

Successor Features for Transfer in Reinforcement Learning

Michaela Hanajikova, Nicola Maritan, Maria Sgaravatto

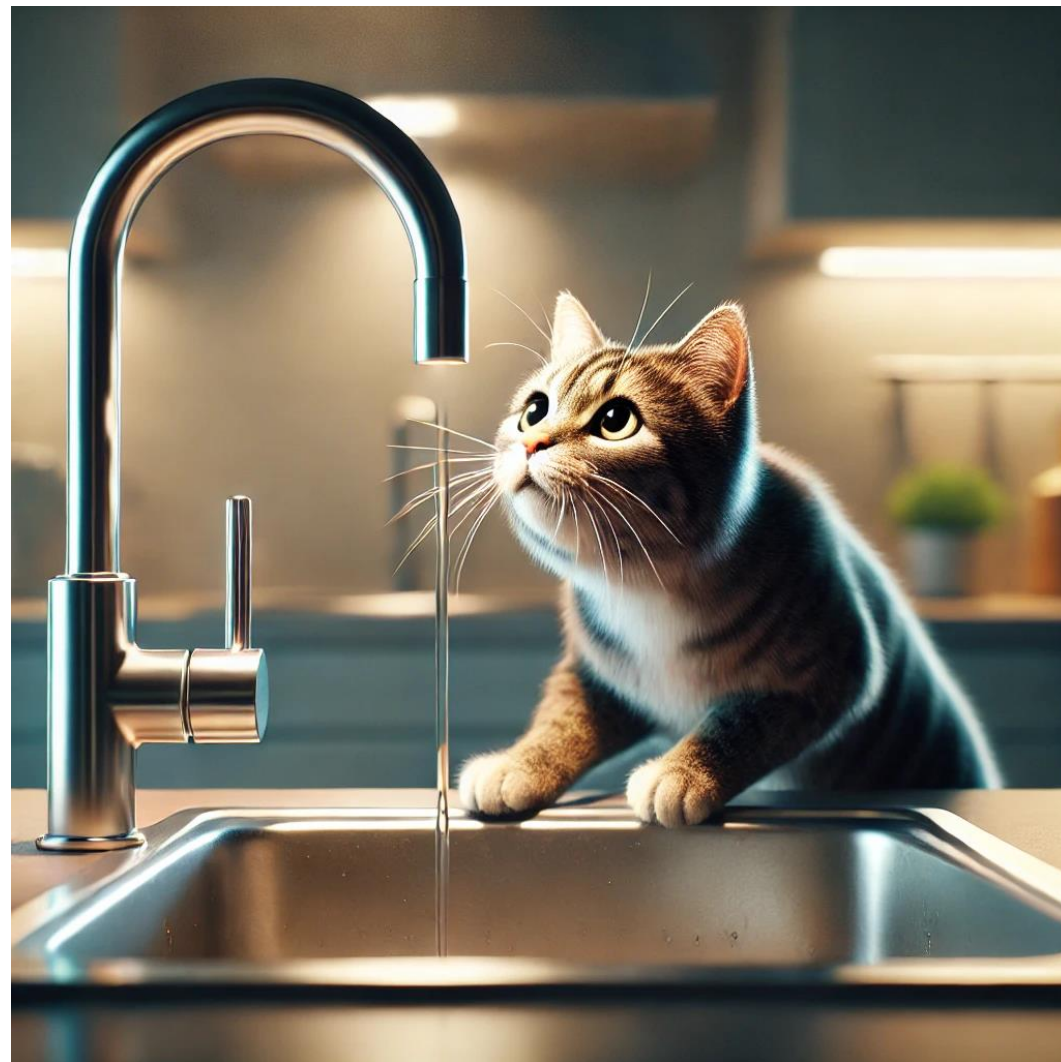
What is *transfer* in Reinforcement Learning?

Transfer in Reinforcement Learning

- Process of leveraging knowledge acquired in one or more tasks to improve performance on other tasks.
- In this paper, *a task is a Markov Decision Process (MDP)*.
- In particular, we are interested in the transfer between tasks:
 - Having the same components of an MDP...
 - ... except for the reward function.

Example

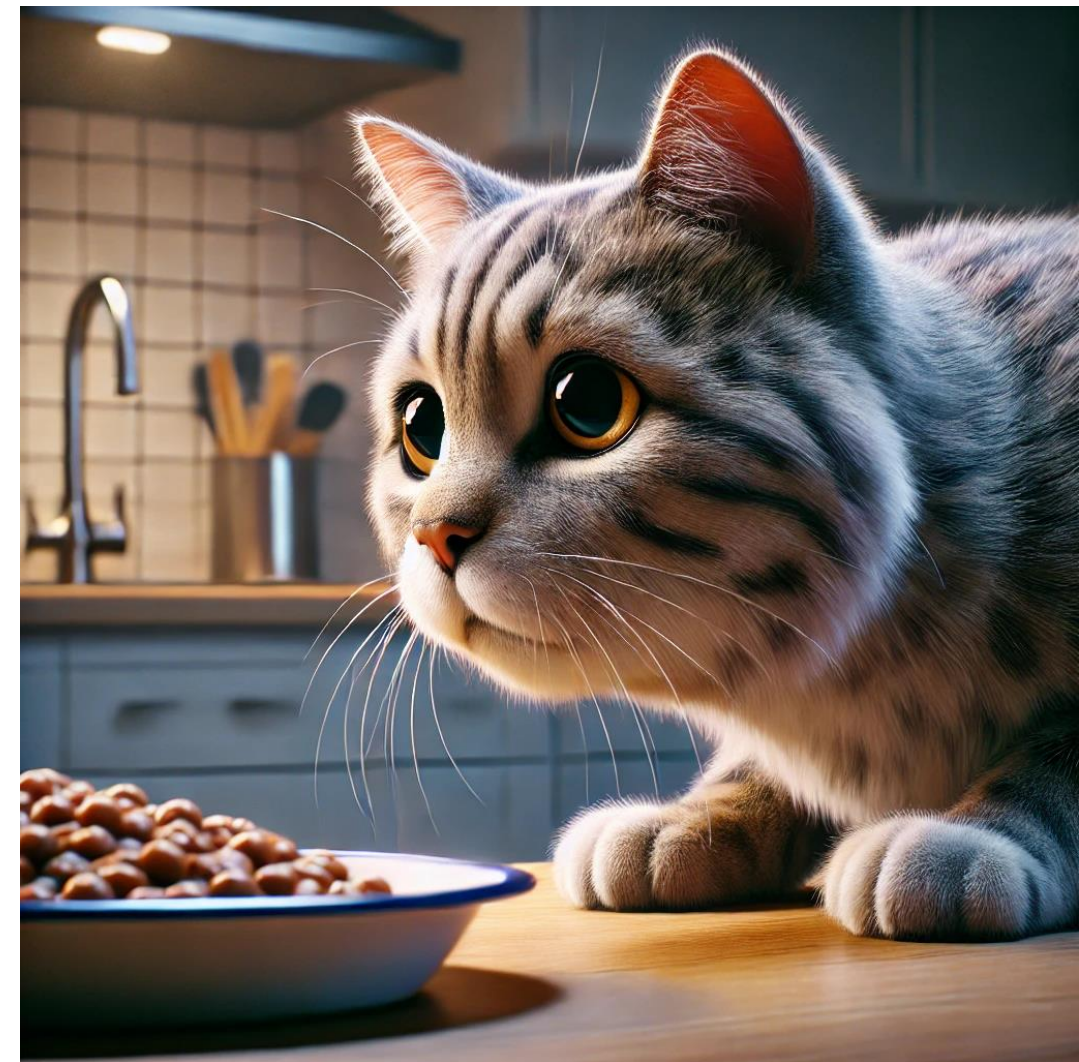
$$M_{\text{thirsty}} = (\mathcal{S}, \mathcal{A}, p, r_{\text{thirsty}}, \gamma)$$



Transfer



$$M_{\text{hungry}} = (\mathcal{S}, \mathcal{A}, p, r_{\text{hungry}}, \gamma)$$



Same MDP except for the reward function.

Successor features

Successor features

Consider an MDP

$$M = (\mathcal{S}, \mathcal{A}, p, r, \gamma).$$

We can express the one-step reward associated with transition (s, a, s') as

$$r(s, a, s') = \boldsymbol{\phi}(s, a, s')^\top \mathbf{w}$$

without loss of generality.

Successor features

Let $\boldsymbol{\phi}_t \triangleq \boldsymbol{\phi}(s_t, a_t, s_{t+1})$. For a given policy π on M we have:

$$\begin{aligned} Q^\pi(s, a) &\triangleq \mathbb{E}^\pi[r_{t+1} + \gamma r_{t+2} + \cdots \mid s_t = s, a_t = a] \\ &= \mathbb{E}^\pi[\boldsymbol{\phi}_{t+1}^\top \mathbf{w} + \gamma \boldsymbol{\phi}_{t+2}^\top \mathbf{w} + \cdots \mid s_t = s, a_t = a] \\ &= \mathbb{E}^\pi\left[\sum_{i=t}^{\infty} \gamma^{i-t} \boldsymbol{\phi}_{i+1} \mid s_t = s, a_t = a\right]^\top \mathbf{w} \\ &= \boldsymbol{\psi}^\pi(s, a)^\top \mathbf{w}. \end{aligned}$$

Successor features

Consider the decomposition

$$Q^{\pi}(s, a) = \boldsymbol{\psi}^{\pi}(s, a)^{\top} \mathbf{w}.$$

We call $\boldsymbol{\psi}^{\pi}(s, a)$ the *successor features* (SFs) of (s, a) under policy π .

Successor features

Consider the decomposition

$$Q^{\pi}(s, a) = \boldsymbol{\psi}^{\pi}(s, a)^{\top} \mathbf{w}.$$

Main idea is having:

- $\boldsymbol{\psi}^{\pi}$ to summarize M 's dynamics induced by π on the environment.
- \mathbf{w} to capture M 's rewards.

Transfer via SFs

Family of MDPs

Suppose ϕ is fixed. We define

$$\begin{aligned}\mathcal{M}^\phi(\mathcal{S}, \mathcal{A}, p, \gamma) &\triangleq \\ &\triangleq \{M(\mathcal{S}, \mathcal{A}, p, r, \gamma) \mid r(s, a, s') = \phi(s, a, s')^\top \mathbf{w}\}.\end{aligned}$$

The task $M_i \in \mathcal{M}^\phi$ is entirely defined by the corresponding \mathbf{w}_i .

Generalized Policy Improvement with SFs

Let:

- π_i^* be the optimal policy for M_i and $Q_i^{\pi_i^*}$ its action-value function.
- $Q_j^{\pi_i^*}$ the action-value function of π_i^* when executed in $M_j \in \mathcal{M}\Phi$.

Generalized Policy Improvement with SFs

Suppose you know

$$\left\{ Q_i^{\pi_j^*} \right\}_{j=1, \dots, n}$$

for $M_1, \dots, M_n \in \mathcal{M}^\Phi$ and $M_i \in \mathcal{M}^\Phi$.

Greedy choose

$$\pi(s) \in \operatorname{argmax}_a \max_j Q_i^{\pi_j^*}(s, a).$$

Generalized Policy Improvement with SFs

$$\pi(s) \in \operatorname{argmax}_a \max_j Q_i^{\pi_j^\star}(s, a)$$

Then we are guaranteed that

$$Q_i^{\pi_i^\star}(s, a) - Q_i^\pi(s, a) \leq \frac{2}{1 - \gamma} \phi_{\max} \min_j \|\mathbf{w}_i - \mathbf{w}_j\|$$

where $\phi_{\max} = \max_{s,a} \|\phi(s, a)\|$.

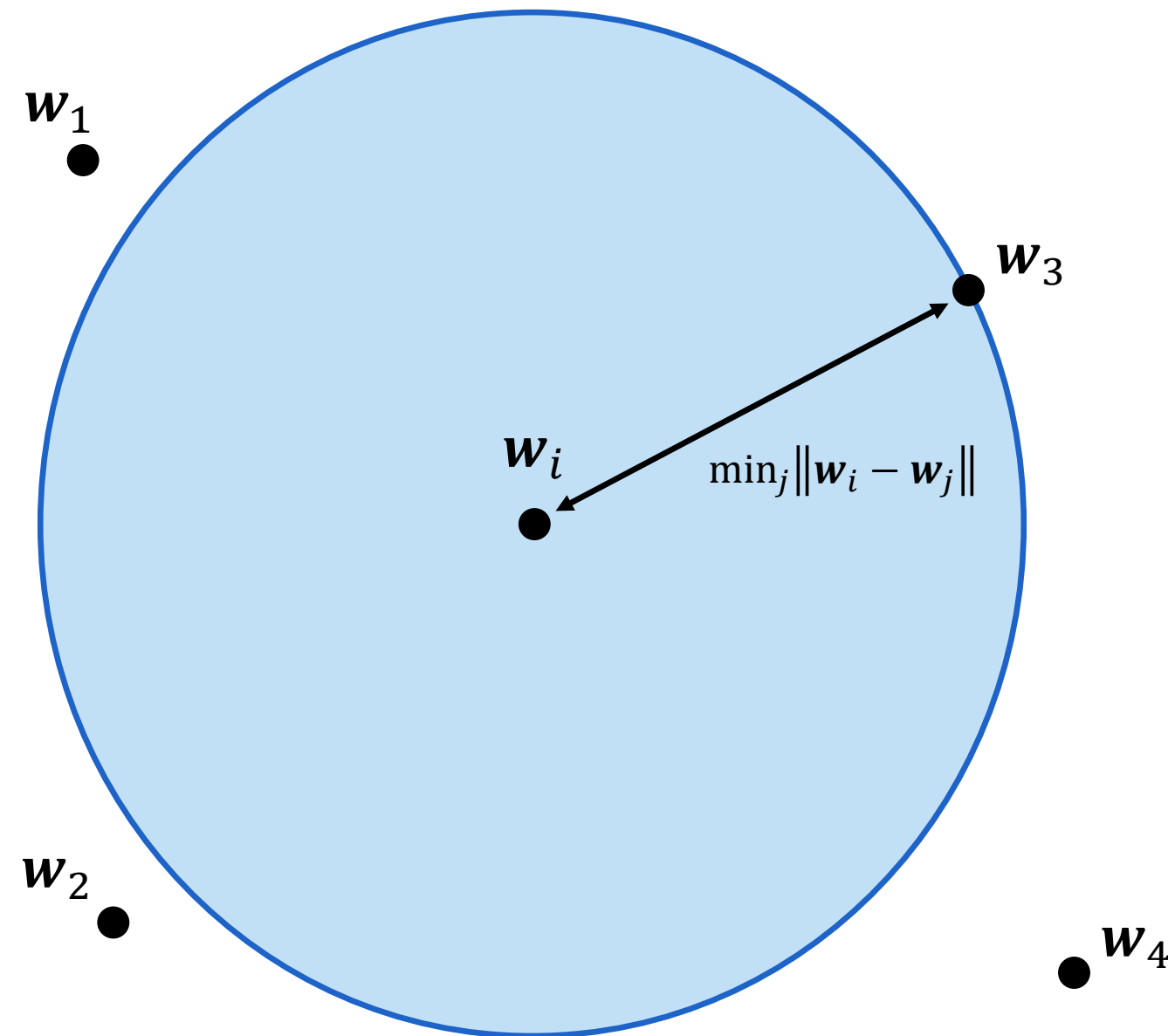
What does this mean?

$$Q_i^{\pi_i^*}(s, a) - Q_i^{\pi}(s, a) \leq \frac{2}{1-\gamma} \phi_{\max} \min_j \|w_i - w_j\|$$

In this example:

$j \in \{1, 2, 3, 4\}$

$w_i, w_j \in \mathbb{R}^2$



π performance on task M_i will **not** be close to π_i^*

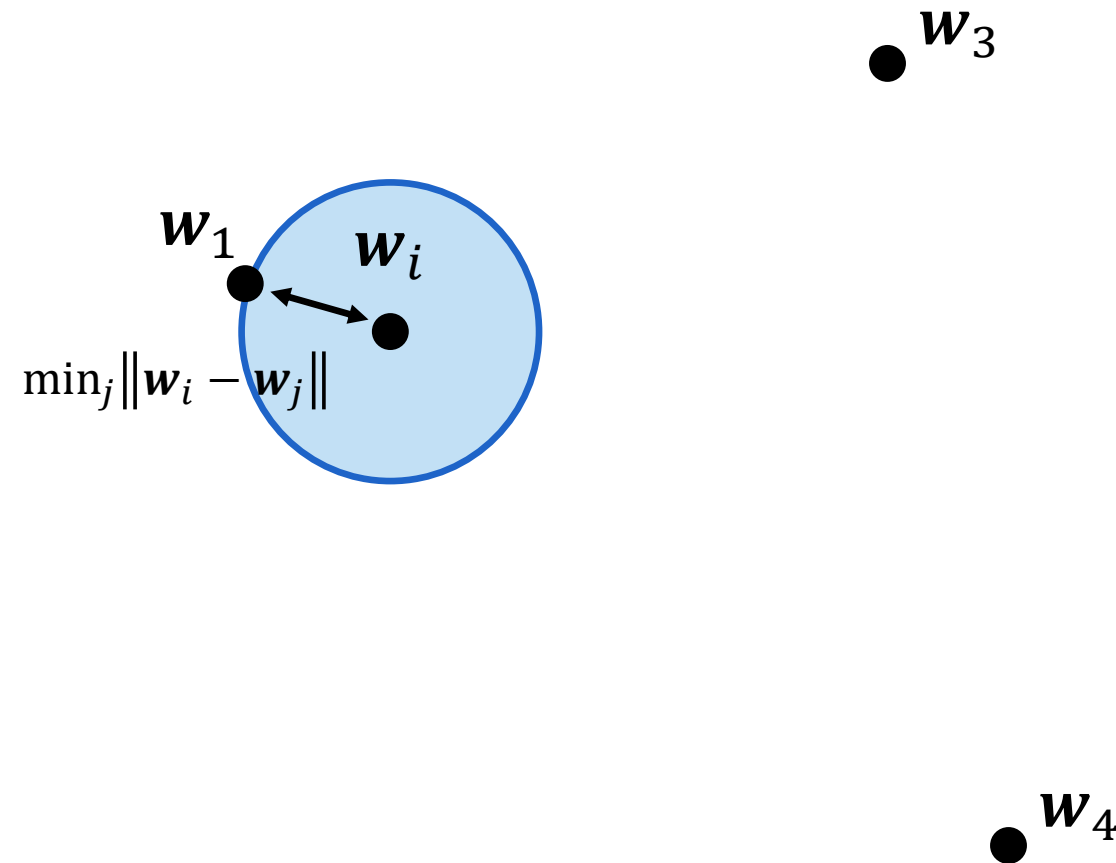
What does this mean?

$$Q_i^{\pi_i^*}(s, a) - Q_i^{\pi}(s, a) \leq \frac{2}{1 - \gamma} \phi_{\max} \min_j \|\mathbf{w}_i - \mathbf{w}_j\|$$

In this example:

$$j \in \{1, 2, 3, 4\}$$

$$\mathbf{w}_i, \mathbf{w}_j \in \mathbb{R}^2$$



π performance on task M_i will be close to π_i^*

Transfer via SFs – In practice

The two versions of SFQL

1. SFQL- ϕ

- assumption:
 - ϕ is known

2. SFQL- h

- assumption:
 - ϕ is *not* known
 - $\phi \in \mathbb{R}^h$

SFQL- ϕ

Given initial tasks $M_1, \dots, M_n \in \mathcal{M}^\phi$, learn and store

$$\left\{ \boldsymbol{\psi}^{\pi_j^\star}, \mathbf{w}_j \right\}_{j=1, \dots, n}$$

Given a new task $M_i \in \mathcal{M}^\phi$:

1. Learn \mathbf{w}_i .
2. Compute π using the GPI with SFs Theorem.
 - This becomes *trivial* using SFs

$$Q_i^{\pi_j^\star}(s, a) = \boldsymbol{\psi}^{\pi_j^\star}(s, a)^\top \mathbf{w}_i$$

SFQL- ϕ

Given initial tasks $M_1, \dots, M_n \in \mathcal{M}^\phi$, learn and store

$$\left\{ \boldsymbol{\psi}^{\pi_j^\star}, \mathbf{w}_j \right\}_{j=1, \dots, n}$$

Given a new task $M_i \in \mathcal{M}^\phi$:

1. Learn \mathbf{w}_i .
2. Compute π using the GPI with SFs Theorem.
3. Learn $\boldsymbol{\psi}^{\pi_i^\star}$ starting from π .
4. Store $\left\{ \boldsymbol{\psi}^{\pi_i^\star}, \mathbf{w}_i \right\}$.

SFQL- ϕ – Learning ψ^π

It is easy to see that

$$\psi^\pi(s, a) = \phi_{t+1} + \gamma \mathbb{E}^\pi [\psi^\pi(s_{t+1}, \pi(s_{t+1})) \mid s_t = s, a_t = a]$$

- SFs satisfy a Bellman equation.
- Any RL method can be used.
- The paper uses Q-Learning (QL), hence the name SFQL.

SFQL- ϕ – Learning w

Since we assume ϕ is known and

$$r(s, a, s') = \phi(s, a, s')^\top \mathbf{w}$$

This is an easy supervised learning problem.

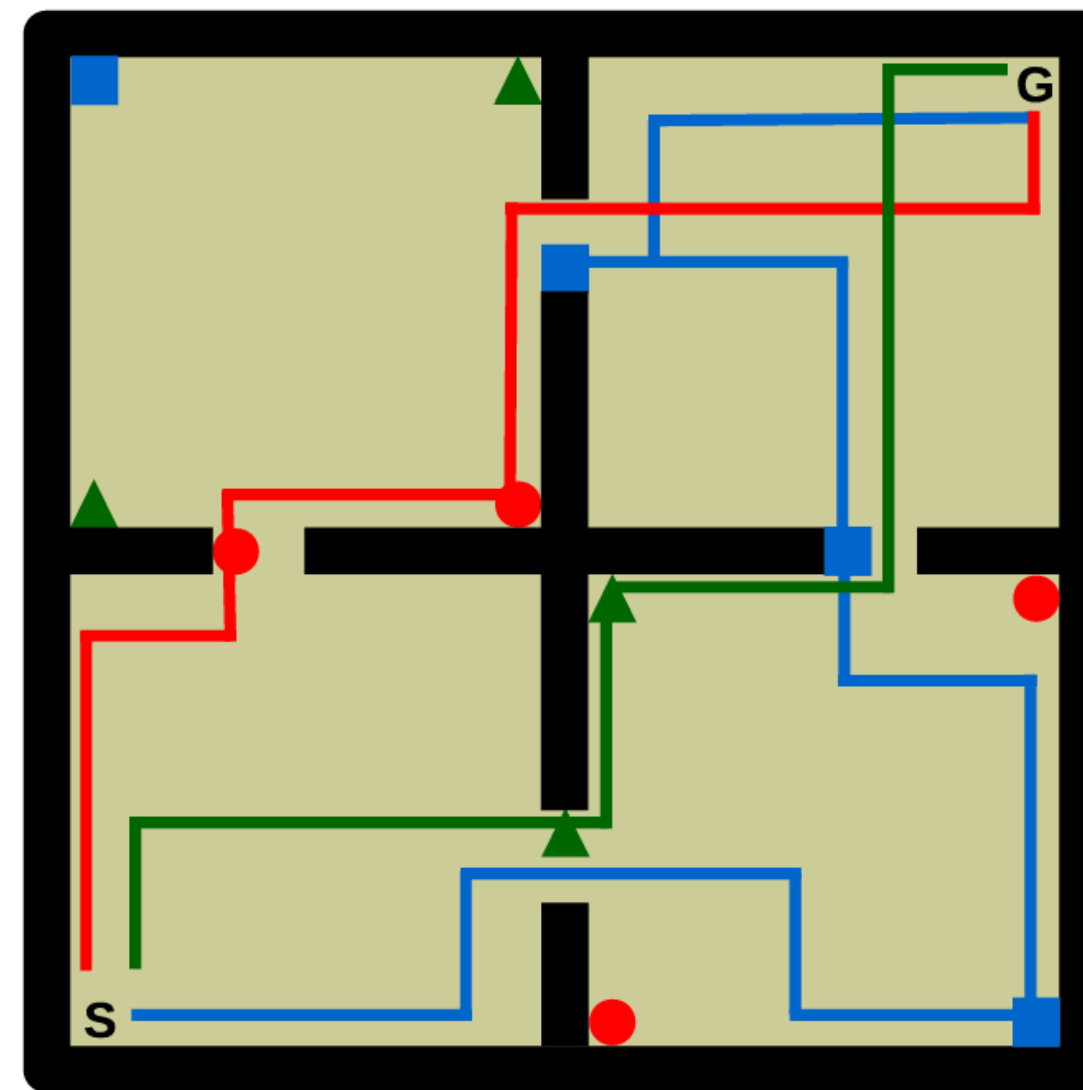
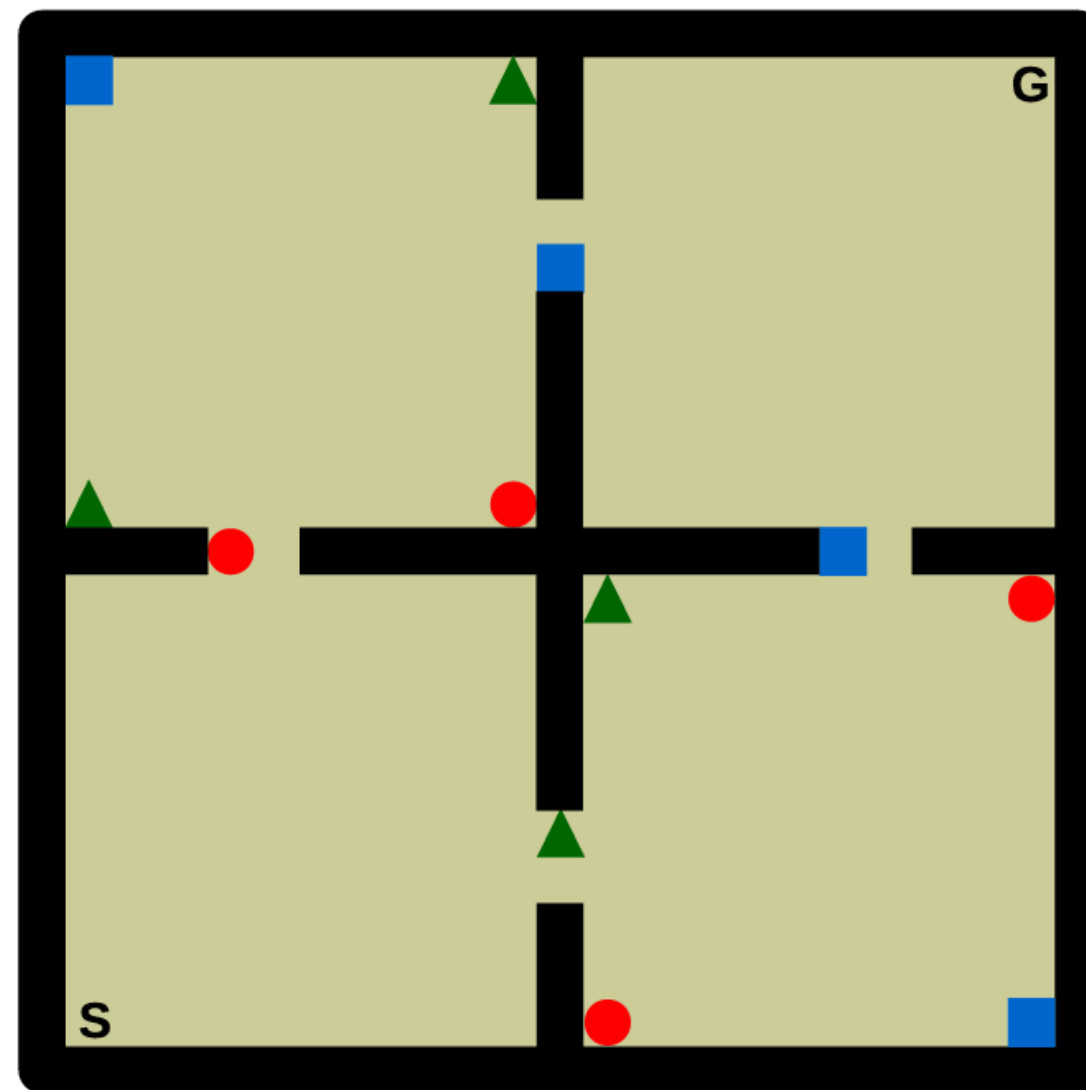
SFQL-*h*

Then by using the samples collected by QL in the first 20 tasks, we simultaneously learn $\mathbf{w}_1, \dots, \mathbf{w}_{20}$ and ϕ .

Experiment

Environment layout

Layout and optimal trajectories associated with specific tasks (i.e. colors).



I Expected step size

Results

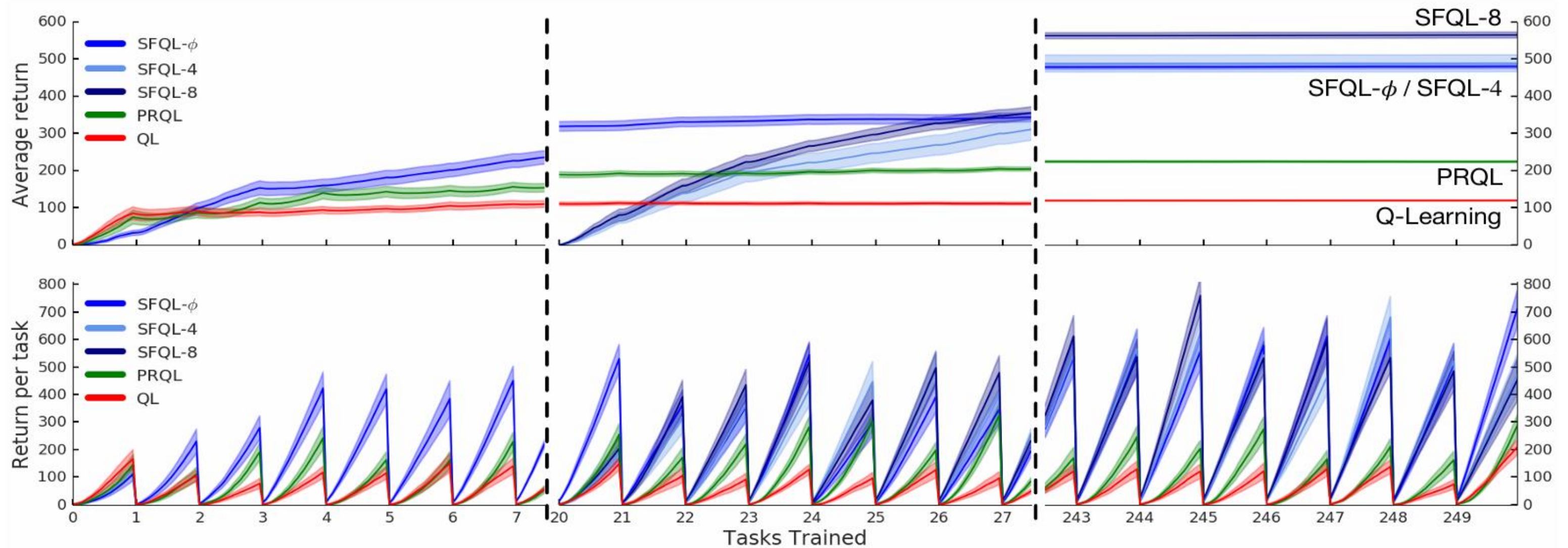


Figure 2: Average and cumulative return per task in the four-room domain. SFQL- h receives no reward during the first 20 tasks while learning $\tilde{\phi}$. Error-bands show one standard error over 30 runs.

Results

- All versions of SFQL significantly outperform the other two methods.
- SFQL- h seems to achieve good overall performance faster than SFQL- ϕ .
 - One possible explanation could be that the learnt ϕ may be less sparse than the real one, thus facilitating learning.

Thank you for your
attention