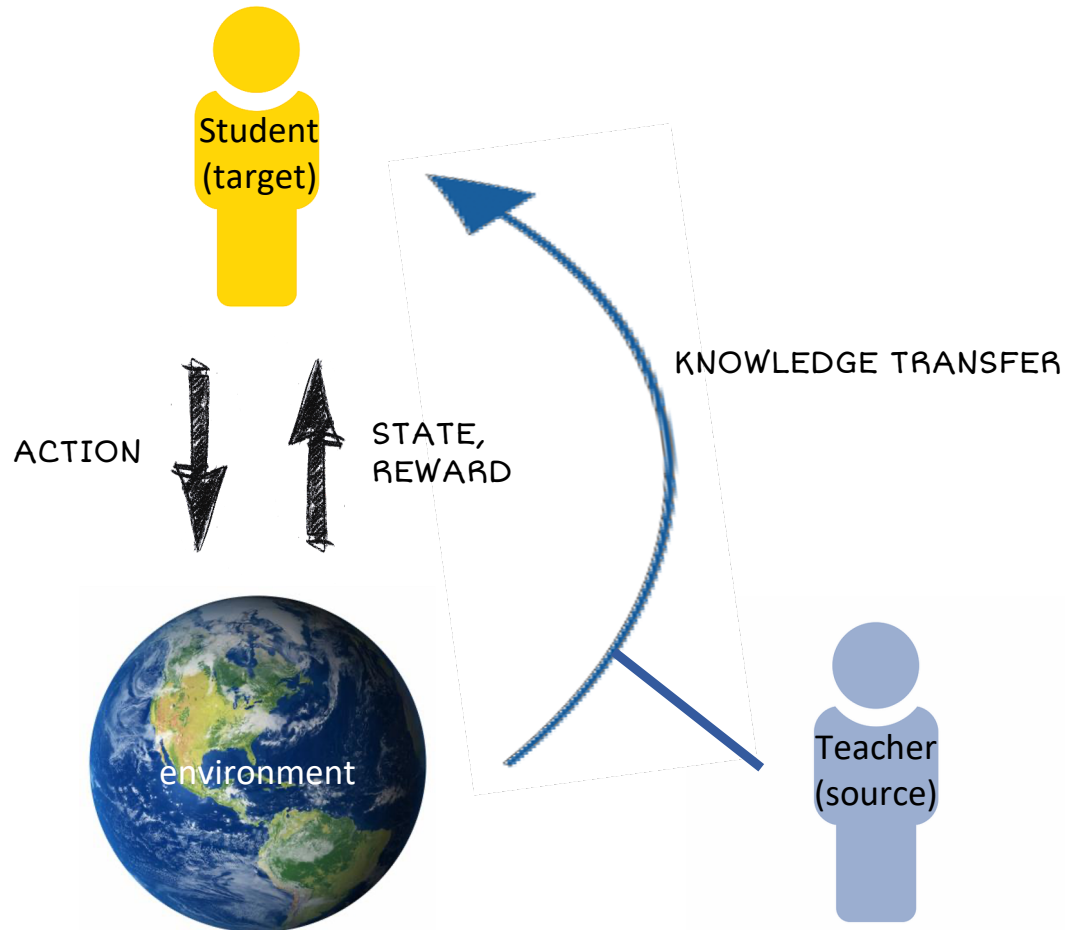# REPAINT: Knowledge Transfer in Deep Reinforcement Learning

**Yunzhe Tao**, Sahika Genc, Jonathan Chung, Tao Sun, Sunil Mallya

Amazon Web Services

# Transfer Learning in RL



Parameter transfer
- [Taylor '08; Mehta '08; Rajendran '15; Gupta '17; etc.]

Representation transfer
- [Konidaris '12; Parisotto '15; Schaul '15; Duan '16; Yin '17; Borsa '18; Zhang '18; Schmitt '18; Ma '18; Barreto '19; etc.]

Instance transfer
- [Lazaric '08; Taylor '08; Tirinzoni '18; etc.]

Our paper:

**REP**resentation **A**nd **IN**stance **T**ransfer (**REPAINT**)

- On-policy representation transfer
- Off-policy instance transfer
- Handles generic cases of source/target task similarity

# On-policy Representation Transfer

- Kickstarting Deep RL [Schmitt et al. 2018]
- Allow a student network (target policy) to exploit access to expert teachers:

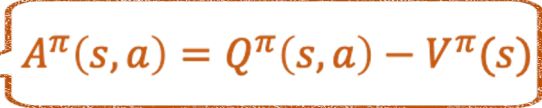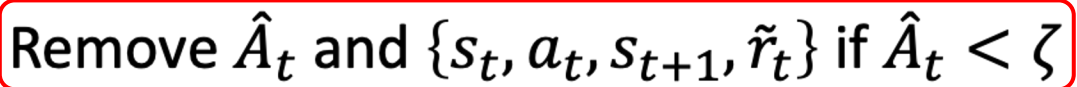$$L_{\mathrm{aux}}(\theta) = H(\pi_{\mathrm{teacher}}(a|s)||\pi_\theta(a|s))$$

cross-entropy

- The objective ($k$ is the iteration number):

$$L_{\mathrm{rep}}^k(\theta) = L_{\mathrm{RL}}(\theta) - \beta_k L_{\mathrm{aux}}(\theta)$$

RL objective          vanishing as $k$ increases

# Off-policy Instance Transfer

- Policy distillation works well only when source/target tasks are *similar*

- Idea: select "*good*" samples and update policy using those "*good*" samples

- "good" ≈ high **advantage** estimates $\longleftarrow$ $A^\pi(s,a) = Q^\pi(s,a) - V^\pi(s)$

- Advantage-based experience selection:
  1. Collect trajectories $\{s_i, a_i, s_{i+1}\}$ following **teacher policy** $\boldsymbol{\pi}_{\text{teacher}}$
  2. Compute rewards using **current reward function**: $\{s_i, a_i, s_{i+1}, \tilde{r}_i\}$
  3. Compute advantage estimates $\hat{A}_1, \hat{A}_2, \dots, \hat{A}_T$
  4. Remove $\hat{A}_t$ and $\{s_t, a_t, s_{t+1}, \tilde{r}_t\}$ if $\hat{A}_t < \zeta$

     Alternative: select top X% samples
  5. Update policy using selected samples
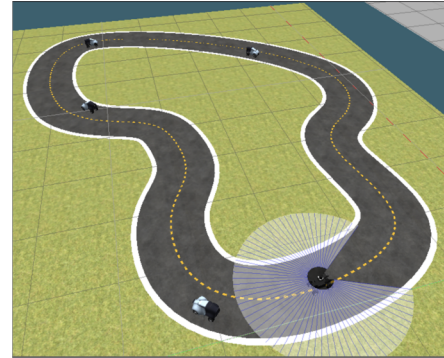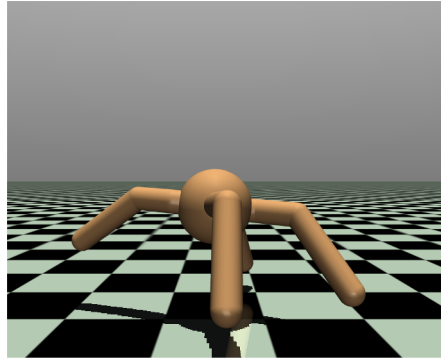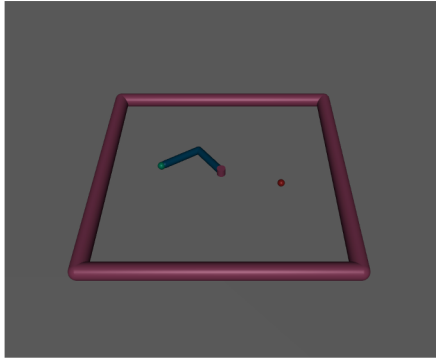
# REPAINT Algorithm with Actor-Critic RL

1. Collect $\mathcal{S} = \{s, a, s', r\}$ following $\pi_{\theta_{\mathrm{old}}}(\cdot)$

2. Collect $\tilde{\mathcal{S}} = \{\tilde{s}, \tilde{a}, \widetilde{s'}, \tilde{r}\}$ following $\pi_{\mathrm{teacher}}(\cdot)$

3. Update critic using $\mathcal{S}$

4. Perform advantage-based experience selection on $\tilde{\mathcal{S}}$

5. Update actor by:

$$\theta \leftarrow \theta + \alpha_1 \nabla_\theta L^k_{\mathrm{rep}}(\theta) + \alpha_2 \nabla_\theta L_{\mathrm{ins}}(\theta)$$

<span style="color:red">using $\mathcal{S}$</span>   <span style="color:red">using $\tilde{\mathcal{S}}$</span>   <span style="color:red">off-policy RL objective</span>

# Summary of Experimental Results



| Env. | Teacher type | Target score | $K_{\text{Baseline}}$ | $K_{\text{KS}}$ (pct. reduced) | $K_{\text{IT}}$ (pct. reduced) | $K_{\text{REPAINT}}$ (pct. reduced) | Best scores | | |
|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | KS | IT | REPAINT |
| Reacher | similar | -7.4 | 173 | 51 (71%) | 97 (44%) | 42 (76%) | -5.3 | -5.9 | -5.4 |
| | different | | | 73 (58%) | 127 (27%) | 51 (71%) | -6.9 | -6.4 | -5.2 |
| Ant | similar | 3685 | 997 | 363 (64%) | 623 (38%) | 334 (66%) | 5464 | 5172 | 5540 |
| Single-car | different | 394 | 18 | Not achieved | Not achieved | 13 (28%) | 331 | 388 | 396 |
| | different | 345 | 22 | Not achieved | Not achieved | 15 (32%) | 300 | 319 | 354 |
| Multi-car | sub-task | 1481 | 100 | 34 (66%) | 75 (25%) | 29 (71%) | 1542 | 1610 | 1623 |
| | diff/sub-task | 2.7 | 77 | 66 (14%) | 53 (31%) | 25 (68%) | 4.9 | 4.2 | 6.1 |
| StarCraft II | sub-task | 112 | 95 | 92 (3%) | 24 (75%) | 6 (94%) | 125 | 312 | 276 |

training time reduction