



## $\Psi\Phi$ -Learning: Reinforcement Learning with Demonstrations using Successor Features and Inverse Temporal Difference Learning

---



Angelos  
Filos <sup>$\phi$</sup>



Clare  
Lyle <sup>$\phi$</sup>



Yarin  
Gal <sup>$\phi$</sup>



Sergey  
Levine <sup>$\psi\gamma$</sup>



Natasha  
Jaques <sup>$*\psi\gamma$</sup>



Gregory  
Farquhar <sup>$*\delta$</sup>

## Social Reinforcement Learning [ $\alpha$ . Problem Setting]

ENV



Control Markov Process

$\mathcal{C}$

$\langle S, A, P, \gamma, d_0 \rangle$

# Social Reinforcement Learning [ $\alpha$ . Problem Setting]

ENV



Control Markov Process

$C$

$\langle S, A, P, \gamma, d_0 \rangle$

EGO



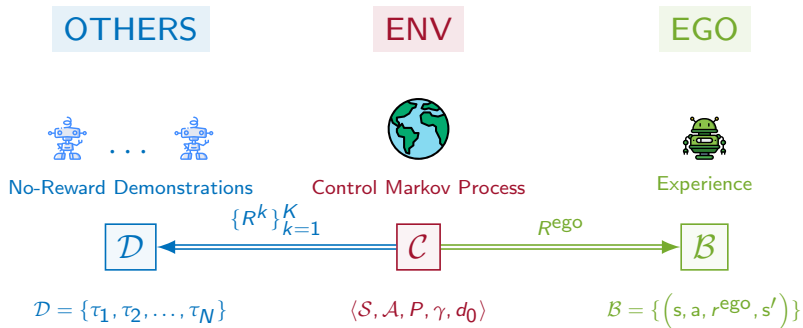
Experience

$B$

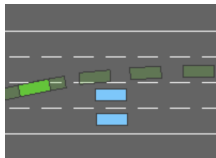
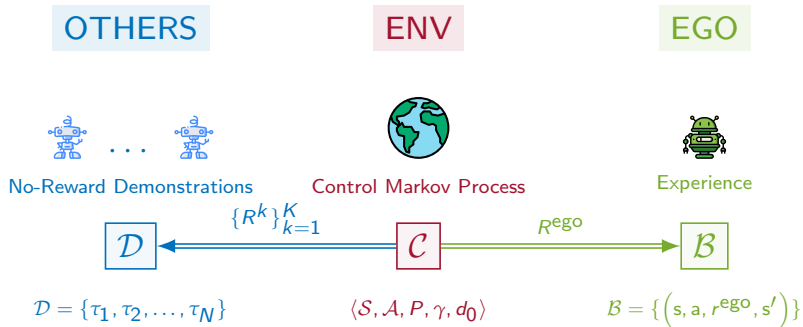
$B = \{(s, a, r^{\text{ego}}, s')\}$

$R^{\text{ego}}$

# Social Reinforcement Learning [ $\alpha$ . Problem Setting]

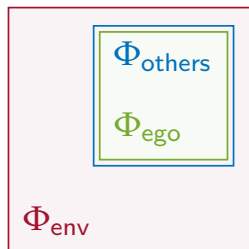


# Social Reinforcement Learning [ $\alpha$ . Problem Setting]



Working Example

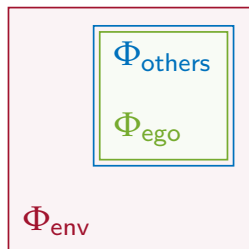




**Single-Task**  
 $R^{ego} = \{R^k\}$

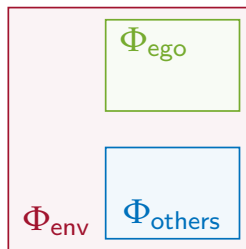


# Social Reinforcement Learning [ $\beta$ . Taxonomy of Tasks]



**Single-Task**

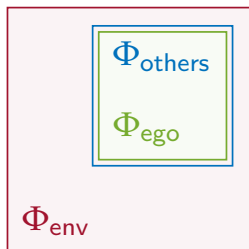
$$R^{ego} = \{R^k\}$$



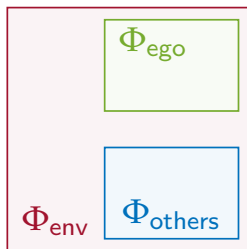
**Adversarial Task**

$$\Phi(R^{ego}) \cap \Phi(\{R^k\}) = \emptyset$$

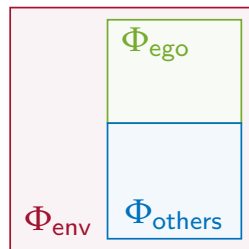
# Social Reinforcement Learning [ $\beta$ . Taxonomy of Tasks]



**Single-Task**  
 $R^{ego} = \{R^k\}$



**Adversarial Task**  
 $\Phi(R^{ego}) \cap \Phi(\{R^k\}) = \emptyset$

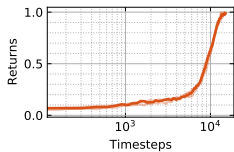


**Multi-Task**  
 $\Phi(R^{ego}) \cap \Phi(\{R^k\}) \neq \emptyset$

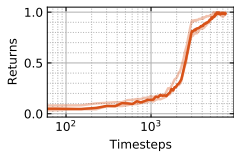
# Baselines [ $\alpha$ . Reinforcement Learning]

— RL

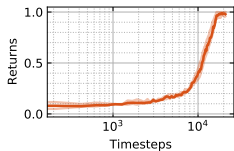
Single Task



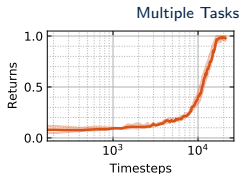
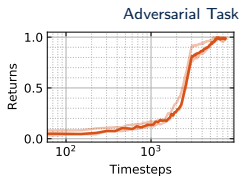
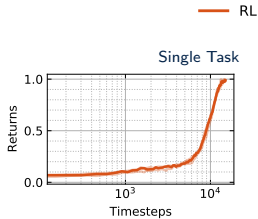
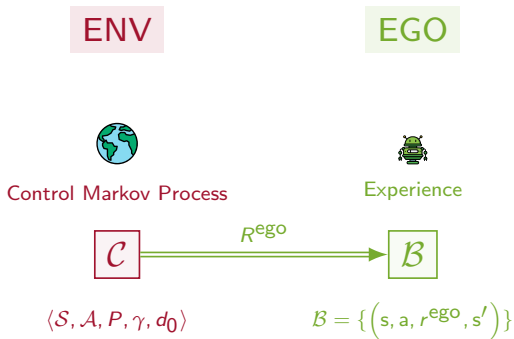
Adversarial Task



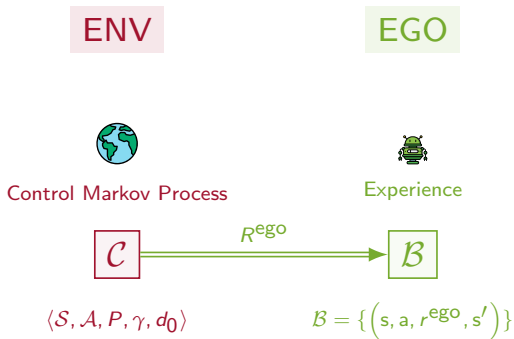
Multiple Tasks



# Baselines [ $\alpha$ . Reinforcement Learning]



# Baselines [ $\alpha$ . Reinforcement Learning]

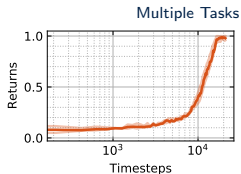
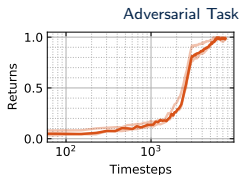
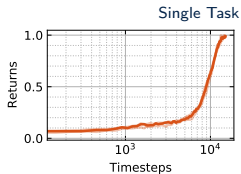


ignores OTHERS' decisions



fails to accelerate its learning from demonstrations

— RL



# Baselines [ $\beta$ . Learning From Demonstrations]

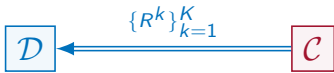
OTHERS

ENV



No-Reward Demonstrations

Control Markov Process



$\mathcal{D} = \{\tau_1, \tau_2, \dots, \tau_N\}$

$\langle \mathcal{S}, \mathcal{A}, P, \gamma, d_0 \rangle$

# Baselines [ $\beta$ . Learning From Demonstrations]

OTHERS



No-Reward Demonstrations



$$\mathcal{D} = \{\tau_1, \tau_2, \dots, \tau_N\}$$

ENV



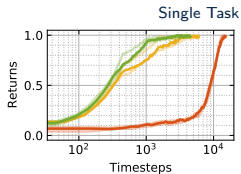
Control Markov Process



$$\langle S, \mathcal{A}, P, \gamma, d_0 \rangle$$

$$\{R^k\}_{k=1}^K$$

— RL — BC — SQLv2



# Baselines [ $\beta$ . Learning From Demonstrations]

OTHERS



No-Reward Demonstrations



$$\mathcal{D} = \{\tau_1, \tau_2, \dots, \tau_N\}$$

ENV

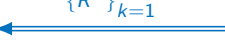


Control Markov Process

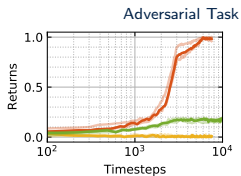
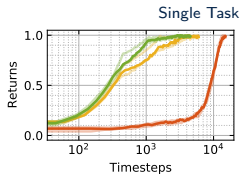


$$\langle S, A, P, \gamma, d_0 \rangle$$

$$\{R^k\}_{k=1}^K$$



— RL — BC — SQLv2





# Baselines [ $\beta$ . Learning From Demonstrations]

OTHERS

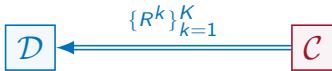


ENV



No-Reward Demonstrations

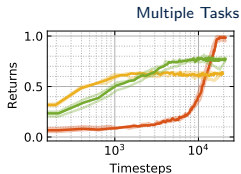
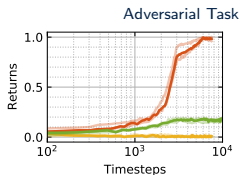
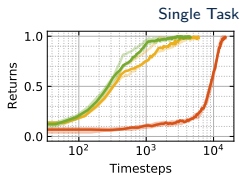
Control Markov Process



$$\mathcal{D} = \{\tau_1, \tau_2, \dots, \tau_N\}$$

$$\langle \mathcal{S}, \mathcal{A}, P, \gamma, d_0 \rangle$$

— RL — BC — SQLv2



# Baselines [ $\beta$ . Learning From Demonstrations]

OTHERS



No-Reward Demonstrations



$$D = \{\tau_1, \tau_2, \dots, \tau_N\}$$

ENV



Control Markov Process



$$\langle S, A, P, \gamma, d_0 \rangle$$

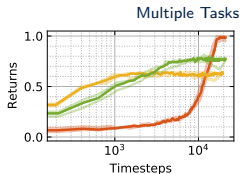
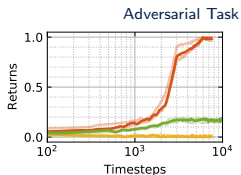
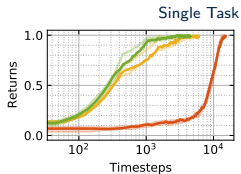
$$\{R^k\}_{k=1}^K$$

ignores EGO experience



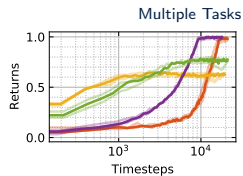
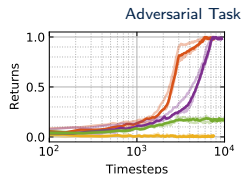
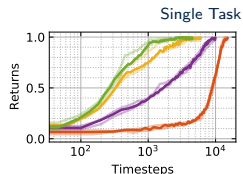
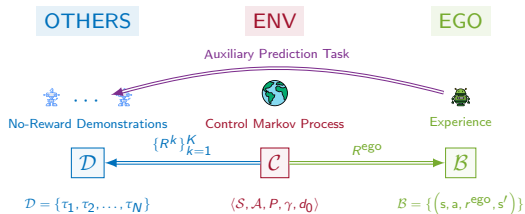
bottlenecked by the quality/relevance  
of the demonstrations

— RL — BC — SQLv2

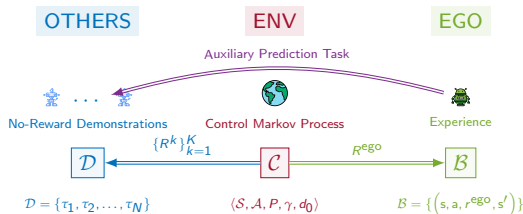


# Baselines [ $\gamma$ . Behavioural Cloning as an Auxiliary Task]

— RL — RL + BC-Aux — SQILv2  
— BC



# Baselines [ $\gamma$ . Behavioural Cloning as an Auxiliary Task]

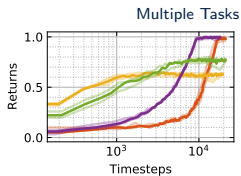
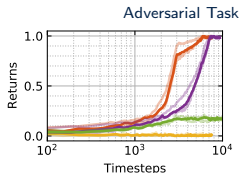
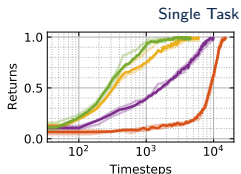


ignores structure in demonstrations

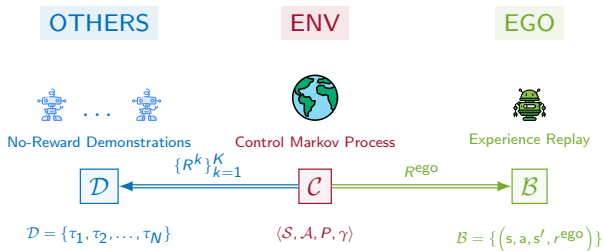


helps only indirectly via representation learning

— RL — RL + BC-Aux — SQILv2  
— BC



# Desiderata for Social Reinforcement Learners



$\alpha$ . Preserve unbiased asymptotic performance of RL.

\* RL, BC, IRL, RL+BC-Aux

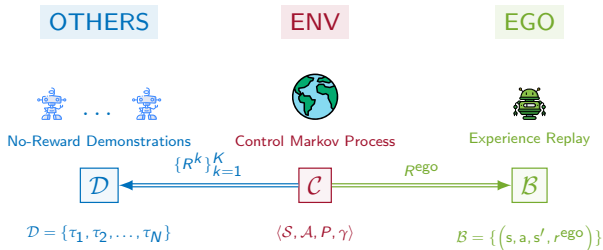
$\beta$ . Accelerate RL from (no-reward) demonstrations.

\* RL, BC, IRL, RL+BC-Aux

$\gamma$ . Focus on *actionable representations* that inform action selection.

\* RL, BC, IRL, RL+BC-Aux

# Desiderata for Social Reinforcement Learners



$\alpha$ . Preserve unbiased asymptotic performance of RL.

\* RL, BC, IRL, RL+BC-Aux

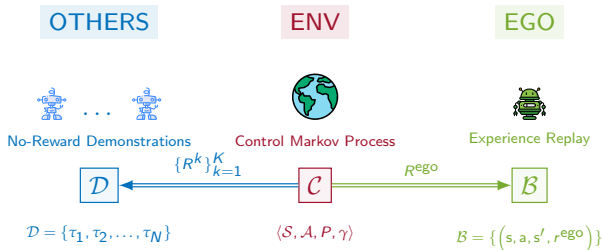
$\beta$ . Accelerate RL from (no-reward) demonstrations.

\* RL, BC, IRL, RL+BC-Aux

$\gamma$ . Focus on *actionable representations* that inform action selection.

\* RL, BC, IRL, RL+BC-Aux

# Desiderata for Social Reinforcement Learners



$\alpha$ . Preserve unbiased asymptotic performance of RL.

\* RL, BC, IRL, RL+BC-Aux

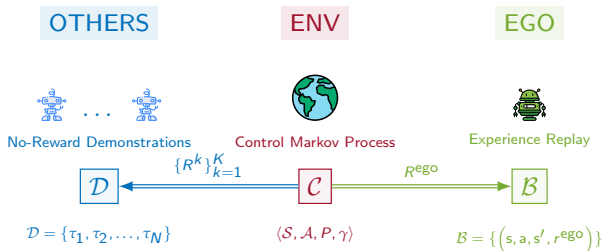
$\beta$ . Accelerate RL from (no-reward) demonstrations.

\* RL, BC, IRL, RL+BC-Aux

$\gamma$ . Focus on *actionable representations* that inform action selection.

\* RL, BC, IRL, RL+BC-Aux

# Desiderata for Social Reinforcement Learners



$\alpha$ . Preserve unbiased asymptotic performance of RL.

\* RL, BC, IRL, RL+BC-Aux

$\beta$ . Accelerate RL from (no-reward) demonstrations.

\* RL, BC, IRL, RL+BC-Aux

$\gamma$ . Focus on *actionable representations* that inform action selection.

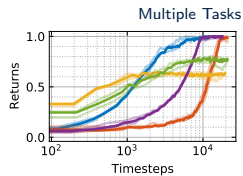
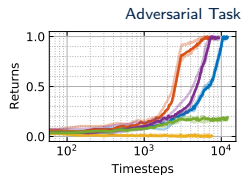
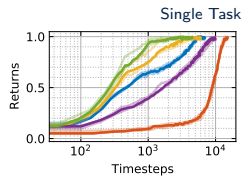
\* RL, BC, IRL, RL+BC-Aux





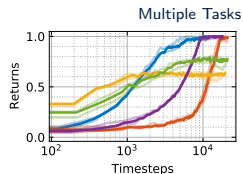
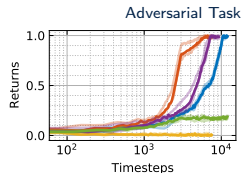
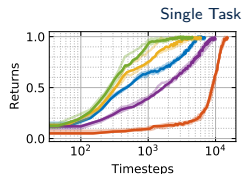
# $\Psi\Phi$ -Learning [ $\alpha$ . Principle & Results]

—  $\Psi\Phi$ L (ours) — BC — SQILv2  
— RL — RL + BC-Aux



## Modelling Principle

The *OTHER* agents are goal-directed and optimal for some task. Their behaviour should be integrated to the *EGO* agent's policy improvement only when it is relevant to its task.

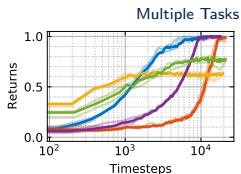
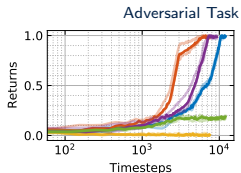
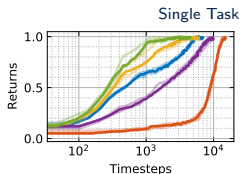


## Modelling Principle

The *OTHER* agents are goal-directed and optimal for some task. Their behaviour should be integrated to the *EGO* agent's policy improvement only when it is relevant to its task.

## Reasoning about Tasks

- $\alpha$ . Task space.
- $\beta$ . Task inference.



## Modelling Principle

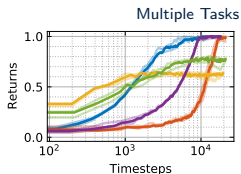
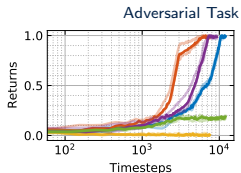
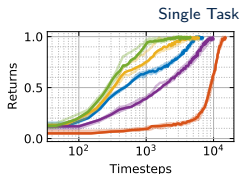
The *OTHER* agents are goal-directed and optimal for some task. Their behaviour should be integrated to the *EGO* agent's policy improvement only when it is relevant to its task.

## Reasoning about Tasks

- $\alpha$ . Task space.
- $\beta$ . Task inference.

## Task-Aware Policy Ops

- $\alpha$ . Fast policy evaluation.
- $\beta$ . Fast policy improvement.



$Q(s, a)$

$$Q^{\pi, w}(s, a)$$

$$Q^{\pi, \mathbf{w}}(s, a) \triangleq \Psi^\pi(s, a)^\top \mathbf{w}$$

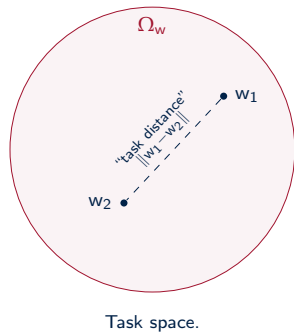


$$Q^{\pi, w}(s, a) \triangleq \Psi^{\pi}(s, a)^{\top} w$$

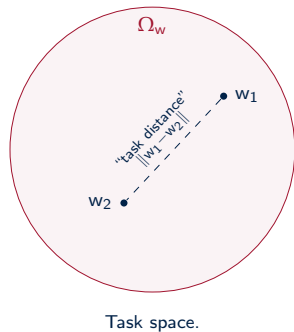
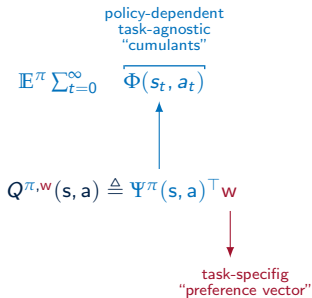
↓  
task-specific  
"preference vector"

$$Q^{\pi, w}(s, a) \triangleq \Psi^{\pi}(s, a)^{\top} w$$

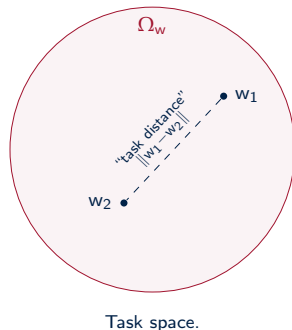
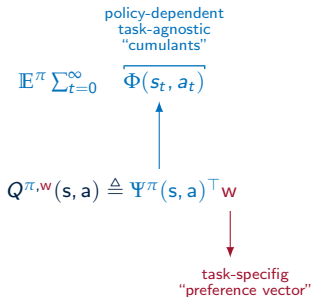
↓  
task-specific  
"preference vector"



# $\Psi\Phi$ -Learning [ $\beta$ . Successor Features Reparametrisation]



# $\Psi\Phi$ -Learning [ $\beta$ . Successor Features Reparametrisation]

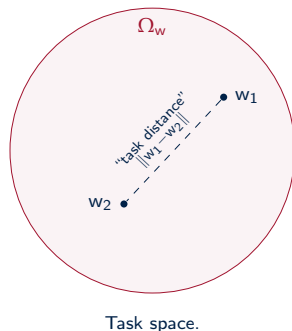
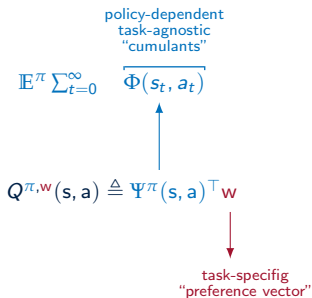


## Generalised Policy Improvement (GPI, Barreto et al., 2017)

Given the **EGO** agent's preference vector  $w^{\text{ego}}$  and the **OTHER** agents' successor features  $\{\Psi^k\}_{k=1}^K$ , we can improve the **EGO** agent's policy by acting according to:

$$\pi'(s) = \arg \max_a \max_{i=[K], \text{ego}} \underbrace{\Psi^i(s, a)^{\top} w^{\text{ego}}}_{\text{evaluate policy } i \text{ under task } w^{\text{ego}}} \succeq \pi^{\text{ego}}.$$

# $\Psi\Phi$ -Learning [ $\beta$ . Successor Features Reparametrisation]



## Generalised Policy Improvement (GPI, Barreto et al., 2017)

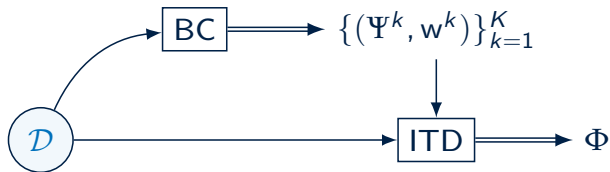
Given the **EGO** agent's preference vector  $w^{\text{ego}}$  and the **OTHER** agents' successor features  $\{\Psi^k\}_{k=1}^K$ , we can improve the **EGO** agent's policy by acting according to:

$$\pi'(s) = \arg \max_a \max_{i=[K], \text{ego}} \underbrace{\Psi^i(s, a)^{\top} w^{\text{ego}}}_{\text{evaluate policy } i \text{ under task } w^{\text{ego}}} \succeq \pi^{\text{ego}}.$$

Where do cumulants  $\Phi$  and preference vectors  $w$  come from?

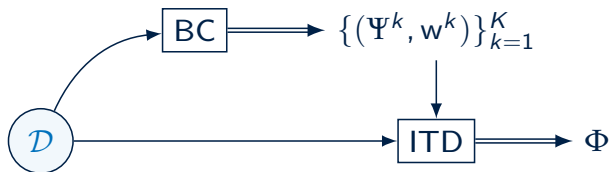


# $\Psi\Phi$ -Learning [ $\gamma$ . Inverse Temporal Difference Learning]



I/O schematic for inverse TD learning algorithm.

# $\Psi\Phi$ -Learning [ $\gamma$ . Inverse Temporal Difference Learning]



I/O schematic for inverse TD learning algorithm.

## Successor Features & Bellman Consistency $\Rightarrow$ Cumulants

$$\mathcal{L}_{\text{BC-Q}}(\theta_{\Psi^k}, \mathbf{w}^k) \triangleq -\mathbb{E} \log \frac{\exp(\Psi(s_t, \mathbf{a}_t; \theta_{\Psi^k})^\top \mathbf{w}^k)}{\sum_a \exp(\Psi(s_t, a; \theta_{\Psi^k})^\top \mathbf{w}^k)}$$

$$\mathcal{L}_{\text{ITD}}(\theta_\Phi, \theta_{\Psi^k}) \triangleq \mathbb{E} \|\Psi(s_t, \mathbf{a}_t; \theta_{\Psi^k}) - \Phi(s_t, \mathbf{a}_t; \theta_\Phi) - \gamma \Psi(s_{t+1}, \mathbf{a}_{t+1}; \theta_{\Psi^k})\|.$$



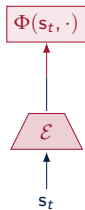
ENV



Control Markov Process

$\mathcal{C}$

$\langle \mathcal{S}, \mathcal{A}, P, \gamma \rangle$



# $\Psi\Phi$ -Learning [ $\delta$ . Implementation]

ENV



Control Markov Process

$C$

$\langle S, \mathcal{A}, P, \gamma \rangle$

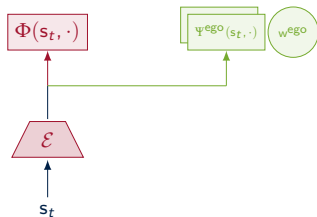
EGO



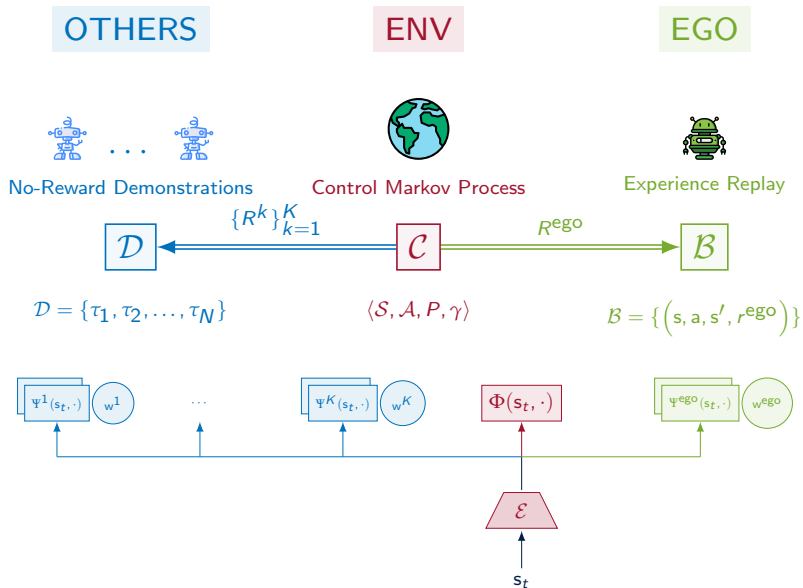
Experience Replay

$B$

$B = \{(s, a, s', r^{\text{ego}})\}$

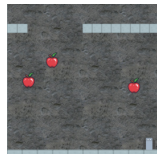
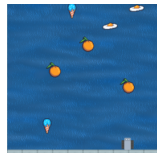
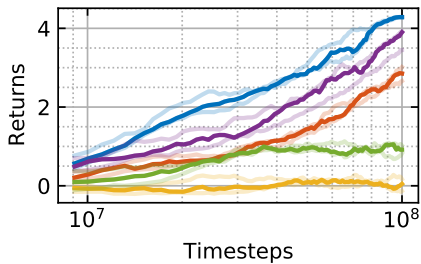


# $\Psi\Phi$ -Learning [ $\delta$ . Implementation]



# Experiments [ $\alpha$ . Deep RL from RGB Observations]

—  $\Psi\Phi L$  (ours) — RL — BC — RL + BC-Aux — SQILv2



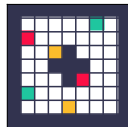
## Experiments [ $\beta$ . Few-Shot Transfer to New Reward Functions]

# Experiments [ $\beta$ . Few-Shot Transfer to New Reward Functions]

Methods	0-shot		1-shot		100-shot	
	R-G	-R-G	R-G	-R-G	R-G	-R-G
SQILv2 <sup>*</sup> (Reddy <i>et al.</i> , 2019)	0.0 $\pm$ 0.0	-1.0 $\pm$ 0.0	0.0 $\pm$ 0.0	-1.0 $\pm$ 0.0	1.0 $\pm$ 0.0	1.0 $\pm$ 0.0
$\Psi\Phi$ -learning $\diamond$	0.2 $\pm$ 0.1	-0.4 $\pm$ 0.2	1.0 $\pm$ 0.0	1.0 $\pm$ 0.0	1.0 $\pm$ 0.0	1.0 $\pm$ 0.0



# Experiments [ $\beta$ . Few-Shot Transfer to New Reward Functions]

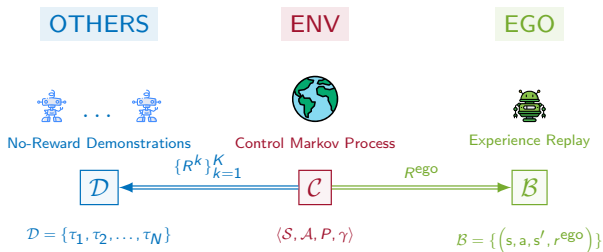


Methods	0-shot		1-shot		100-shot	
	R-G	-R-G	R-G	-R-G	R-G	-R-G
SQILv2 <sup>*</sup> (Reddy <i>et al.</i> , 2019)	0.0 $\pm$ 0.0	-1.0 $\pm$ 0.0	0.0 $\pm$ 0.0	-1.0 $\pm$ 0.0	1.0 $\pm$ 0.0	1.0 $\pm$ 0.0
$\Psi\Phi$ -learning $\diamond$	0.2 $\pm$ 0.1	-0.4 $\pm$ 0.2	1.0 $\pm$ 0.0	1.0 $\pm$ 0.0	1.0 $\pm$ 0.0	1.0 $\pm$ 0.0

## Zero-Shot Generalisation Bound of $\Psi\Phi$ -Learning)

$$\underbrace{Q^*(s, a) - Q^\pi(s, a)}_{\text{sub-optimality gap}} \leq \frac{2}{1 - \gamma} \left[ \underbrace{(\phi_{\max} \|w_j - w'\|)}_{\text{relevance of demonstrations}} + 2\delta_r \right] + \underbrace{\|w'\| \delta_\Psi + \frac{1}{(1 - \gamma)} \delta_r}_{\text{approximation error}}$$

# Summary

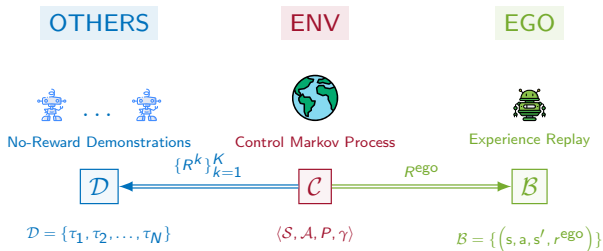


## $\Psi\Phi$ -Learning

- Combination of our **offline** IRL method with generalised policy improvement.
- Utilisation of no-reward demonstrations for accelerating RL.
- Graceful fallback to solitary RL when demos are of “poor” quality (theorem).
- Few-shot transfer to new reward functions.
- Scalable for deep RL, RGB observations.

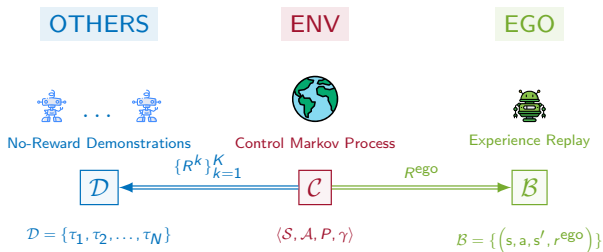


# Summary



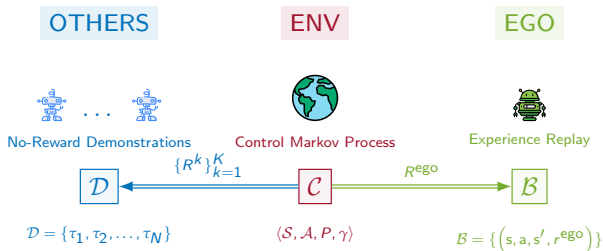
## $\Psi\Phi$ -Learning

- $\alpha$ . Combination of our **offline** IRL method with generalised policy improvement.
- $\beta$ . Utilisation of no-reward demonstrations for accelerating RL.
- $\gamma$ . Graceful fallback to solitary RL when demos are of “poor” quality (theorem).
- $\delta$ . Few-shot transfer to new reward functions.
- $\epsilon$ . Scalable for deep RL, RGB observations.



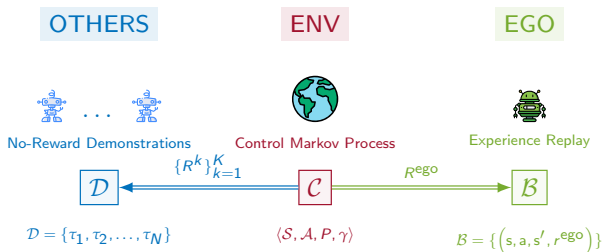
## $\Psi\Phi$ -Learning

- $\alpha$ . Combination of our **offline** IRL method with generalised policy improvement.
- $\beta$ . Utilisation of no-reward demonstrations for accelerating RL.
- $\gamma$ . Graceful fallback to solitary RL when demos are of “poor” quality (theorem).
- $\delta$ . Few-shot transfer to new reward functions.
- $\epsilon$ . Scalable for deep RL, RGB observations.



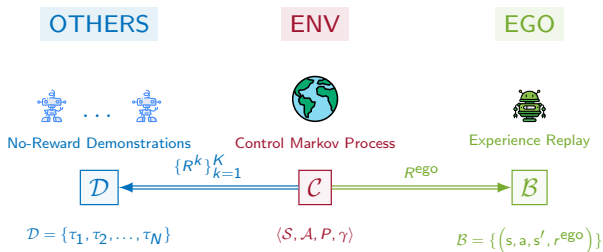
## $\Psi\Phi$ -Learning

- $\alpha$ . Combination of our **offline** IRL method with generalised policy improvement.
- $\beta$ . Utilisation of no-reward demonstrations for accelerating RL.
- $\gamma$ . Graceful fallback to solitary RL when demos are of "poor" quality (theorem).
- $\delta$ . Few-shot transfer to new reward functions.
- $\epsilon$ . Scalable for deep RL, RGB observations.



## $\Psi\Phi$ -Learning

- $\alpha$ . Combination of our **offline** IRL method with generalised policy improvement.
- $\beta$ . Utilisation of no-reward demonstrations for accelerating RL.
- $\gamma$ . Graceful fallback to solitary RL when demos are of "poor" quality (theorem).
- $\delta$ . Few-shot transfer to new reward functions.
- $\epsilon$ . Scalable for deep RL, RGB observations.



## $\Psi\Phi$ -Learning

- $\alpha$ . Combination of our **offline** IRL method with generalised policy improvement.
- $\beta$ . Utilisation of no-reward demonstrations for accelerating RL.
- $\gamma$ . Graceful fallback to solitary RL when demos are of "poor" quality (theorem).
- $\delta$ . Few-shot transfer to new reward functions.
- $\epsilon$ . Scalable for deep RL, RGB observations.