

# Risk-Aware Transfer in Reinforcement Learning using Successor Features

---

Michael Gimelfarb<sup>1,3</sup>   André Barreto<sup>2</sup>   Scott Sanner<sup>1,3</sup>   Chi-Guhn Lee<sup>1</sup>

<sup>1</sup>University of Toronto

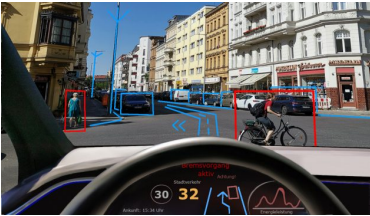
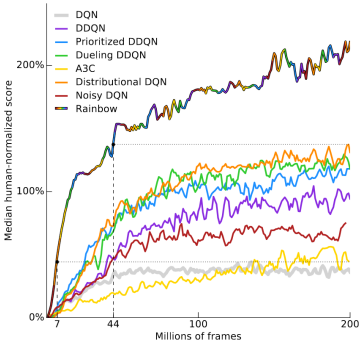
<sup>2</sup>DeepMind

<sup>3</sup>Vector Institute (Affiliate Program)

# Introduction



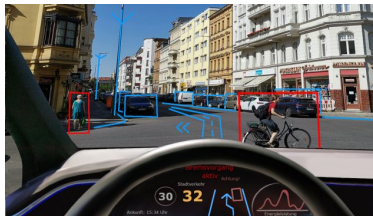
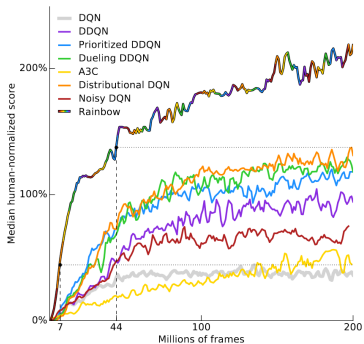
# Motivation



**Figure 1:** Sample Efficiency (Source: original paper on Rainbow DQN)

**Figure 2:** Risk-Awareness (Source: Wikimedia Commons)

# Motivation



**Figure 1:** Sample Efficiency (Source: original paper on Rainbow DQN)

**Figure 2:** Risk-Awareness (Source: Wikimedia Commons)

- **transfer learning**

- replace  $\mathbb{E}[\cdot]$  by non-linear **utility**  $\mathcal{U}[\cdot]$

# Motivation

---

Our goals:

# Motivation

---

Our goals:

- transfer between tasks with **shared dynamics and different goals**

# Motivation

---

Our goals:

- transfer between tasks with **shared dynamics and different goals**
- borrow **GPI/GPE** (e.g. successor features) from the risk-neutral setting<sup>1</sup>

<sup>1</sup>Barreto, André, et al. "Successor Features for Transfer in Reinforcement Learning." NIPS. 2017.



# Motivation

---

Our goals:

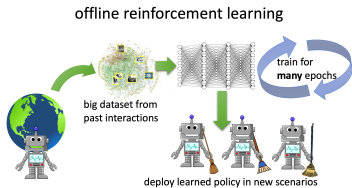
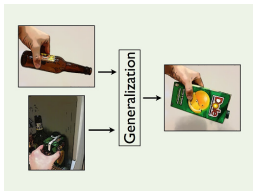
- transfer between tasks with **shared dynamics and different goals**
- borrow **GPI/GPE** (e.g. successor features) from the risk-neutral setting<sup>1</sup>
- provide **task generalization** by exploiting the structure of the task/reward space

<sup>1</sup>Barreto, André, et al. "Successor Features for Transfer in Reinforcement Learning." NIPS. 2017.

# Motivation

Our goals:

- transfer between tasks with **shared dynamics and different goals**
- borrow **GPI/GPE** (e.g. successor features) from the risk-neutral setting<sup>1</sup>
- provide **task generalization** by exploiting the structure of the task/reward space
- design a method suitable for offline RL<sup>2</sup>



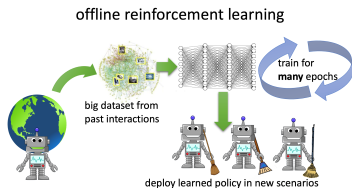
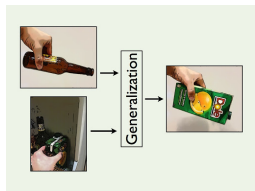
<sup>1</sup>Barreto, André, et al. "Successor Features for Transfer in Reinforcement Learning." NIPS. 2017.

<sup>2</sup>Levine, Sergey, et al. "Offline reinforcement learning..." arXiv. 2020.

# Motivation

Our goals:

- transfer between tasks with **shared dynamics and different goals**
- borrow **GPI/GPE** (e.g. successor features) from the risk-neutral setting<sup>1</sup>
- provide **task generalization** by exploiting the structure of the task/reward space
- design a method suitable for offline RL<sup>2</sup>



- incorporate risk explicitly by e.g. penalizing the variance of returns

<sup>1</sup>Barreto, André, et al. "Successor Features for Transfer in Reinforcement Learning." NIPS. 2017.

<sup>2</sup>Levine, Sergey, et al. "Offline reinforcement learning..." arXiv. 2020.

Introduce **Risk-Aware Successor Features** (RaSF)

# Motivation

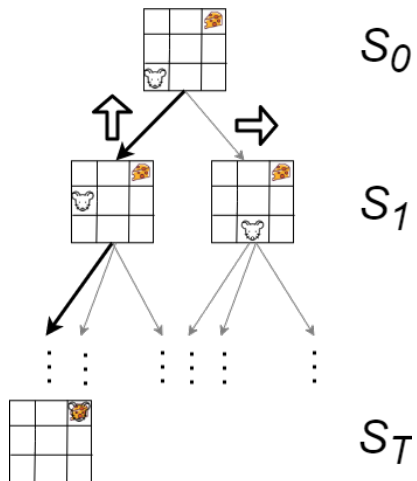
---

Introduce **Risk-Aware Successor Features** (RaSF)

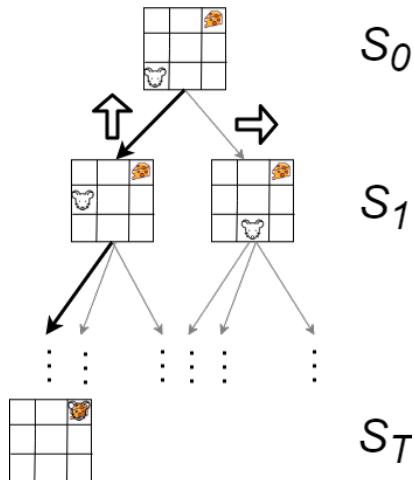
	Transfers Skills	Exploits Task Structure	Risk-Sensitive
Risk-Aware RL	✗	✗	✓
Risk-Aware Transfer	✓	✗	✓
Successor Features	✓	✓	✗
<b>RaSF (Ours)</b>	✓	✓	✓

## **Preliminaries – Successor Features**

# Policy Evaluation



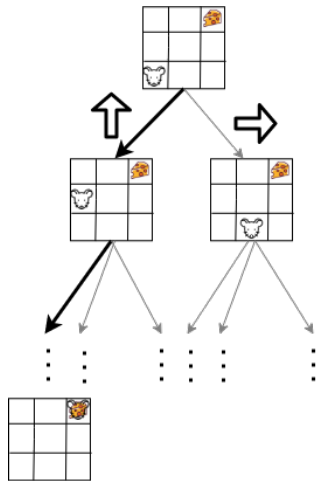
## Policy Evaluation



- $\mathbb{E}[R_0 + \gamma R_1 + \dots]$  requires averaging all possible future outcomes of the world – **curse of dimensionality**



## Policy Evaluation



$S_0$

- $\mathbb{E}[R_0 + \gamma R_1 + \dots]$  requires averaging all possible future outcomes of the world – **curse of dimensionality**

$S_1$

- use cached  $Q(s', a')$  in each successor state to **bootstrap** the estimated Q-values in state  $s$

$$Q(s, a) = \mathbb{E}_{S' \sim P(\cdot | s, a)} [R_t + \gamma Q(S', \pi(S'))]$$

$S_T$



## Policy Improvement

---

Suppose an initial policy  $\pi$  is given:

# Policy Improvement

---

Suppose an initial policy  $\pi$  is given:

- compute the value of  $\pi$ , e.g.

$$Q^\pi(s, a) = \mathbb{E}_{S' \sim P(\cdot | s, a)} [R_t^\pi + \gamma Q(S', \pi(S'))]$$

◁ policy evaluation

# Policy Improvement

---

Suppose an initial policy  $\pi$  is given:

- compute the value of  $\pi$ , e.g.

$$Q^\pi(s, a) = \mathbb{E}_{S' \sim P(\cdot | s, a)} [R_t^\pi + \gamma Q(S', \pi(S'))] \quad \triangleleft \text{policy evaluation}$$

- construct a new policy  $\pi'$  according to

$$\pi'(s) \in \arg \max_{a \in \mathcal{A}} Q^\pi(s, a) \quad \triangleleft \text{policy improvement}$$

# Policy Improvement

---

Suppose an initial policy  $\pi$  is given:

- compute the value of  $\pi$ , e.g.

$$Q^\pi(s, a) = \mathbb{E}_{S' \sim P(\cdot|s, a)}[R_t^\pi + \gamma Q(S', \pi(S'))] \quad \triangleleft \text{policy evaluation}$$

- construct a new policy  $\pi'$  according to

$$\pi'(s) \in \arg \max_{a \in \mathcal{A}} Q^\pi(s, a) \quad \triangleleft \text{policy improvement}$$

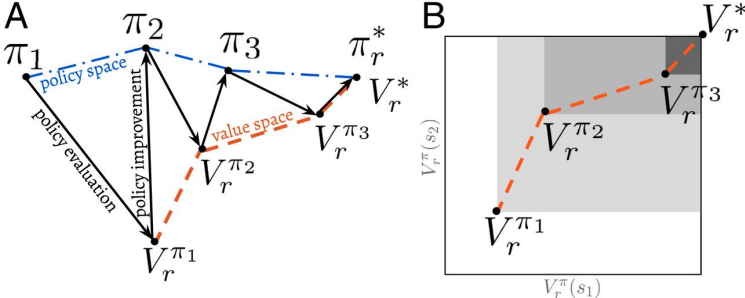
**Policy Improvement Theorem:**  $\pi'$  is “better” than  $\pi$ , e.g.  $Q^{\pi'}(s, a) \geq Q^\pi(s, a)$

# Generalized Policy Iteration

---

# Generalized Policy Iteration

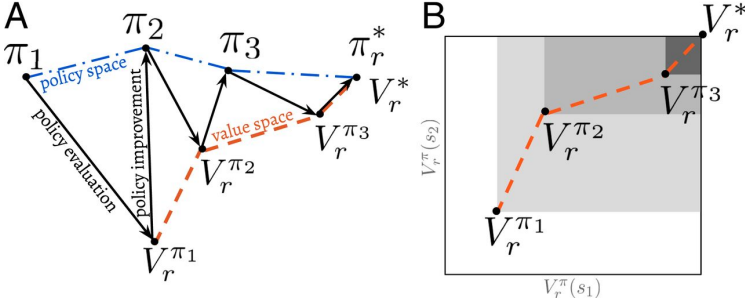
Alternating between evaluation and improvement leads to an optimal policy





# Generalized Policy Iteration

Alternating between evaluation and improvement leads to an optimal policy



**Key Idea:** Replace  $\pi$  with multiple source policies  $\pi_1 \dots \pi_n$ .

# Generalized Policy Improvement 2.0

---

## Generalized Policy Improvement 2.0

---

Suppose  $\pi_1, \dots, \pi_n$  are given<sup>1</sup>:

<sup>1</sup>Barreto, André, et al. "Successor Features for Transfer in Reinforcement Learning." NeurIPS 2017.

## Generalized Policy Improvement 2.0

---

Suppose  $\pi_1, \dots, \pi_n$  are given<sup>1</sup>:

- compute the values

$$Q^{\pi_1}(s, a), \dots, Q^{\pi_n}(s, a)$$

◁ generalized policy evaluation (GPE)

<sup>1</sup>Barreto, André, et al. "Successor Features for Transfer in Reinforcement Learning." NeurIPS 2017.

## Generalized Policy Improvement 2.0

---

Suppose  $\pi_1, \dots, \pi_n$  are given<sup>1</sup>:

- compute the values

$$Q^{\pi_1}(s, a), \dots, Q^{\pi_n}(s, a) \quad \triangleleft \text{generalized policy evaluation (GPE)}$$

- pick the policy with the best return

$$i^* \in \arg \max_i Q^{\pi_i}(s, a)$$

<sup>1</sup>Barreto, André, et al. "Successor Features for Transfer in Reinforcement Learning." NeurIPS 2017.

## Generalized Policy Improvement 2.0

Suppose  $\pi_1, \dots, \pi_n$  are given<sup>1</sup>:

- compute the values

$$Q^{\pi_1}(s, a), \dots, Q^{\pi_n}(s, a) \quad \triangleleft \text{generalized policy evaluation (GPE)}$$

- pick the policy with the best return

$$i^* \in \arg \max_i Q^{\pi_i}(s, a)$$

- construct  $\pi'$  as usual but w.r.t. the “best” policy  $\pi_{i^*}$

$$\pi'(s) \in \arg \max_{a \in \mathcal{A}} Q^{\pi_{i^*}}(s, a) = \arg \max_{a \in \mathcal{A}} \max_{i=1 \dots n} Q^{\pi_i}(s, a) \quad \triangleleft \text{generalized policy improvement (GPI)}$$

<sup>1</sup>Barreto, André, et al. "Successor Features for Transfer in Reinforcement Learning." NeurIPS 2017.

## Generalized Policy Improvement 2.0

Suppose  $\pi_1, \dots, \pi_n$  are given<sup>1</sup>:

- compute the values

$$Q^{\pi_1}(s, a), \dots, Q^{\pi_n}(s, a) \quad \triangleleft \text{generalized policy evaluation (GPE)}$$

- pick the policy with the best return

$$i^* \in \arg \max_i Q^{\pi_i}(s, a)$$

- construct  $\pi'$  as usual but w.r.t. the “best” policy  $\pi_{i^*}$

$$\pi'(s) \in \arg \max_{a \in \mathcal{A}} Q^{\pi_{i^*}}(s, a) = \arg \max_{a \in \mathcal{A}} \max_{i=1 \dots n} Q^{\pi_i}(s, a) \quad \triangleleft \text{generalized policy improvement (GPI)}$$

**Key result:**  $\pi'$  is better than  $\pi_{i^*}$ .

<sup>1</sup>Barreto, André, et al. "Successor Features for Transfer in Reinforcement Learning." NeurIPS 2017.

# **Preliminaries – Risk-Aversion in MDPs using Entropic Utility Functions**





## Time-Consistency

---

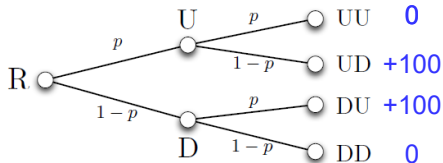
Optimizing risk measures in sequential problems is hard<sup>2</sup>:

<sup>2</sup>Chow, Yin-Lam, and Marco Pavone. "A framework for time-consistent, risk-averse model predictive control: Theory and algorithms." 2014 American Control Conference. IEEE, 2014.

## Time-Consistency

Optimizing risk measures in sequential problems is hard<sup>2</sup>:

- considering optimizing a final cost  $Z$ :

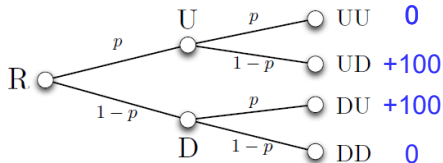


<sup>2</sup>Chow, Yin-Lam, and Marco Pavone. "A framework for time-consistent, risk-averse model predictive control: Theory and algorithms." 2014 American Control Conference. IEEE, 2014.

## Time-Consistency

Optimizing risk measures in sequential problems is hard<sup>2</sup>:

- considering optimizing a final cost  $Z$ :



- consider the dynamic risk measure

$$\rho_{k,N}(Z) = \max_{p \in \{0.4, 0.6\}} \mathbb{E}_p[Z | \mathcal{F}_k],$$

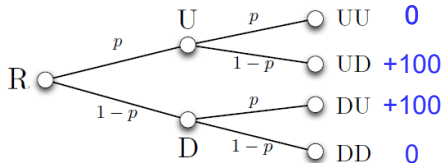
for  $k = 0, 1, 2$

<sup>2</sup>Chow, Yin-Lam, and Marco Pavone. "A framework for time-consistent, risk-averse model predictive control: Theory and algorithms." 2014 American Control Conference. IEEE, 2014.

## Time-Consistency

Optimizing risk measures in sequential problems is hard<sup>2</sup>:

- considering optimizing a final cost  $Z$ :



- it is easy to check that  $\rho_1(Z)(\omega) = 60$  for all  $\omega$ , so  $Z$  is riskier than a deterministic cash flow of  $W = 50$  at time 1

- consider the dynamic risk measure

$$\rho_{k,N}(Z) = \max_{p \in \{0.4, 0.6\}} \mathbb{E}_p[Z | \mathcal{F}_k],$$

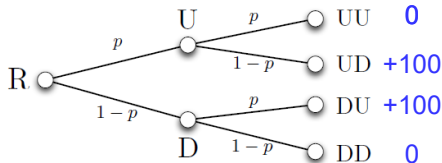
for  $k = 0, 1, 2$

<sup>2</sup>Chow, Yin-Lam, and Marco Pavone. "A framework for time-consistent, risk-averse model predictive control: Theory and algorithms." 2014 American Control Conference. IEEE, 2014.

## Time-Consistency

Optimizing risk measures in sequential problems is hard<sup>2</sup>:

- considering optimizing a final cost  $Z$ :



- consider the dynamic risk measure

$$\rho_{k,N}(Z) = \max_{p \in \{0.4, 0.6\}} \mathbb{E}_p[Z | \mathcal{F}_k],$$

for  $k = 0, 1, 2$

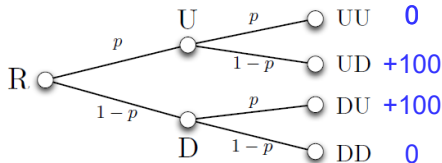
<sup>2</sup>Chow, Yin-Lam, and Marco Pavone. "A framework for time-consistent, risk-averse model predictive control: Theory and algorithms." 2014 American Control Conference. IEEE, 2014.

- it is easy to check that  $\rho_1(Z)(\omega) = 60$  for all  $\omega$ , so  $Z$  is riskier than a deterministic cash flow of  $W = 50$  at time 1
- yet,  $\rho_0(Z)(\omega) = 40$  and so  $Z$  is less risky than  $W$  at time 0!

## Time-Consistency

Optimizing risk measures in sequential problems is hard<sup>2</sup>:

- considering optimizing a final cost  $Z$ :



- consider the dynamic risk measure

$$\rho_{k,N}(Z) = \max_{p \in \{0.4, 0.6\}} \mathbb{E}_p[Z | \mathcal{F}_k],$$

for  $k = 0, 1, 2$

<sup>2</sup>Chow, Yin-Lam, and Marco Pavone. "A framework for time-consistent, risk-averse model predictive control: Theory and algorithms." 2014 American Control Conference. IEEE, 2014.

- it is easy to check that  $\rho_1(Z)(\omega) = 60$  for all  $\omega$ , so  $Z$  is riskier than a deterministic cash flow of  $W = 50$  at time 1
- yet,  $\rho_0(Z)(\omega) = 40$  and so  $Z$  is less risky than  $W$  at time 0!

**Key idea:**  $Z$  has become riskier just because time has passed!

## Variance Reduction in MDPs with Entropic Utilities

---



## Variance Reduction in MDPs with Entropic Utilities

---

We use the **entropic utility** to measure risk:

## Variance Reduction in MDPs with Entropic Utilities

---

We use the **entropic utility** to measure risk:

- defined in terms of the moment-generating function

$$U_{\beta}[R] = \frac{1}{\beta} \log \mathbb{E} \left[ e^{\beta R} \right]$$

## Variance Reduction in MDPs with Entropic Utilities

---

We use the **entropic utility** to measure risk:

- defined in terms of the moment-generating function

$$U_{\beta}[R] = \frac{1}{\beta} \log \mathbb{E} \left[ e^{\beta R} \right]$$

- has the Taylor expansion

$$U_{\beta}[R] = \mathbb{E}[R] + \frac{\beta}{2} \mathbb{V}[R] + O(\beta^2) \quad \triangleleft \beta \text{ is a level of risk-aversion}$$

## Variance Reduction in MDPs with Entropic Utilities

---

We use the **entropic utility** to measure risk:

- defined in terms of the moment-generating function

$$U_{\beta}[R] = \frac{1}{\beta} \log \mathbb{E} \left[ e^{\beta R} \right]$$

- has the Taylor expansion

$$U_{\beta}[R] = \mathbb{E}[R] + \frac{\beta}{2} \mathbb{V}[R] + O(\beta^2) \quad \triangleleft \beta \text{ is a level of risk-aversion}$$

- connected to the mean-variance optimization in MDPs<sup>1</sup>

<sup>1</sup>Mannor, Shie, and John N. Tsitsiklis. "Mean-variance optimization in Markov decision processes." ICML. 2011.

## Variance Reduction in MDPs with Entropic Utilities

---

We incorporate entropic utility into MDPs:

# Variance Reduction in MDPs with Entropic Utilities

We incorporate entropic utility into MDPs:

- **dynamic programming:** has a Bellman equation formulation

$$\mathcal{Q}_{h,\beta}^{\pi}(s, a) = U_{\beta} \left[ \sum_{t=h}^T r(s_t, \pi_t(s_t), s_{t+1}) \right] \iff \begin{aligned} \mathcal{Q}_{h,\beta}^{\pi}(s, a) &= U_{\beta} [r(s, a, s') + \mathcal{Q}_{h+1,\beta}^{\pi}(s', \pi_{h+1}(s'))] \\ &= \frac{1}{\beta} \log \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ e^{\beta \{r(s, a, s') + \mathcal{Q}_{h+1,\beta}^{\pi}(s', \pi_{h+1}(s'))\}} \right]. \end{aligned}$$

# Variance Reduction in MDPs with Entropic Utilities

We incorporate entropic utility into MDPs:

- **dynamic programming:** has a Bellman equation formulation

$$\mathcal{Q}_{h,\beta}^{\pi}(s, a) = U_{\beta} \left[ \sum_{t=h}^T r(s_t, \pi_t(s_t), s_{t+1}) \right] \iff \begin{aligned} \mathcal{Q}_{h,\beta}^{\pi}(s, a) &= U_{\beta} [r(s, a, s') + \mathcal{Q}_{h+1,\beta}^{\pi}(s', \pi_{h+1}(s'))] \\ &= \frac{1}{\beta} \log \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ e^{\beta \{r(s, a, s') + \mathcal{Q}_{h+1,\beta}^{\pi}(s', \pi_{h+1}(s'))\}} \right]. \end{aligned}$$

- **recursive property:** behaves similar to expectation in total-reward episodic MDPs<sup>1</sup> (and discounted MDPs with simple modifications)

<sup>1</sup>Osogami, Takayuki. "Robustness and risk-sensitivity in Markov decision processes." NeurIPS 2012.

# Variance Reduction in MDPs with Entropic Utilities

We incorporate entropic utility into MDPs:

- **dynamic programming:** has a Bellman equation formulation

$$\boxed{Q_{h,\beta}^{\pi}(s, a) = U_{\beta} \left[ \sum_{t=h}^T r(s_t, \pi_t(s_t), s_{t+1}) \right]} \longleftrightarrow \boxed{\begin{aligned} Q_{h,\beta}^{\pi}(s, a) &= U_{\beta} [r(s, a, s') + Q_{h+1,\beta}^{\pi}(s', \pi_{h+1}(s'))] \\ &= \frac{1}{\beta} \log \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ e^{\beta \{r(s, a, s') + Q_{h+1,\beta}^{\pi}(s', \pi_{h+1}(s'))\}} \right]. \end{aligned}}$$

- **recursive property:** behaves similar to expectation in total-reward episodic MDPs<sup>1</sup> (and discounted MDPs with simple modifications)
- **convex/concave:** satisfies properties that can be seen as rational decision making

<sup>1</sup>Osogami, Takayuki. "Robustness and risk-sensitivity in Markov decision processes." NeurIPS 2012.



# Variance Reduction in MDPs with Entropic Utilities

We incorporate entropic utility into MDPs:

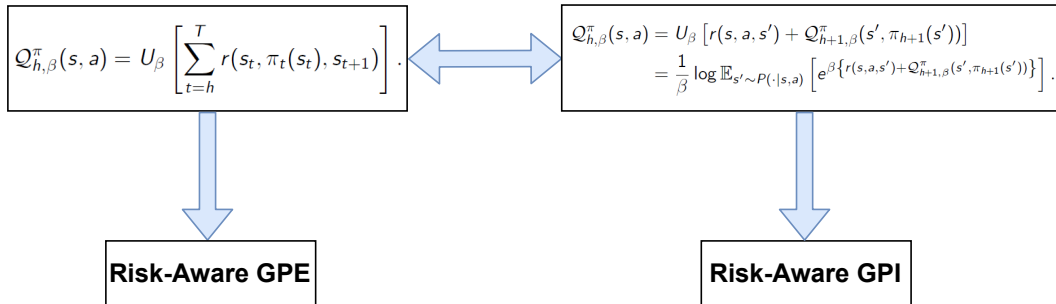
- **dynamic programming:** has a Bellman equation formulation

$$\mathcal{Q}_{h,\beta}^{\pi}(s, a) = U_{\beta} \left[ \sum_{t=h}^T r(s_t, \pi_t(s_t), s_{t+1}) \right] \iff \begin{aligned} \mathcal{Q}_{h,\beta}^{\pi}(s, a) &= U_{\beta} [r(s, a, s') + \mathcal{Q}_{h+1,\beta}^{\pi}(s', \pi_{h+1}(s'))] \\ &= \frac{1}{\beta} \log \mathbb{E}_{s' \sim P(\cdot | s, a)} \left[ e^{\beta \{r(s, a, s') + \mathcal{Q}_{h+1,\beta}^{\pi}(s', \pi_{h+1}(s'))\}} \right]. \end{aligned}$$

- **recursive property:** behaves similar to expectation in total-reward episodic MDPs<sup>1</sup> (and discounted MDPs with simple modifications)
- **convex/concave:** satisfies properties that can be seen as rational decision making
- **time consistency:** can focus on Markov policies

<sup>1</sup>Osogami, Takayuki. "Robustness and risk-sensitivity in Markov decision processes." NeurIPS 2012.

# Roadmap



# Theory

## Motivating Example

---

## Motivating Example

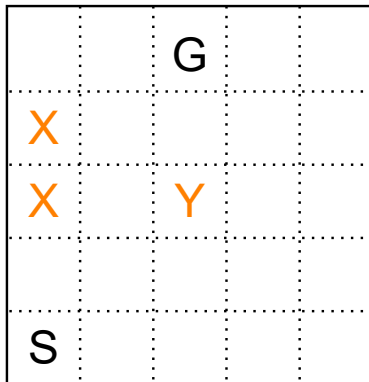
---

Why is the problem non-trivial?

## Motivating Example

---

Why is the problem non-trivial?



## Motivating Example

---

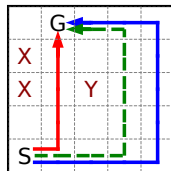
1. Define a family of tasks:
  - two source tasks:  
low failure cost + high failure cost
  - one target task:  
only X has high failure cost

## Motivating Example

1. Define a family of tasks:

- two source tasks:  
low failure cost + high failure cost
- one target task:  
only X has high failure cost

2. Solve them with VI for **fixed**  $\beta$



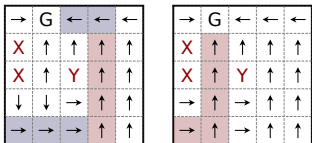


# Motivating Example

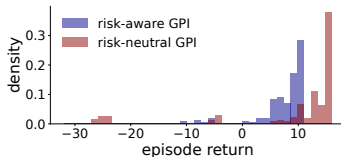
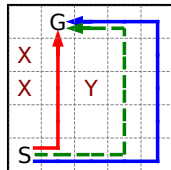
1. Define a family of tasks:

- two source tasks:  
low failure cost + high failure cost
- one target task:  
only X has high failure cost

3. Apply risk-aware and risk-neutral GPI:



2. Solve them with VI for **fixed**  $\beta$

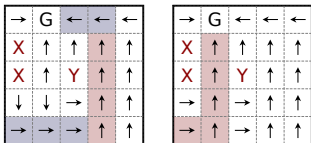


# Motivating Example

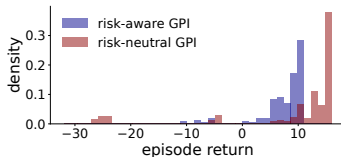
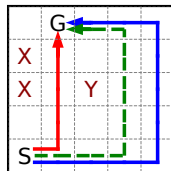
1. Define a family of tasks:

- two source tasks:  
low failure cost + high failure cost
- one target task:  
only X has high failure cost

3. Apply risk-aware and risk-neutral GPI:



2. Solve them with VI for **fixed**  $\beta$



**Conclusion:** only risk-aware GPI results in the correct target policy

## Key Theoretical Results

---

## Key Theoretical Results

---

Armed with this knowledge, we prove that risk-aware GPI:

## Key Theoretical Results

---

Armed with this knowledge, we prove that risk-aware GPI:

- is a strict policy improvement operator

**Theorem 1 (GPI for Entropic Utility).** *Let  $\pi_1, \dots, \pi_n$  be arbitrary deterministic Markov policies with utilities  $\tilde{Q}_{h,\beta}^{\pi_1}, \dots, \tilde{Q}_{h,\beta}^{\pi_n}$  evaluated in an arbitrary task  $M$ , such that  $|\tilde{Q}_{h,\beta}^{\pi_i}(s, a) - Q_{h,\beta}^{\pi_i}(s, a)| \leq \varepsilon$  for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}$ ,  $i = 1 \dots n$  and  $h \in \mathcal{T}$ . Define*

$$\pi_h(s) \in \arg \max_{a \in \mathcal{A}} \max_{i=1 \dots n} \tilde{Q}_{h,\beta}^{\pi_i}(s, a), \quad \forall s \in \mathcal{S}. \quad (4)$$

Then,

$$Q_{h,\beta}^{\pi}(s, a) \geq \max_i Q_{h,\beta}^{\pi_i}(s, a) - 2(T - h + 1)\varepsilon, \quad h \leq T.$$

## Key Theoretical Results

Armed with this knowledge, we prove that risk-aware GPI:

- is a strict policy improvement operator

**Theorem 1 (GPI for Entropic Utility).** *Let  $\pi_1, \dots, \pi_n$  be arbitrary deterministic Markov policies with utilities  $\tilde{Q}_{h,\beta}^{\pi_1}, \dots, \tilde{Q}_{h,\beta}^{\pi_n}$  evaluated in an arbitrary task  $M$ , such that  $|\tilde{Q}_{h,\beta}^{\pi_i}(s, a) - Q_{h,\beta}^{\pi_i}(s, a)| \leq \varepsilon$  for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}$ ,  $i = 1 \dots n$  and  $h \in \mathcal{T}$ . Define*

$$\pi_h(s) \in \arg \max_{a \in \mathcal{A}} \max_{i=1 \dots n} \tilde{Q}_{h,\beta}^{\pi_i}(s, a), \quad \forall s \in \mathcal{S}. \quad (4)$$

Then,

$$Q_{h,\beta}^{\pi}(s, a) \geq \max_i Q_{h,\beta}^{\pi_i}(s, a) - 2(T - h + 1)\varepsilon, \quad h \leq T.$$

- is optimal up to an irreducible task discrepancy gap

**Theorem 2.** *Let  $Q_{h,\beta}^{\pi_i^*}$  be the utilities of optimal Markov policies  $\pi_i^*$  from task  $M_i$  but evaluated in task  $M$  with reward function  $r(s, a, s')$ . Furthermore, let  $\tilde{Q}_{h,\beta}^{\pi_i^*}$  be such that  $|\tilde{Q}_{h,\beta}^{\pi_i^*}(s, a) - Q_{h,\beta}^{\pi_i^*}(s, a)| < \varepsilon$  for all  $s \in \mathcal{S}$ ,  $a \in \mathcal{A}$ ,  $h \in \mathcal{T}$  and  $i = 1 \dots n$ , and let  $\pi$  be the corresponding policy in (4). Finally, let  $\delta_r = \min_{i=1 \dots n} \sup_{s,a,s'} |r(s, a, s') - r_i(s, a, s')|$ . Then,*

$$|Q_{h,\beta}^{\pi}(s, a) - Q_{h,\beta}^*(s, a)| \leq 2(T - h + 1)(\delta_r + \varepsilon), \quad h \leq T.$$



# Generalized Policy Evaluation

---

Assume **linear reward**:

$$r(s, a, s') = \phi(s, a, s')^T \mathbf{w}$$



# Generalized Policy Evaluation

Assume **linear reward**:

$$r(s, a, s') = \phi(s, a, s')^T \mathbf{w}$$

Now:

$$\begin{aligned} Q_{\mathbf{w}}^{\pi}(s, a) &= \mathbb{E} \left[ \sum_t \gamma^t R_t \mid S_0 = s, A_0 = a, A_t \sim \pi(S_t) \right] \\ &= \mathbb{E} \left[ \sum_t \gamma^t \phi_t^T \mathbf{w} \mid S_0 = s, A_0 = a, A_t \sim \pi(S_t) \right] \\ &= \mathbb{E} \left[ \underbrace{\sum_t \gamma^t \phi_t}_{\psi^{\pi}(s, a)} \mid S_0 = s, A_0 = a, A_t \sim \pi(S_t) \right]^T \mathbf{w} \end{aligned}$$



Can generalize GPE to **distributions of return**:

$$Q_{h,\beta}^{\pi}(s, a) = U_{\beta} \left[ \sum_{t=h}^T r(s_t, \pi_t(s_t), s_{t+1}) \right] = U_{\beta} [\Psi_h^{\pi}(s, a)^T \mathbf{w}]$$

Can generalize GPE to **distributions of return**:

$$Q_{h,\beta}^{\pi}(s, a) = U_{\beta} \left[ \sum_{t=h}^T r(s_t, \pi_t(s_t), s_{t+1}) \right] = U_{\beta} [\Psi_h^{\pi}(s, a)^T w]$$

One simple trick is to Taylor expand to the second moment:

$$\begin{aligned} U_{\beta} [\Psi_h^{\pi}(s, a)^T w] &= \mathbb{E}_P[\Psi_h^{\pi}(s, a)^T w] + \frac{\beta}{2} \text{Var}_P[\Psi_h^{\pi}(s, a)^T w] + O(\beta^2) \\ &\approx \psi_h^{\pi}(s, a)^T w + \frac{\beta}{2} w^T \text{Var}_P[\Psi_h^{\pi}(s, a)] w = \tilde{Q}_{h,\beta}^{\pi}(s, a) \end{aligned}$$

Can generalize GPE to **distributions of return**:

$$Q_{h,\beta}^{\pi}(s, a) = U_{\beta} \left[ \sum_{t=h}^T r(s_t, \pi_t(s_t), s_{t+1}) \right] = U_{\beta} [\Psi_h^{\pi}(s, a)^T w]$$

One simple trick is to Taylor expand to the second moment:

$$\begin{aligned} U_{\beta} [\Psi_h^{\pi}(s, a)^T w] &= \mathbb{E}_P[\Psi_h^{\pi}(s, a)^T w] + \frac{\beta}{2} \text{Var}_P[\Psi_h^{\pi}(s, a)^T w] + O(\beta^2) \\ &\approx \psi_h^{\pi}(s, a)^T w + \frac{\beta}{2} w^T \text{Var}_P[\Psi_h^{\pi}(s, a)] w = \tilde{Q}_{h,\beta}^{\pi}(s, a) \end{aligned}$$

Reduces to a (simpler) problem of estimating **sufficient statistics** of the feature occupancy

$\psi^{\pi_i}$  and  $\Sigma^{\pi_i}$

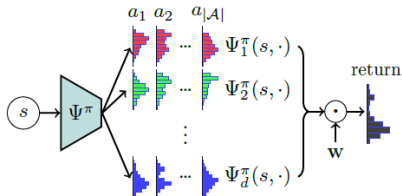
**Bellman principle**

**Distributional RL**

$$\bar{\Sigma}_h^{\pi}(s, a) = \mathbb{E}_{s' \sim P(\cdot|s, a)} [\bar{\delta}_h \bar{\delta}_h^{\top} + \bar{\Sigma}_{h+1}^{\pi}(s', \pi_{h+1}(s')) \mid s_h = s, a_h = a]$$

**Theorem 3 (Convergence of Covariance).** Let  $\|\cdot\|$  be a matrix-compatible norm, and suppose there exists  $\varepsilon : \mathcal{S} \times \mathcal{A} \times \mathcal{T} \rightarrow [0, \infty)$  such that  $\|\bar{\psi}_h^{\pi}(s, a) - \psi_h^{\pi}(s, a)\|^2 \leq \varepsilon_h(s, a)$  and  $\|\mathbb{E}_{s' \sim P(\cdot|s, a)}[\bar{\delta}_h(\bar{\psi}_h^{\pi}(s', \pi_{h+1}(s')) - \psi_h^{\pi}(s', \pi_{h+1}(s')))^{\top}]\| \leq \varepsilon_h(s, a)$ . Then,

$$\left\| \bar{\Sigma}_h^{\pi}(s, a) - \mathbb{E}_{s' \sim P(\cdot|s, a)} \left[ \bar{\delta}_h \bar{\delta}_h^{\top} + \bar{\Sigma}_{h+1}^{\pi}(s', \pi_{h+1}(s')) \right] \right\| \leq 3\varepsilon_h(s, a).$$



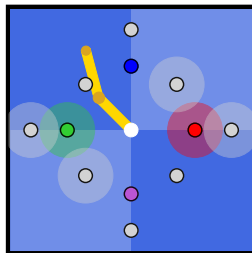
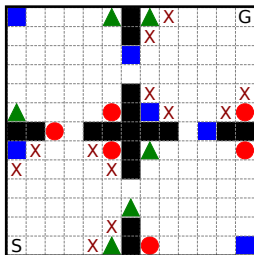
# Experiments





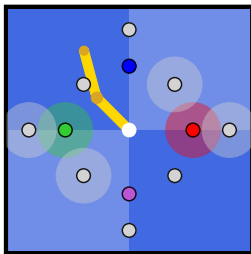
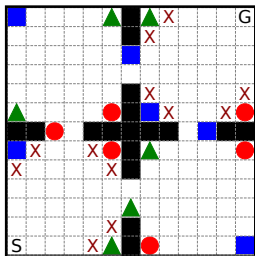
# Domains

Two domains from Barreto et al., 2017:



# Domains

Two domains from Barreto et al., 2017:



Introduce **reward volatility**:

- traps **X** for four-room
- action noise + danger zones for reacher



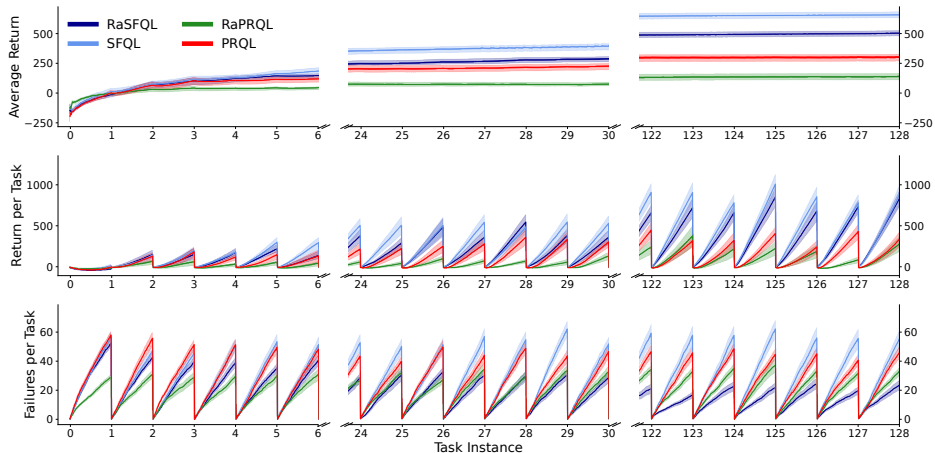
## Four-Room

---

Train on a sequence of 128 random task instances, for 20,000 steps each

# Four-Room

Train on a sequence of 128 random task instances, for 20,000 steps each



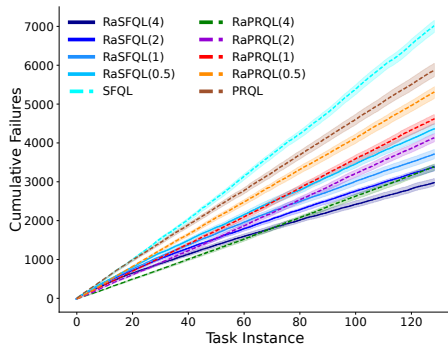
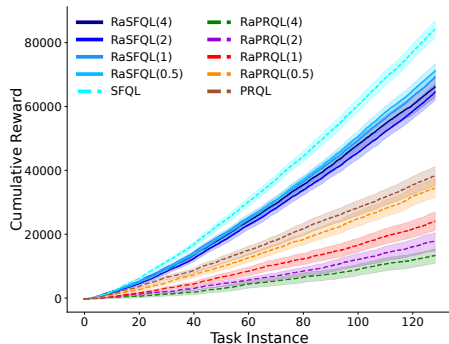
## Four-Room

---

Sensitivity to  $\beta$  parameter:

# Four-Room

Sensitivity to  $\beta$  parameter:

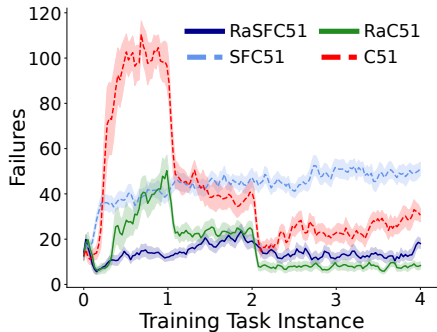
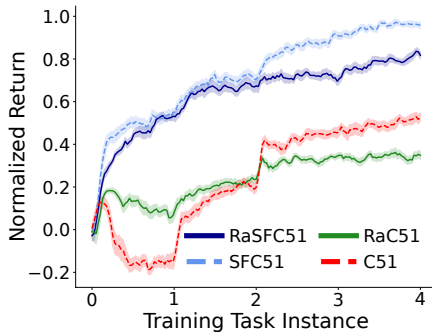






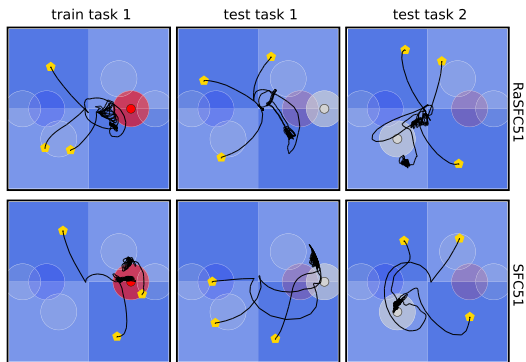
Train on four source tasks, test periodically on 8 unseen test tasks:

Train on four source tasks, test periodically on 8 unseen test tasks:



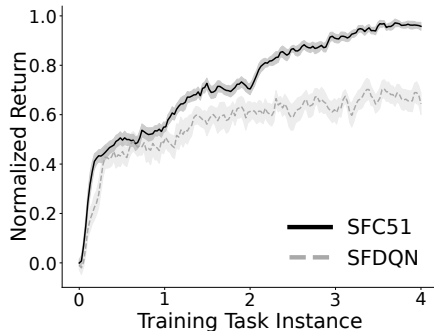
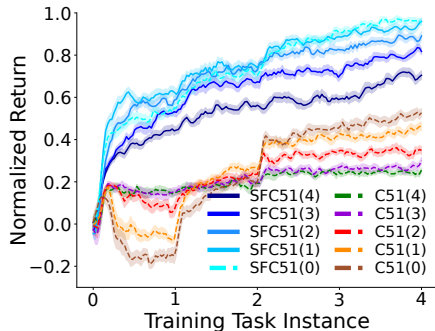
Does the agent learn risk-sensitive behavior?

Does the agent learn risk-sensitive behavior?



How sensitive is the agent to  $\beta$ ? Does the C51 method help in learning SFs?

How sensitive is the agent to  $\beta$ ? Does the C51 method help in learning SFs?



## Conclusion

---

## Conclusion

---

In conclusion:



## Conclusion

---

In conclusion:

- we presented Risk-aware Successor Features (RaSFs) for realizing policy transfer in domains where tasks have different goals

# Conclusion

---

In conclusion:

- we presented Risk-aware Successor Features (RaSFs) for realizing policy transfer in domains where tasks have different goals
- we extended generalized policy improvement to the risk-aware setting with entropic utilities

# Conclusion

---

In conclusion:

- we presented Risk-aware Successor Features (RaSFs) for realizing policy transfer in domains where tasks have different goals
- we extended generalized policy improvement to the risk-aware setting with entropic utilities
- we then extended the notion of generalized policy evaluation via the Taylor expansion of the entropic utility

# Conclusion

---

In conclusion:

- we presented Risk-aware Successor Features (RaSFs) for realizing policy transfer in domains where tasks have different goals
- we extended generalized policy improvement to the risk-aware setting with entropic utilities
- we then extended the notion of generalized policy evaluation via the Taylor expansion of the entropic utility
- together, risk-aware GPI and GPE are shown to inherit the superior task generalization abilities of successor features, while also learning to avoid risky situations

**Thank you.**