

Offline Multitask Representation Learning for Reinforcement Learning

Haque Ishfaq^{1,2}, Thanh Nguyen-Tang³, Songtao Feng⁴, Raman Arora³, Mengdi Wang⁵, Ming Yin⁵, Doina Precup^{1,2}

¹Mila – Quebec AI Institute ²McGill University ³Johns Hopkins University ⁴University of Florida ⁵Princeton University

Highlights

- We study offline multitask representation learning in reinforcement learning.
- Learner is provided with offline datasets from different tasks with shared representation.
- We prove our proposed algorithm can learn near-accurate model and near-optimal policies.
- We show theoretical benefits of using learned representation in downstream reward-free, offline and online RL tasks.

Setting

- We consider **low-rank** episodic MDPs $(\mathcal{S}, \mathcal{A}, H, \mathbb{P}, r)$: two unknown embedding functions $\phi_h^* : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}^d$ and $\mu_h^* : \mathcal{S} \rightarrow \mathbb{R}^d$ such that for all $s, s' \in \mathcal{S}$ and $a \in \mathcal{A}$, $P_h^*(s' | s, a) = \langle \phi_h^*(s, a), \mu_h^*(s') \rangle$.
- **Value function** of policy π :

$$V_{h,P,r}^\pi(s) = \mathbb{E}_{(s_{h'}, a_{h'}) \sim (P, \pi)} \left[\sum_{h'=h}^H r_{h'}(s_{h'}, a_{h'}) | s_h = s \right].$$

- Upstream offline multitask learning: T tasks, each task t : $\mathcal{M}^t = (\mathcal{S}, \mathcal{A}, H, P^t, r^t)$.

$$P_h^{(*,t)}(s' | s, a) = \langle \phi_h^*(s, a), \mu_h^{(*,t)}(s') \rangle, \quad \forall s, s' \in \mathcal{S}, a \in \mathcal{A}.$$

- We have access to offline dataset $\mathcal{D} = \bigcup_{t \in [T], h \in [H]} \mathcal{D}_h^{(t)}$, where $\mathcal{D}_h^{(t)} = \{(s_h^{(i,t)}, a_h^{(i,t)}, r_h^{(i,t)}, s_{h+1}^{(i,t)}) | i \in [n]\}$ with $s_{h+1}^{(i,t)} \sim P_h^{(*,t)}(\cdot | s_h^{(i,t)}, a_h^{(i,t)})$ and $\mathcal{D}_h^{(t)}$ was collected using a *fixed behavior policy* π_t^b .

- Downstream target task $T+1$ with

$$P_h^{(*,T+1)}(s' | s, a) = \langle \phi_h^*(s, a), \mu_h^{(*,T+1)}(s') \rangle, \quad \forall s, s' \in \mathcal{S}, a \in \mathcal{A}.$$

Goal of Upstream and Downstream Tasks

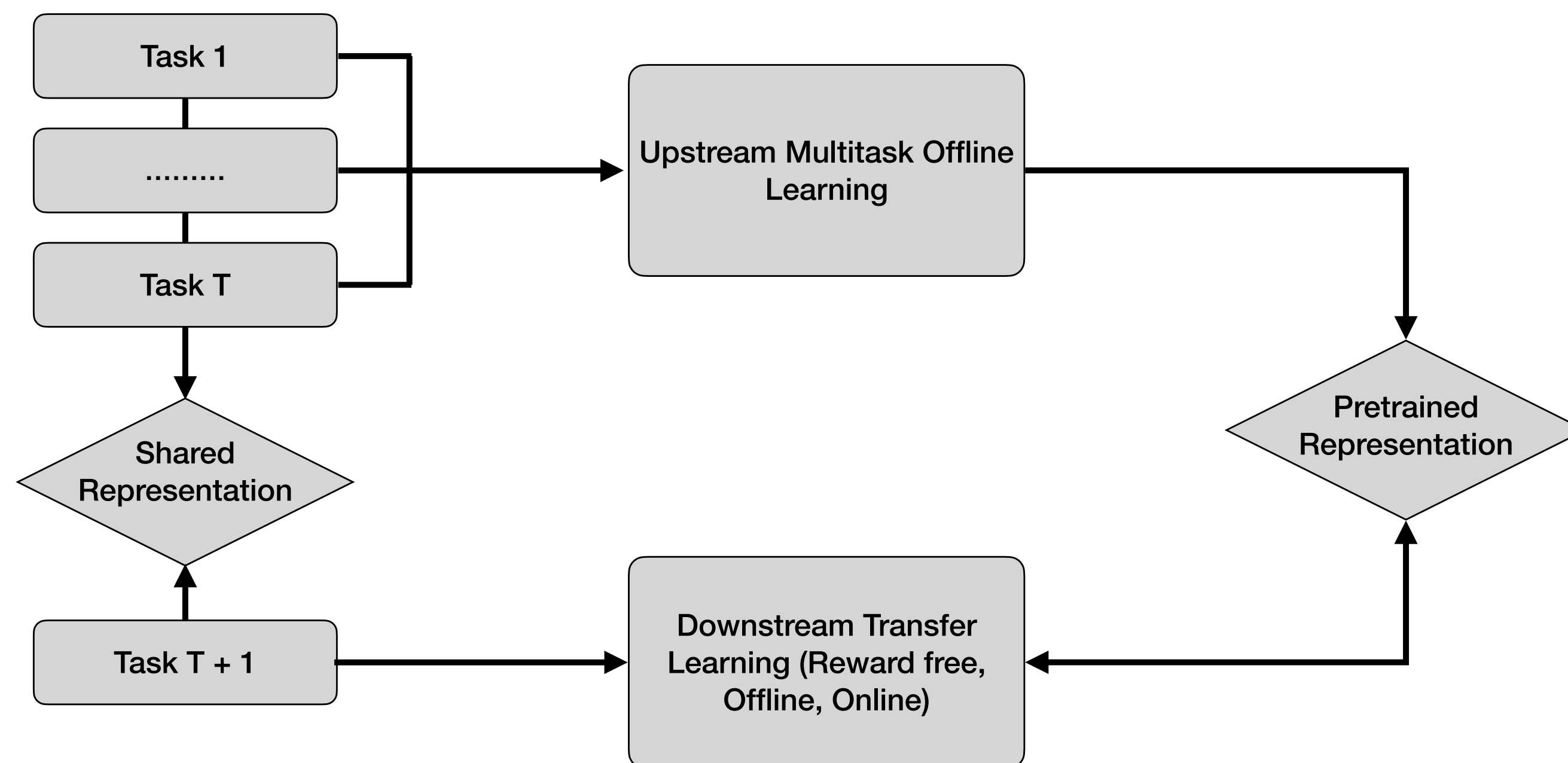


Figure 1. Upstream and Downstream Task Overview

Algorithm: Multitask Offline Representation Learning (MORL)

- Learning jointly via offline Multitask Maximum Likelihood Estimation (MLE) oracle

$$(\hat{\phi}_h, \hat{\mu}_h^{(1)}, \dots, \hat{\mu}_h^{(T)}) = \operatorname{argmax}_{\phi_h \in \Phi, \mu_h^{(1)}, \dots, \mu_h^{(T)} \in \Psi} \sum_{i=1}^n \sum_{t=1}^T \log \left(\langle \phi_h(s_h^{(i,t)}, a_h^{(i,t)}), \mu_h^{(i,t)}(s_{h+1}^{(i,t)}) \rangle \right).$$

- for each $t \in \{1, \dots, T\}$ define:

- ▶ estimated transition kernel: $\hat{P}_h^{(t)}(s' | s, a) = \langle \hat{\phi}_h(s, a), \hat{\mu}_h^{(t)}(s') \rangle$.
- ▶ empirical covariance matrix: $\hat{\Sigma}_{h, \hat{\phi}}^{(t)} = \sum_{i=1}^n \hat{\phi}_h(s_h^{(i,t)}, a_h^{(i,t)}) \hat{\phi}_h(s_h^{(i,t)}, a_h^{(i,t)})^\top + \lambda I$.

- ▶ penalty term: $\hat{b}_h^{(t)}(s_h, a_h) = \min \left\{ \alpha \|\hat{\phi}_h(s_h, a_h)\|_{(\hat{\Sigma}_{h, \hat{\phi}}^{(t)})^{-1}}, 1 \right\}$.

- ▶ Get policy $\hat{\pi}_t = \operatorname{argmax}_\pi V_{\hat{P}^{(t)}, r^t} \hat{b}^{(t)}$

- **Output:** $\hat{\phi}, \hat{P}^{(1)}, \dots, \hat{P}^{(T)}, \hat{\pi}_1, \dots, \hat{\pi}_T$

Theoretical Result on Upstream Task

Definition 1 (Multi-task relative condition number). For task t and time step h , we define $C_{t,h}^*(\pi_t, \pi_t^b)$ as the relative condition number under ϕ_h^* :

$$C_{t,h}^*(\pi_t, \pi_t^b) := \sup_{x \in \mathbb{R}^d} \frac{x^\top \mathbb{E}_{(s_h, a_h) \sim (P^{(*,t)}, \pi_t)} [\phi_h^*(s_h, a_h) \phi_h^*(s_h, a_h)^\top] x}{x^\top \mathbb{E}_{(s_h, a_h) \sim (P^{(*,t)}, \pi_t^b)} [\phi_h^*(s_h, a_h) \phi_h^*(s_h, a_h)^\top] x}.$$

We define $C_t^* := \max_{h \in [H]} C_{t,h}^*(\pi_t, \pi_t^b)$ and $C^* := \max_{t \in [T]} C_t^*$.

Theorem 1. Under realizability assumption, with probability at least $1 - \delta$, for any step $h \in [H]$, we have

$$\frac{1}{T} \sum_{t=1}^T \mathbb{E}_{(s_h, a_h) \sim (P^{(*,t)}, \pi_t^b)} \left[\left\| \hat{P}_h^{(t)}(\cdot | s_h, a_h) - P_h^{(*,t)}(\cdot | s_h, a_h) \right\|_{TV} \right] \leq \sqrt{\frac{2 \log(2|\Phi| |\Psi|^T n H / \delta)}{nT}},$$

where $\hat{\phi}, \hat{P}^{(1)}, \dots, \hat{P}^{(T)}$ be the output of MORL.

In addition, if we set $\alpha = \sqrt{2n\omega\zeta_n + \lambda d}$, $\lambda = cd \log(|\Phi| |\Psi|^T n H / \delta)$ with $\zeta_n := \frac{2 \log(2|\Phi| |\Psi|^T n H / \delta)}{n}$ and c being a constant, where we assume that $\omega := \max_t \max_{s,a} (1/\pi_t^b(a | s)) < \infty$, then under realizability assumption, with probability at least $1 - \delta$, we have

$$\frac{1}{T} \sum_{t=1}^T \left[V_{P^{(*,t)}, r^t}^{\hat{\pi}_t} - V_{P^{(*,t)}, r^t}^{\pi_t} \right] \leq \omega \alpha d H \sqrt{\frac{C^*}{n}} + 2dH^2 \sqrt{\frac{\lambda C^*}{n}} + \omega H^2 \sqrt{\frac{dC^* \zeta_n}{T}} + \alpha \sqrt{\frac{d}{n}} + 2H \sqrt{\frac{\omega \zeta_n}{T}},$$

where $\{\hat{\pi}_t\}_{t \in [T]}$ is the output of the algorithm MORL.

Connecting Upstream and Downstream tasks

- **Assumptions:** reachability of behavior policies, compact state space, smoothness of transition probabilities, and approximate linear combination.

Lemma 2. Under the above assumptions, the output $\hat{\phi}$ of MORL is a ξ_{down} -approximate feature for MDP \mathcal{M}^{T+1} where $\xi_{\text{down}} = \xi + \frac{C_L C_{RV}}{\kappa} \sqrt{\frac{2T \log(2|\Phi| |\Psi|^T n H / \delta)}{n}}$, i.e. there exist a time-dependent unknown (signed) measure $\hat{\mu}^*$ over \mathcal{S} such that for any $(s, a) \in \mathcal{S} \times \mathcal{A}$, we have

$$\|P_h^{(*,T+1)}(\cdot | s, a) - \langle \hat{\phi}_h(s, a), \hat{\mu}_h^*(\cdot) \rangle\|_{TV} \leq \xi_{\text{down}}.$$

Downstream RL: Reward-free Exploration

Theorem 3. Under the above assumptions, after collecting K_{RFE} trajectories during the exploration phase, with probability at least $1 - \delta$, the output of the planning phase, policy π satisfies

$$\mathbb{E}_{s_1 \sim \mu} [V_1^*(s_1, r) - V_1^\pi(s_1, r)] \leq c' \sqrt{d^3 H^4 \log(dK_{RFE} H / \delta) / K_{RFE}} + 6H^2 \xi_{\text{down}}. \quad (1)$$

If the linear combination misspecification error ξ satisfies $\tilde{O}(\sqrt{d^3 / K_{RFE}})$ and the number of trajectories in the offline dataset for each upstream task is at least $\tilde{O}(TK_{RFE}/d^3)$, then, provided K_{RFE} is at least $O(H^4 d^3 \log(dH\delta^{-1}\epsilon^{-1})/\epsilon^2)$, with probability $1 - \delta$, the policy π will be an ϵ -optimal policy for any given reward during the planning phase.

Algorithm	Sample Complexity	Task
FLAMBE [Agarwal et al., 2020]	$\tilde{O}(\frac{H^{22} d^7 K^9}{\epsilon^{10}})$	Single task
MOFFLE [Modi et al., 2021]	$\tilde{O}(\frac{H^7 d^{11} K^{14}}{\min\{\epsilon^2, \eta^3\}})$	Single task
RAFFLE [Cheng et al., 2023]	$\tilde{O}(\frac{H^9 d^4 K}{\epsilon^2})$	Single task
This work	$\tilde{O}(\frac{H^4 d^3}{\epsilon^2})$	Multi-task

Table 1. Sample complexities of different approaches to learning an ϵ -optimal policy for the reward-free RL setting with low-rank MDPs.

Downstream RL: Offline RL and Online RL

Downstream Offline Task:

Assumption 4 (Feature coverage). There exists an absolute constant κ_ρ such that for all $h \in [H]$ and $\phi_h \in \Phi_h$, $\lambda_{\min}(\mathbb{E}_\rho[\phi_h(s_h, a_h) \phi_h(s_h, a_h)^\top | s_1 = s]) \geq \kappa_\rho$.

Theorem 5 (Downstream offline task). Under the above assumptions and the sample size $N_{\text{off}} \geq 40/\kappa_\rho \cdot \log(4dH/\delta)$, with probability at least $1 - \delta$, the suboptimality gap of offline downstream task is at most

$$V_{P^{(*,T+1)}, r}^{\pi^*}(s) - V_{P^{(*,T+1)}, r}^{\hat{\pi}}(s) \leq O\left(\kappa_\rho^{-1/2} H^2 d \sqrt{\frac{\log(HdN_{\text{off}} \max(\xi_{\text{down}}, 1)/\delta)}{N_{\text{off}}}} + \kappa_\rho^{-1/2} H^2 d^{1/2} \xi_{\text{down}}\right).$$

Downstream Online Task:

Theorem 6 (Downstream online task). Let $\tilde{\pi}$ be the uniform mixture of $\pi^1, \dots, \pi^{N_{\text{on}}}$. Under the above assumptions, with probability $1 - \delta$, the suboptimality gap of online downstream task satisfies

$$V_{P^{(*,T+1)}, r}^{\pi^*} - V_{P^{(*,T+1)}, r}^{\tilde{\pi}} \leq \tilde{O}(H^2 d^{3/2} N_{\text{on}}^{-1/2} + H^2 d \xi_{\text{down}}).$$

For more details check the paper!

