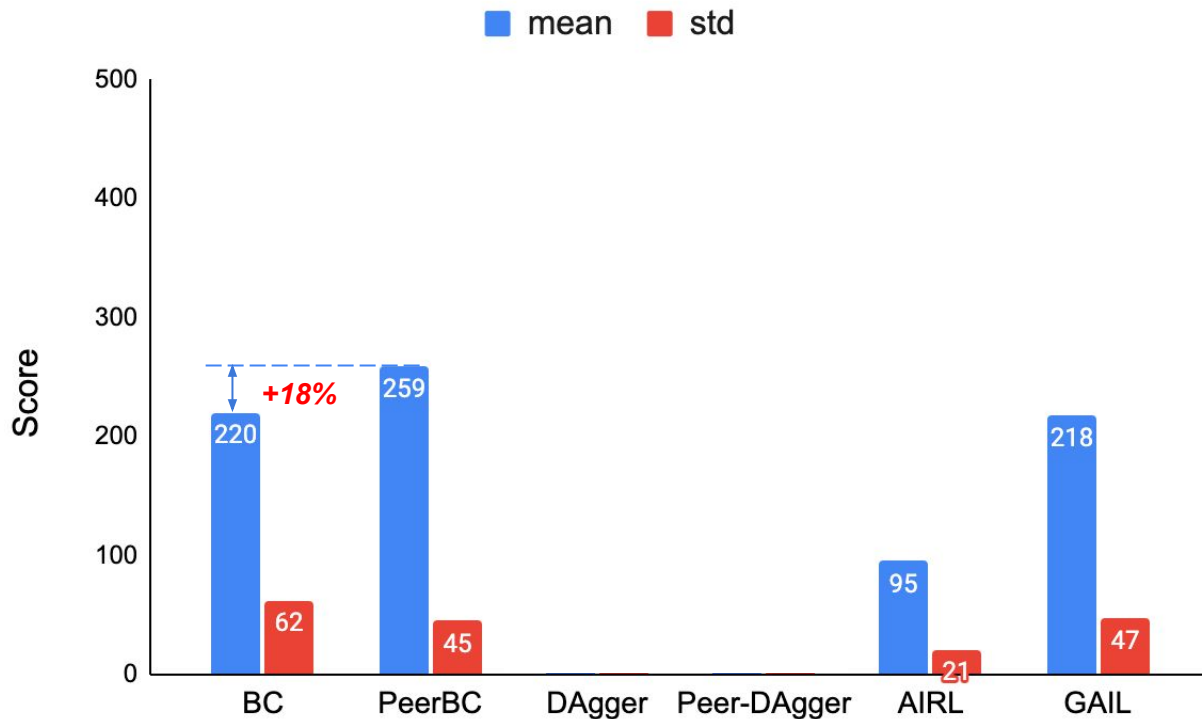noisy reward $\tilde{r}$



peer reward $\tilde{r}_{\text{peer}}$

# PeerRL recovers true reward signals

- ***CartPole***: training DDQN for 10,000 steps on, binary reward: $\{-1, 1\}$
- symmetric noise: $e = e_- = e_+$



(a) $e = 0.1$      (b) $e = 0.2$      (c) $e = 0.3$      (d) $e = 0.4$

true reward     noisy reward     surrogate reward (Wang et al, 2020)     PeerRL (ours)
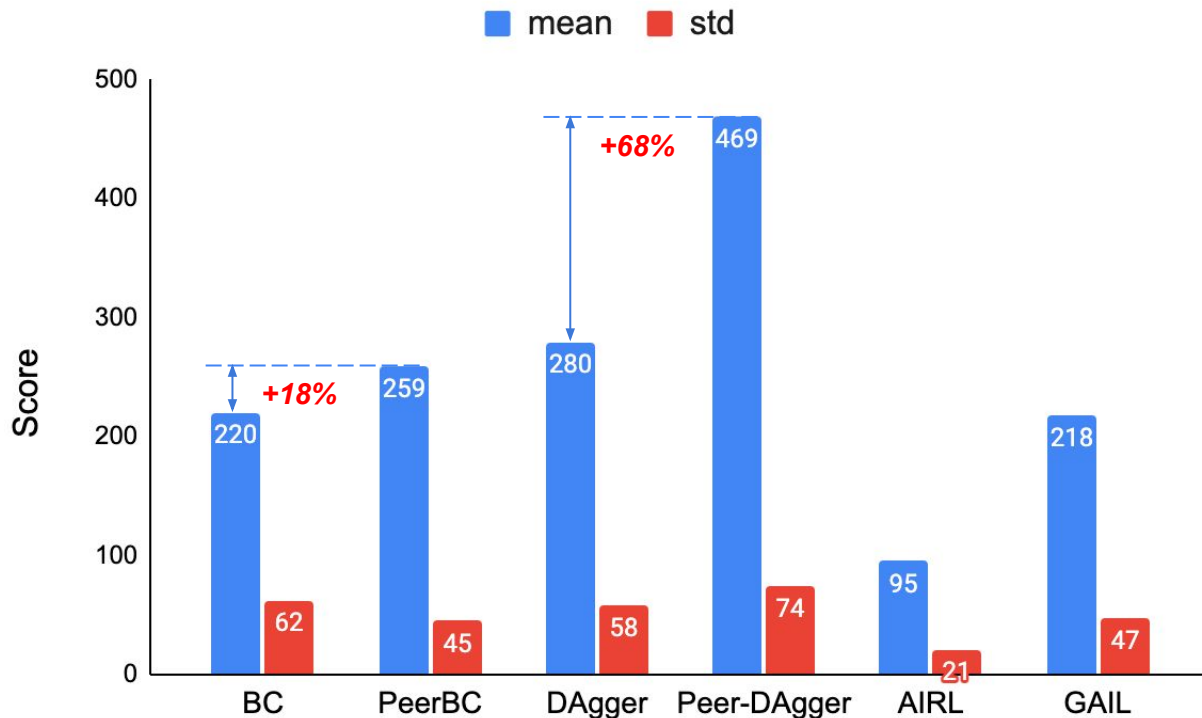
# PeerBC recovers true expert signals

- *CartPole-v1*: train an imperfect RL model with PPO algorithm, unroll 16 episodes
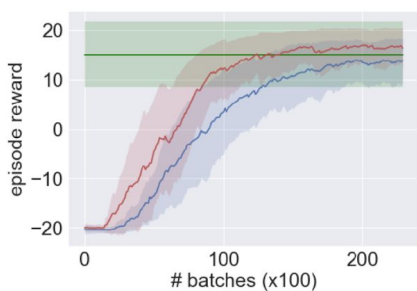
# PeerBC recovers true expert supervision signals

- *CartPole-v1*: train an imperfect RL model with PPO algorithm, unroll 16 episodes
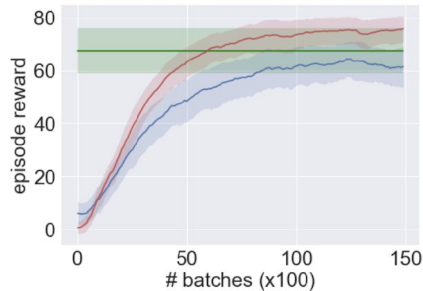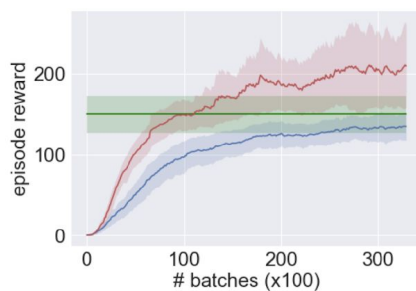
# PeerBC recovers true expert signals

- ***Atari games***: train an imperfect RL model with PPO algorithm
- weak expert = 70%~90% as good as fully converged agent
- collect demonstrations using weak expert and generate 100 trajectories for each environment
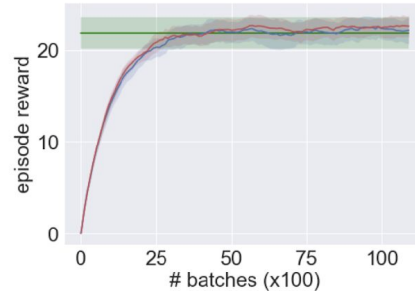- Note that no synthetic noise is added in the experiments
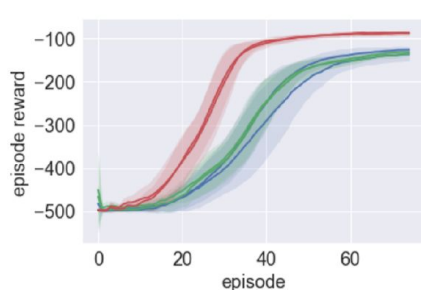


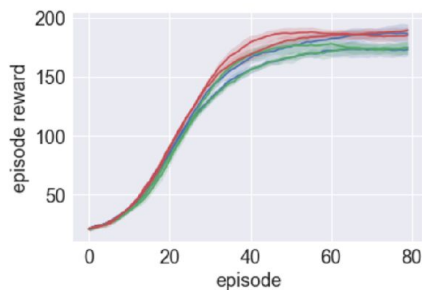(a) Pong     (b) Boxing     (c) Enduro     (d) Freeway

Standard BC     Weak expert     PeerBC (ours)

# PeerCT recovers true reward signals

- ***Continuous Control/Atari***:  adopt the exact same setting as Song et al., 2019 without any synthetic noise included
- removes all even index coordinates in the state vector (view-A) or removing all odd index ones (view-B)
- implies the potential of our approach to deal with natural noise in real-world applications



(a) Acrobot          (b) CartPole          (c) Pong          (d) Breakout

Single view          Co-Training (Song et al., 2019)          Peer Co-Training (ours)

# Sensitivity of over-agreement penalty

*Atari - Pong*



(a) $\xi = 0.1$    (b) $\xi = 0.2$    (c) $\xi = 0.3$    (d) $\xi = 0.4$

(e) $\xi = 0.5$    (f) $\xi = 0.6$    (g) $\xi = 0.7$    (h) $\xi = 0.8$

(i) $\xi = 0.9$    (j) $\xi = 1.0$    (k) $\xi = 1.1$    (l) $\xi = 1.2$

Standard BC    Weak expert    SQIL (Reddy et al., 2019)    PeerBC (ours)

# Sensitivity of over-agreement penalty

*Atari - Pong*

Our method works robustly in a wide range of $\xi$ !



(a) $\xi = 0.1$     (b) $\xi = 0.2$     (c) $\xi = 0.3$     (d) $\xi = 0.4$

(e) $\xi = 0.5$     (f) $\xi = 0.6$     (g) $\xi = 0.7$     (h) $\xi = 0.8$

(i) $\xi = 0.9$     (j) $\xi = 1.0$     (k) $\xi = 1.1$     (l) $\xi = 1.2$

Standard BC     Weak expert     SQIL (Reddy et al., 2019)     PeerBC (ours)
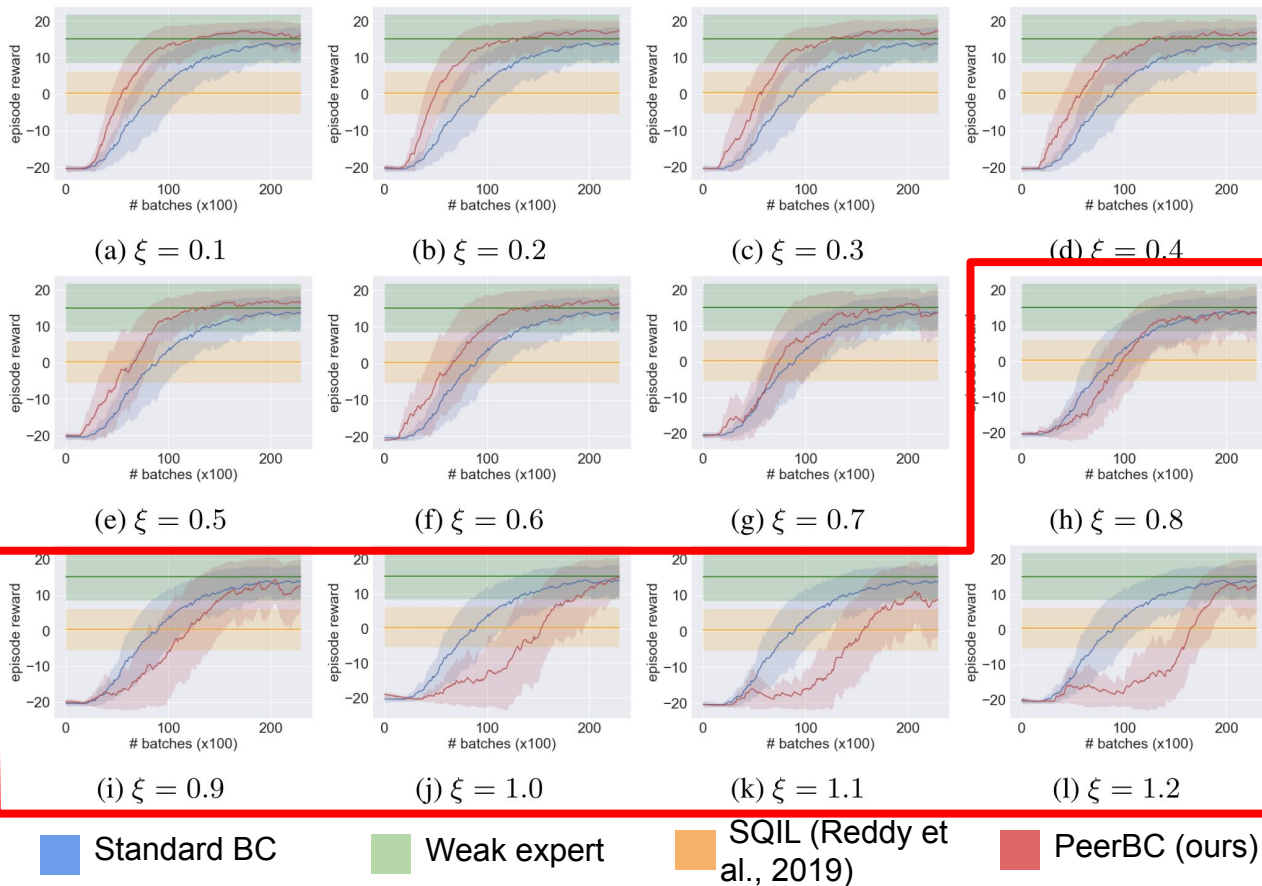
# Sensitivity of over-agreement penalty



**Atari - Pong**

Our method works robustly in a wide range of $\xi$ !

Overly large penalty introduces too much noise

(a) $\xi = 0.1$
(b) $\xi = 0.2$
(c) $\xi = 0.3$
(d) $\xi = 0.4$
(e) $\xi = 0.5$
(f) $\xi = 0.6$
(g) $\xi = 0.7$
(h) $\xi = 0.8$
(i) $\xi = 0.9$
(j) $\xi = 1.0$
(k) $\xi = 1.1$
(l) $\xi = 1.2$

Standard BC    Weak expert    SQIL (Reddy et al., 2019)    PeerBC (ours)

# Conclusion

- We provided a unified formulation of the ***weakly supervised policy learning*** problems

- We proposed PeerPL, a weakly supervised policy learning framework to unify a series of RL/BC problems with low-quality supervision signals

  - RL with perturbed reward

  - BC with imperfect demonstrations

  - Policy Co-Training (Hybrid RL + BC)

- Our method is theoretically guaranteed to recover the optimal policy with sufficient weak supervision signals.