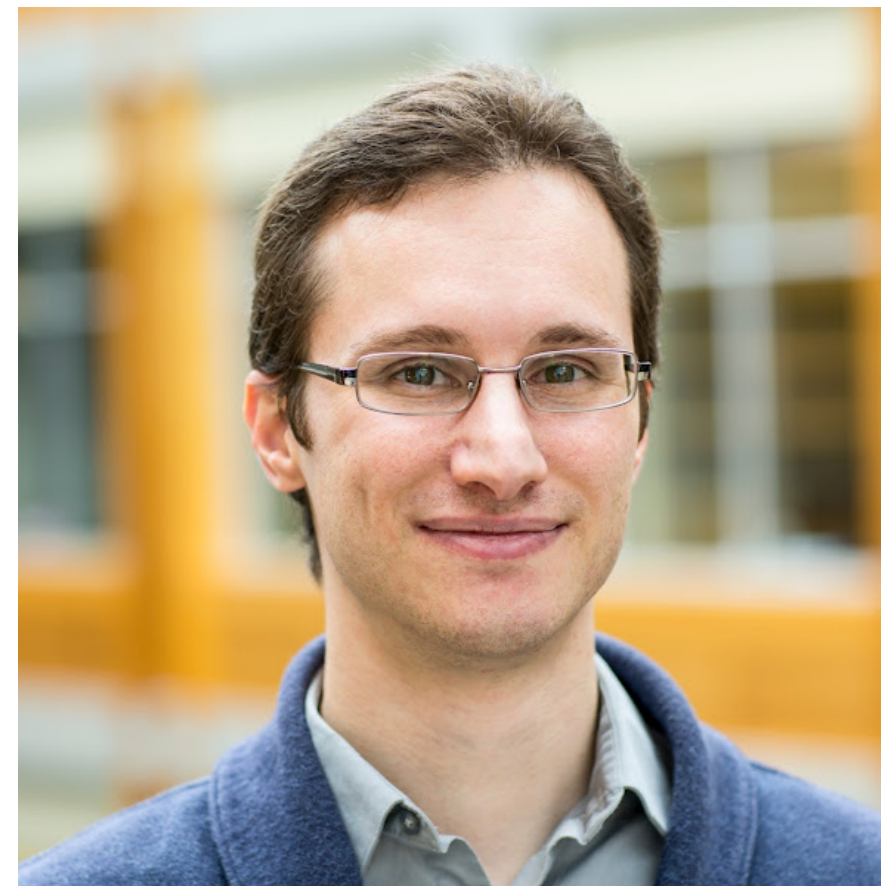


TRAIL: Near-Optimal Imitation Learning with Suboptimal Data

Sherry Yang

Sergey Levine

Ofir Nachum



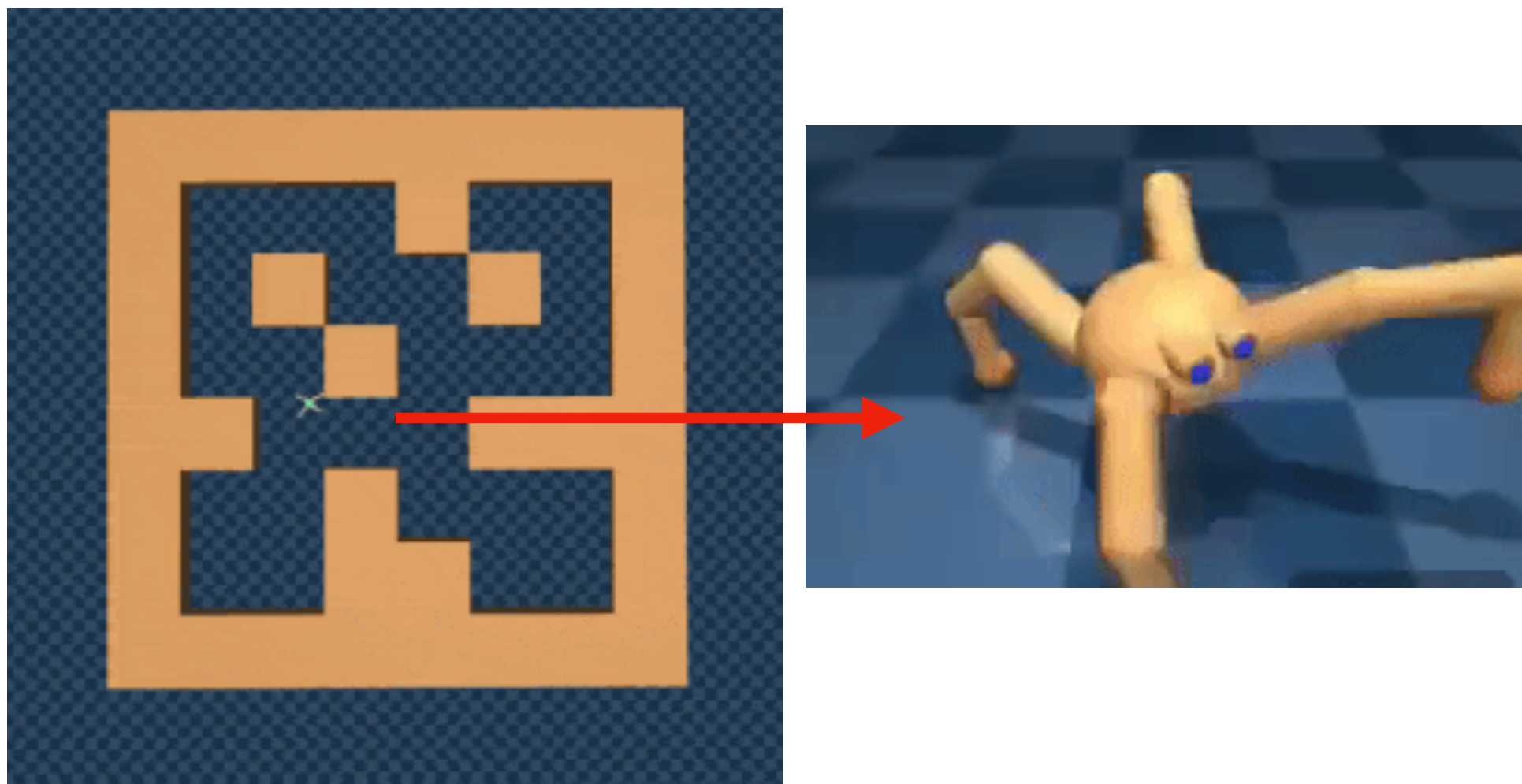
Paper: https://openreview.net/pdf?id=6q_2b6u0BnJ

Code: https://github.com/google-research/google-research/tree/master/rl_repr

Imitation Learning

Given expert demonstrations \mathcal{D}^{π^*}

Learn π that recovers π^* : $\text{Diff}(\pi, \pi_*) = D_{\text{TV}}(d^\pi \| d_*^\pi)$



Behavioral cloning:

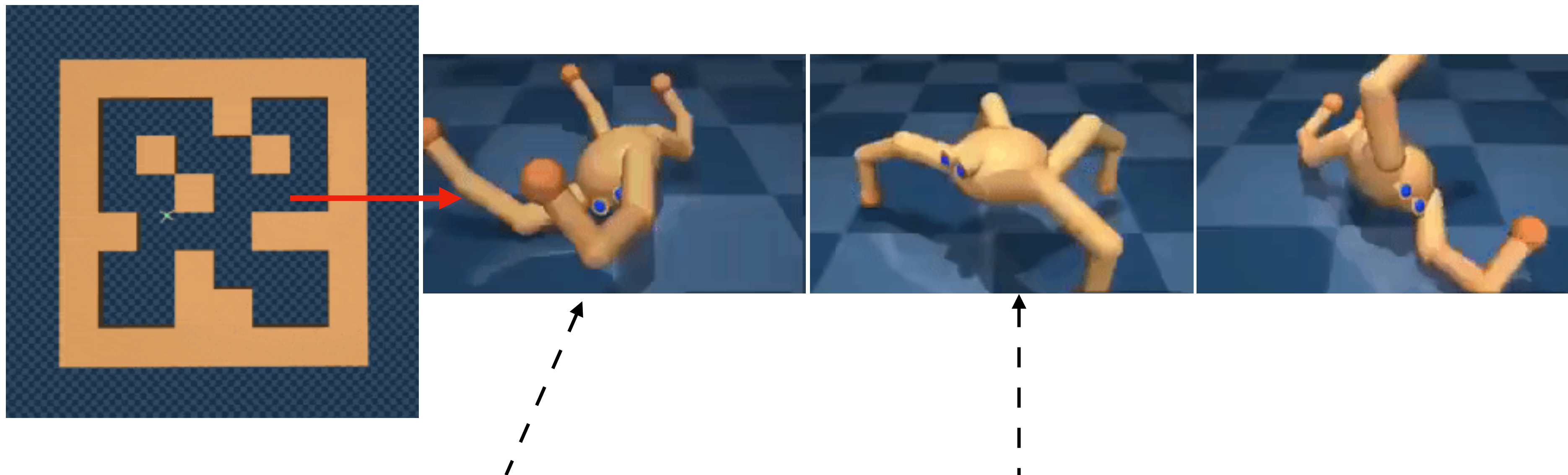
$$J_{\text{BC}}(\pi) := \mathbb{E}_{(s,a) \sim (d^{\pi_*}, \pi_*)} [-\log \pi(a|s)]$$

Limited & Hard to obtain
(e.g., involves human expert)

Suboptimal Offline Data

Large amounts of suboptimal offline data \mathcal{D}^{off}

How can \mathcal{D}^{off} provably facilitate imitation learning?

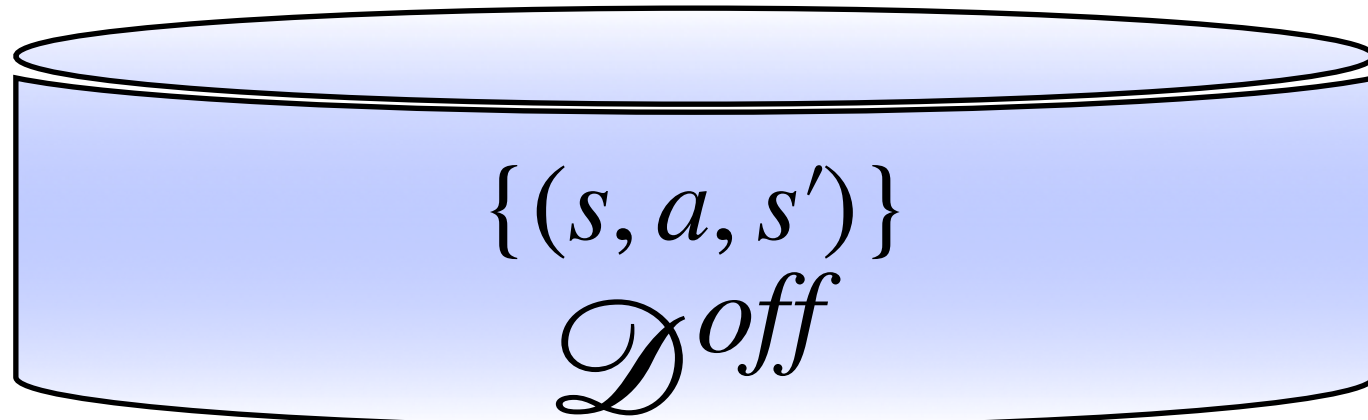


- Highly suboptimal (e.g., random policy)
- Single modal (e.g., collected by one stationary policy)

TRAIL: Transition Reparametrized Actions

Factored transition model

(1) $T_z \circ \phi(s, a)$



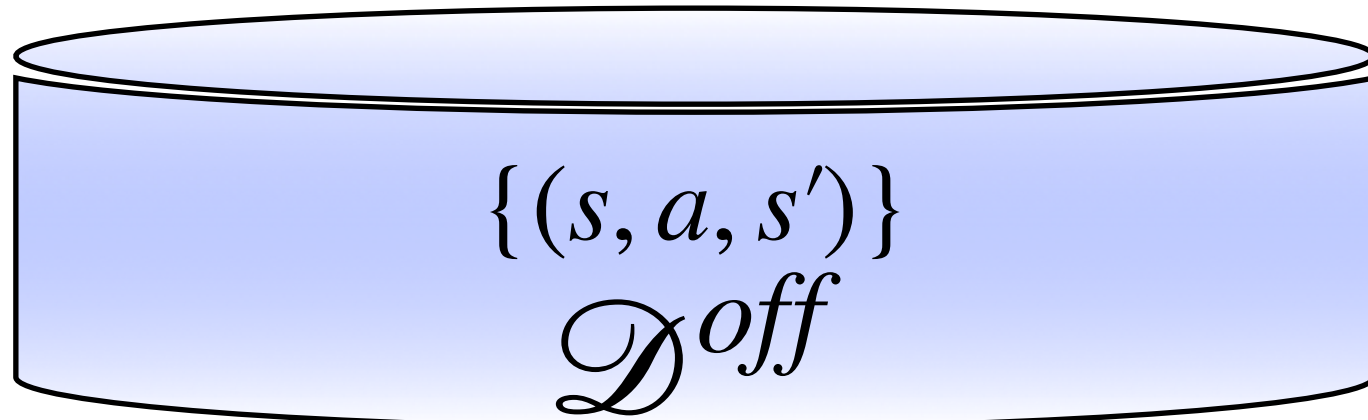
Pretraining

$$\left. \begin{array}{l} \text{Pretraining} \end{array} \right\} \underbrace{\mathbb{E}_{(s,a) \sim d^{\text{off}}} [D_{\text{KL}}(\mathcal{T}(s, a) \parallel \mathcal{T}_Z(s, \phi(s, a)))]}_{= J_{\text{T}}(\mathcal{T}_Z, \phi)} \quad (1)$$

TRAIL: Transition Reparametrized Actions

Action decoder

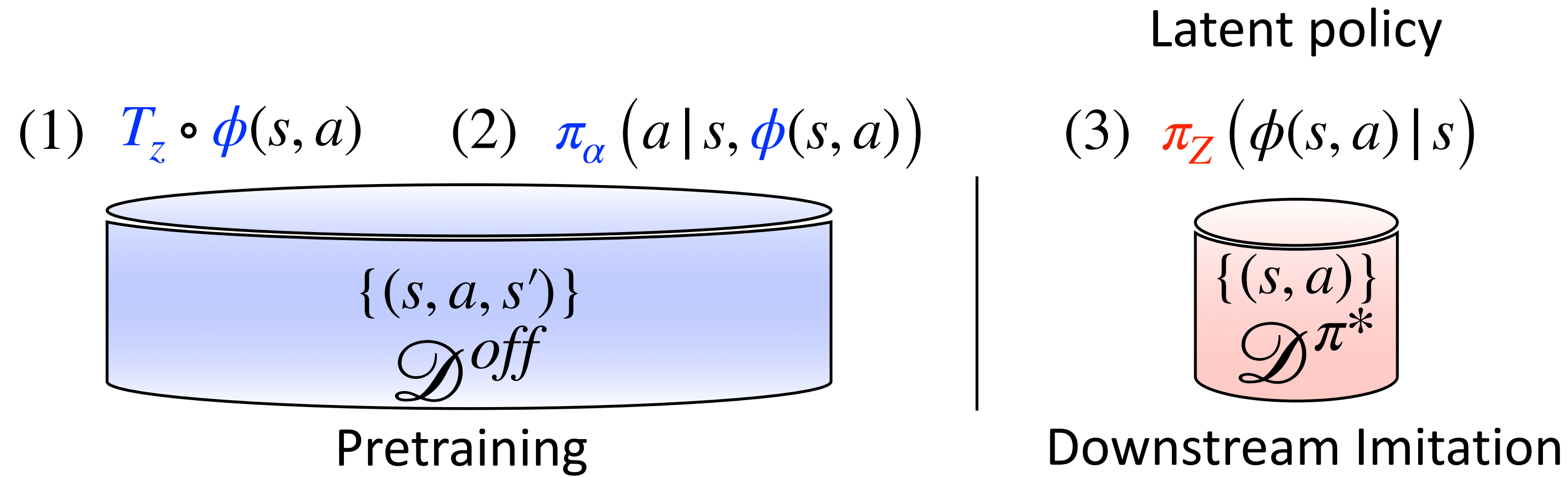
$$(1) \quad T_z \circ \phi(s, a) \quad (2) \quad \pi_\alpha(a | s, \phi(s, a))$$



Pretraining

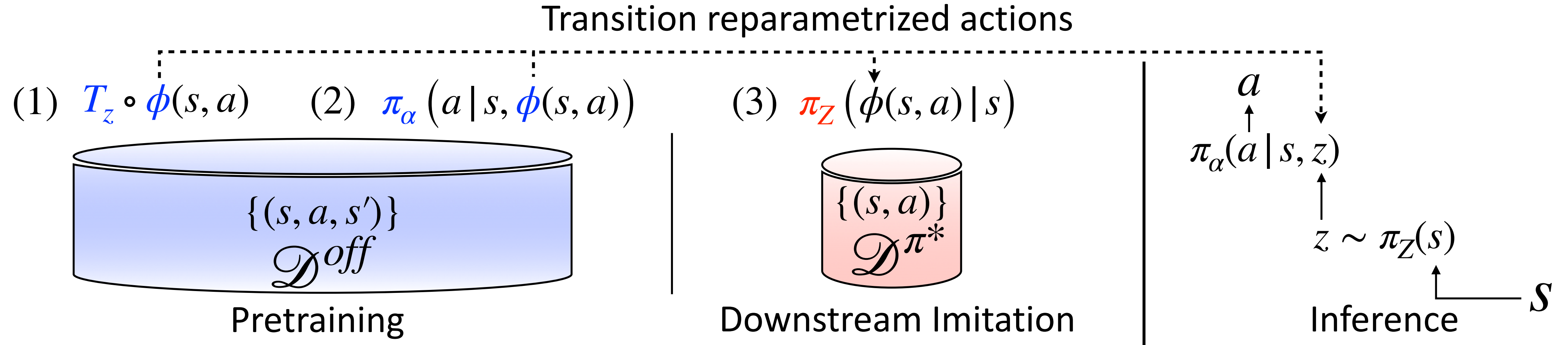
$$\left. \begin{array}{l} \text{Pretraining} \end{array} \right\} \begin{array}{l} \underbrace{\mathbb{E}_{(s,a) \sim d^{\text{off}}} [D_{\text{KL}}(\mathcal{T}(s, a) \| \mathcal{T}_Z(s, \phi(s, a)))]}_{= J_{\text{T}}(\mathcal{T}_Z, \phi)} \quad (1) \\ \underbrace{\mathbb{E}_{s \sim d^{\text{off}}} [\max_{z \in Z} D_{\text{KL}}(\pi_{\alpha^*}(s, z) \| \pi_\alpha(s, z))]}_{\approx \text{const}(d^{\text{off}}, \phi) + J_{\text{DE}}(\pi_\alpha, \phi)} \quad (2) \end{array}$$

TRAIL: Transition Reparametrized Actions



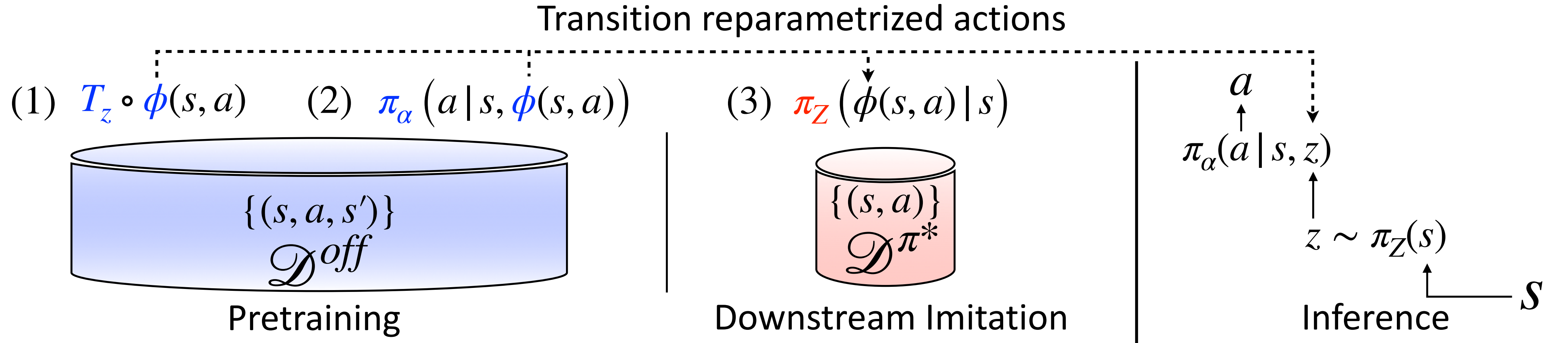
$$\begin{array}{l}
 \left. \begin{array}{l} \text{Pretraining} \end{array} \right\} \begin{array}{l}
 \mathbb{E}_{(s,a) \sim d^{\text{off}}} [D_{\text{KL}}(\mathcal{T}(s, a) \| \mathcal{T}_Z(s, \phi(s, a)))] \\
 \underbrace{\hspace{10em}} \\
 = J_{\text{T}}(\mathcal{T}_Z, \phi)
 \end{array} \quad (1) \\
 \\
 \left. \begin{array}{l} \text{Downstream} \\ \text{Imitation} \end{array} \right\} \begin{array}{l}
 \mathbb{E}_{s \sim d^{\text{off}}} [\max_{z \in Z} D_{\text{KL}}(\pi_{\alpha^*}(s, z) \| \pi_\alpha(s, z))] \\
 \underbrace{\hspace{10em}} \\
 \approx \text{const}(d^{\text{off}}, \phi) + J_{\text{DE}}(\pi_\alpha, \phi)
 \end{array} \quad (2) \\
 \\
 \left. \begin{array}{l} \text{Downstream} \\ \text{Imitation} \end{array} \right\} \begin{array}{l}
 \mathbb{E}_{s \sim d^{\pi^*}} [D_{\text{KL}}(\pi_{*,Z}(s) \| \pi_Z(s))] \\
 \underbrace{\hspace{10em}} \\
 = \text{const}(\pi_*, \phi) + J_{\text{BC},\phi}(\pi_Z)
 \end{array} \quad (3)
 \end{array}$$

TRAIL: Transition Reparametrized Actions



<i>Pretraining</i> {	$\underbrace{\mathbb{E}_{(s,a) \sim d^{\text{off}}} [D_{\text{KL}}(\mathcal{T}(s, a) \ \mathcal{T}_Z(s, \phi(s, a)))]}_{= J_{\text{T}}(\mathcal{T}_Z, \phi)} \quad (1)$
	$\underbrace{\mathbb{E}_{s \sim d^{\text{off}}} [\max_{z \in Z} D_{\text{KL}}(\pi_{\alpha^*}(s, z) \ \pi_\alpha(s, z))]}_{\approx \text{const}(d^{\text{off}}, \phi) + J_{\text{DE}}(\pi_\alpha, \phi)} \quad (2)$
<i>Downstream Imitation</i> {	$\underbrace{\mathbb{E}_{s \sim d^{\pi^*}} [D_{\text{KL}}(\pi_{*,Z}(s) \ \pi_Z(s))]}_{= \text{const}(\pi_*, \phi) + J_{\text{BC}, \phi}(\pi_Z)} \quad (3)$

TRAIL: Transition Reparametrized Actions



$$\text{Diff}(\pi_\alpha \circ \pi_Z, \pi_*) \leq$$

$$\left\{ \begin{array}{l} \text{Pretraining} \left\{ \begin{array}{l} C_1 \cdot \sqrt{\frac{1}{2} \mathbb{E}_{(s,a) \sim d^{\text{off}}} [D_{\text{KL}}(\mathcal{T}(s, a) \| \mathcal{T}_Z(s, \phi(s, a)))]} \\ \qquad \qquad \qquad = J_{\text{T}}(\mathcal{T}_Z, \phi) \end{array} \right. \quad (1) \\ + C_2 \cdot \sqrt{\frac{1}{2} \mathbb{E}_{s \sim d^{\text{off}}} [\max_{z \in Z} D_{\text{KL}}(\pi_{\alpha^*}(s, z) \| \pi_\alpha(s, z))]} \\ \qquad \qquad \qquad \approx \text{const}(d^{\text{off}}, \phi) + J_{\text{DE}}(\pi_\alpha, \phi) \end{array} \right. \quad (2)$$

$$\left\{ \begin{array}{l} \text{Downstream} \\ \text{Imitation} \end{array} \right\} + C_3 \cdot \sqrt{\frac{1}{2} \mathbb{E}_{s \sim d^{\pi^*}} [D_{\text{KL}}(\pi_{*,Z}(s) \| \pi_Z(s))]} \\ = \text{const}(\pi_*, \phi) + J_{\text{BC}, \phi}(\pi_Z)$$

$$\left\{ \begin{array}{l} C_1 = \gamma |A| (1 - \gamma)^{-1} (1 + D_{\chi^2}(d^{\pi^*} \| d^{\text{off}})^{\frac{1}{2}}) \\ C_2 = \gamma (1 - \gamma)^{-1} (1 + D_{\chi^2}(d^{\pi^*} \| d^{\text{off}})^{\frac{1}{2}}) \\ C_3 = \gamma (1 - \gamma)^{-1} \end{array} \right.$$

Sample Complexity of TRAIL

$$\mathbb{E}_{\mathcal{D}^{\pi_*}} [\text{Diff}(\pi_{opt,Z}, \pi_*)] \leq (1)(\phi_{opt}) + (2)(\phi_{opt}) + C_3 \cdot \sqrt{\frac{|Z||S|}{n}}.$$

So far, our analysis is based on tabular latent actions.

What about continuous latent actions and stochastic expert policy?

TRAIL with Linear Transition Dynamics

deterministic linear: $T_z = w(s')^\top \phi(s, a)$

$$\text{Diff}(\pi_\alpha \circ \pi_\theta, \pi_*) \leq (1)(\mathcal{T}_Z, \phi) + (2)(\pi_\alpha, \phi)$$

$$\text{Downstream Imitation} \left\{ + C_4 \cdot \left\| \frac{\partial}{\partial \theta} \mathbb{E}_{s \sim d^{\pi_*}, a \sim \pi_*(s)} [(\theta_s - \phi(s, a))^2] \right\|_1 \right.$$

TRAIL with Linear Transition Dynamics

deterministic linear: $T_z = w(s')^\top \phi(s, a)$

$$\text{Diff}(\pi_\alpha \circ \pi_\theta, \pi_*) \leq (1)(\mathcal{T}_Z, \phi) + (2)(\pi_\alpha, \phi)$$

$$\text{Downstream Imitation} \left\{ + C_4 \cdot \left\| \frac{\partial}{\partial \theta} \mathbb{E}_{s \sim d^{\pi_*}, a \sim \pi_*(s)} [(\theta_s - \phi(s, a))^2] \right\|_1 \right.$$

easier to optimize compared to:

$$\left[C_3 \cdot \sqrt{\frac{1}{2} \mathbb{E}_{s \sim d^{\pi_*}} [D_{\text{KL}}(\pi_{*,Z}(s) \parallel \pi_Z(s))]} \right]$$

$$= \text{const}(\pi_*, \phi) + J_{\text{BC}, \phi}(\pi_Z)$$

Learning TRAIL in Practice

$$(1) \quad T_z \circ \phi(s, a)$$

TRAIL EBM: $\mathcal{T}_Z(s'|s, \phi(s, a)) \propto \rho(s') \exp(-\|\phi(s, a) - \psi(s')\|^2)$.

$$\begin{aligned} \mathbb{E}_{d^{\text{off}}} [-\log \mathcal{T}_Z(s'|s, \phi(s, a))] &= \text{const}(d^{\text{off}}) + \frac{1}{2} \mathbb{E}_{d^{\text{off}}} [\|\phi(s, a) - \psi(s')\|^2] \quad \text{contrastive learning} \\ &\quad + \log \mathbb{E}_{\tilde{s}' \sim \rho} [\exp\{-\frac{1}{2} \|\phi(s, a) - \psi(\tilde{s}')\|^2\}] \end{aligned}$$

Learning TRAIL in Practice

$$(1) \quad T_z \circ \phi(s, a)$$

$$\text{TRAIL EBM: } \mathcal{T}_Z(s'|s, \phi(s, a)) \propto \rho(s') \exp(-\|\phi(s, a) - \psi(s')\|^2).$$

$$\mathbb{E}_{d^{\text{off}}} [-\log \mathcal{T}_Z(s'|s, \phi(s, a))] = \text{const}(d^{\text{off}}) + \frac{1}{2} \mathbb{E}_{d^{\text{off}}} [\|\phi(s, a) - \psi(s')\|^2] \quad \text{contrastive learning} \\ + \log \mathbb{E}_{\tilde{s}' \sim \rho} [\exp\{-\frac{1}{2} \|\phi(s, a) - \psi(\tilde{s}')\|^2\}]$$

$$\text{TRAIL linear: } \bar{\mathcal{T}}(s'|s, a) \propto \rho(s') \exp\{-\|f(s, a) - g(s')\|^2/2\} \propto \bar{\psi}(s')^\top \bar{\phi}(s, a)$$

$$\text{recover } \bar{\phi} \text{ with random Fourier features: } \bar{\phi}(s, a) = \cos(Wf(s, a) + b)$$

Learning TRAIL in Practice

$$(1) \quad T_z \circ \phi(s, a) \quad (2) \quad \pi_\alpha(a | s, \phi(s, a)) \quad (3) \quad \pi_Z(\phi(s, a) | s)$$

TRAIL EBM: $\mathcal{T}_Z(s' | s, \phi(s, a)) \propto \rho(s') \exp(-\|\phi(s, a) - \psi(s')\|^2)$.

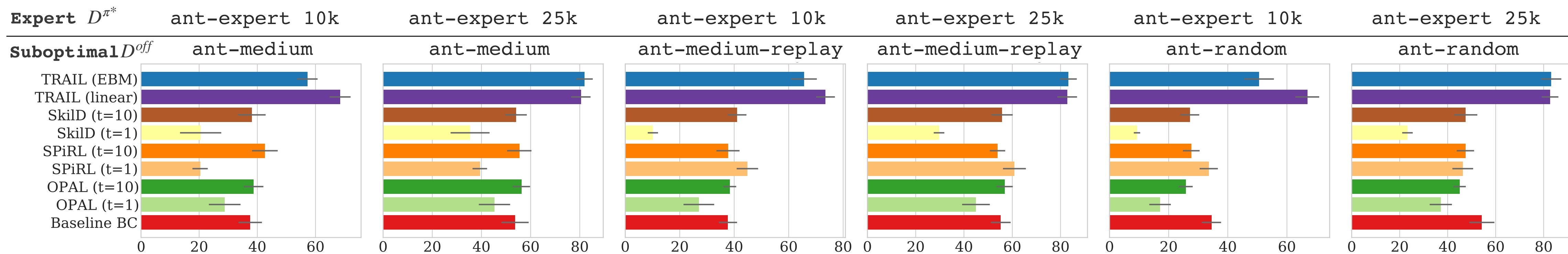
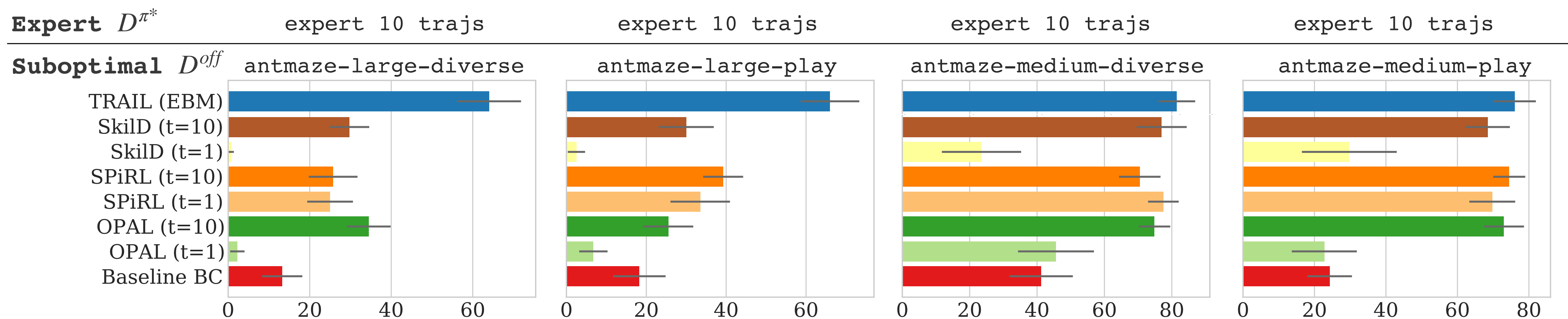
$$\mathbb{E}_{d^{\text{off}}}[-\log \mathcal{T}_Z(s' | s, \phi(s, a))] = \text{const}(d^{\text{off}}) + \frac{1}{2} \mathbb{E}_{d^{\text{off}}}[\|\phi(s, a) - \psi(s')\|^2] \quad \text{contrastive learning} \\ + \log \mathbb{E}_{\tilde{s}' \sim \rho}[\exp\{-\frac{1}{2}\|\phi(s, a) - \psi(\tilde{s}')\|^2\}]$$

TRAIL linear: $\bar{\mathcal{T}}(s' | s, a) \propto \rho(s') \exp\{-\|f(s, a) - g(s')\|^2/2\} \propto \bar{\psi}(s')^\top \bar{\phi}(s, a)$

recover $\bar{\phi}$ with random Fourier features: $\bar{\phi}(s, a) = \cos(Wf(s, a) + b)$

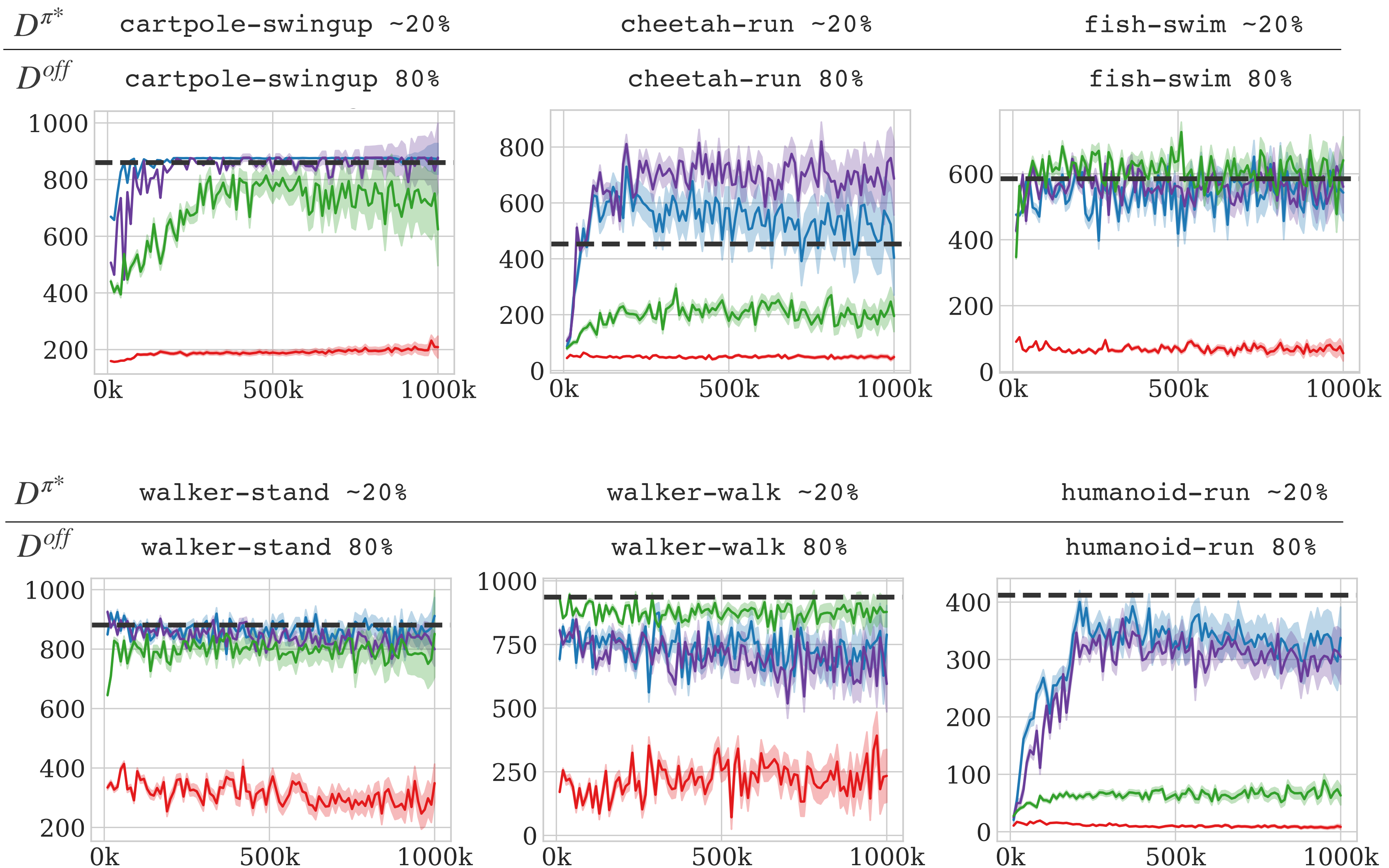
π_α and π_Z are neural-network parametrized Gaussian policies.

Experiments



Experiments - DM Control Suite

— TRAIL (energy) — TRAIL (linear) — Baseline BC — OPAL (t=10) - - - CRR



Recap & Conclusion

- How to utilize additional offline data for imitation learning?
 - Learn action representations.
- What if the offline data is highly suboptimal or unimodal?
 - Learn transition model as opposed to temporal skills.
- Representation learning + imitation learning as an alternative to offline RL?
 - Beneficial especially in the absence of reward labels.

Thank you. Checkout

Paper: https://openreview.net/pdf?id=6q_2b6u0BnJ

Code: https://github.com/google-research/google-research/tree/master/rl_repr