

Behavior Alignment via Reward Function Optimization (Supplemental Material)

Table 1: Notations

Symbol	Description
θ	Parameters for policy π
ϕ	Parameters for reward function
φ	Parameters for learned γ
$\pi_\theta, r_\phi, \gamma_\varphi$	Functional form of policy, reward and γ with their respective parameters
$\alpha_\theta, \alpha_\phi, \alpha_\varphi$	Step sizes for the respective parameters
$\lambda_\theta, \lambda_\phi, \lambda_\varphi$	Regularization for policy, reward and γ function
δ	Number of on-policy samples collected between subsequent updates to ϕ, φ
η	Neumann Approximator Eigen value scaling factor
n	Number of loops used in Neumann Approximation
optim	Any standard optimizer like Adam, RMSprop, SGD, which takes input as gradients and outputs the appropriate update
E	Total Number of episodes to sample from the environment
N_i	Number of updates to be performed for updating the π by Alg
N_0	Number of initial updates to be performed

551 A Proofs for Theoretical Results

552 In this section, we provide proofs for Property 1, Property 2, and Property 3. For the purpose of these proofs, we
 553 introduce some additional notation. To have a unified MDP notation for goal-based and time-based tasks, we
 554 first consider that in the the time-based task, time is a part of the state such that Markovian dynamics is ensured.

555 The (un-normalized) discounted and undiscounted visitation probability is denoted as,

$$d_\gamma^\pi(s, a) := \sum_{t=0}^T \gamma^t \Pr(S_t = s, A_t = a; \pi), \quad (9)$$

$$\bar{d}^\pi(s, a) := \sum_{t=0}^T \Pr(S_t = s, A_t = a; \pi). \quad (10)$$

556 We can normalize so that it is a distribution as follows :

$$d^\pi(s, a) := \frac{\bar{d}^\pi(s, a)}{\sum_{s' \in S, a' \in \mathcal{A}} \bar{d}^\pi(s', a')}$$

557 **Property 1.** $\Delta(\theta, \tilde{r}) = \Delta(\theta, r_p)$.

Proof.

$$\begin{aligned}
\Delta(\theta, \tilde{r}) &= \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \psi_\theta(S_t, A_t) \sum_{j=t}^T \gamma^{j-t} \tilde{r}(S_j, A_j) \right] \\
&= \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \psi_\theta(S_t, A_t) \left(\sum_{j=t}^T \gamma^{j-t} (r_p(S_j, A_j) + \gamma \Phi(S_{j+1}) - \Phi(S_j)) \right) \right] \\
&= \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \psi_\theta(S_t, A_t) \sum_{j=t}^T \gamma^{j-t} r_p(S_j, A_j) \right] + \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \psi_\theta(S_t, A_t) \sum_{j=t}^T \gamma^{j-t} (\gamma \Phi(S_{j+1}) - \Phi(S_j)) \right] \\
&= \Delta(\theta, r_p) + \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \psi_\theta(S_t, A_t) \sum_{j=t}^T \gamma^{j-t} (\gamma \Phi(S_{j+1}) - \Phi(S_j)) \right] \\
&\stackrel{(a)}{=} \Delta(\theta, r_p) + \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \psi_\theta(S_t, A_t) (\gamma^{T-t+1} \Phi(S_{T+1}) - \Phi(S_t)) \right] \\
&\stackrel{(b)}{=} \Delta(\theta, r_p) + \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \psi_\theta(S_t, A_t) (\gamma^{T-t+1} c - \Phi(S_t)) \right] \\
&\stackrel{(c)}{=} \Delta(\theta, r_p) + \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T (\gamma^{T-t+1} c - \Phi(S_t)) \mathbb{E}_{\pi_\theta} [\psi_\theta(S_t, A_t) | S_t] \right] \\
&\stackrel{(d)}{=} \Delta(\theta, r_p),
\end{aligned}$$

558 where (a) holds because on the expansion of future return, intermediate potential values cancel out, (b) holds
559 because S_{T+1} is the terminal state and potential function is defined to be a fixed constant c for the terminal state
560 [42], (c) holds from the law of total expectation, and (d) holds because,

$$\mathbb{E}_{\pi_\theta} [\psi_\theta(S_t, A_t) | S_t] = \sum_{a \in \mathcal{A}} \pi_\theta(S_t, a) \frac{\partial \ln \pi_\theta(S_t, a)}{\partial \theta} = \sum_{a \in \mathcal{A}} \frac{\partial \pi_\theta(S_t, a)}{\partial \theta} = \frac{\partial}{\partial \theta} \sum_{a \in \mathcal{A}} \pi_\theta(S_t, a) = 0.$$

561 In the stochastic setting, i.e., when using sample average estimates instead of the true expectation, $\gamma^{T-t+1} c -$
562 $\Phi(S_t)$ is analogous to a state dependent baseline for the sum of discounted future primary rewards. It may reduce
563 or increase the variance of $\Delta(\theta, r_p)$, depending on this baseline's co-variance with $\sum_{j=t}^T \gamma^{j-t} r_p(S_j, A_j)$. \square

564 **Note:** As we encountered the potential at the terminal state to be c , as it is a constant, we will use the value of
565 $c = 0$ in accordance with [42].

566 $\mathbb{V}_{\pi_\theta} [\hat{\Delta}(\theta, \tilde{r})]$ can be more than $\mathbb{V}_{\pi_\theta} [\hat{\Delta}(\theta, r_p)]$.

567 *Proof.* We will first look at a different variant of the update, i.e., $\Delta(\theta, \tilde{r}) =$
568 $\mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \psi_\theta(S_t, A_t) \sum_{j=t}^T \gamma^{j-t} \tilde{r}(S_j, A_j) \right]$ can also be written as follows:

$$\Delta(\theta, \tilde{r}) = \mathbb{E}_{d^\pi, \pi} [\psi_\theta(S, A) (q^\pi(S, A) - \Phi(S))].$$

569 This is proved as part of the proof of Property 2. Similarly the baseline update (using r_p can be written as,

$$\Delta(\theta, r_p) = \mathbb{E}_{d^\pi, \pi} [\psi_\theta(S, A) (q^\pi(S, A))].$$

570 Let's put $Y = \psi_\theta(S, A) (q^\pi(S, A) - \Phi(S))$, $X = \psi_\theta(S, A) q^\pi(S, A)$ and $Z = \psi_\theta(S, A) q^\pi(S)$ as our random
571 variables and try to understand when the variance of $\hat{\Delta}(\theta, \tilde{r})$ can be higher than $\hat{\Delta}(\theta, r_p)$. Also note that
572 $\mathbb{E}[Y] = \mathbb{E}[X]$

$$\begin{aligned}
\mathbb{V}_{d^\pi}[Y] &= \mathbb{E}_{d^\pi}[Y^T Y] - \mathbb{E}[Y]^T \mathbb{E}[Y] \\
&= \mathbb{E}_{d^\pi} \left[(\psi_\theta(S, A)(q^\pi(S, A) - \Phi(S)))^T (\psi_\theta(S, A)(q^\pi(S, A) - \Phi(S))) \right] - \mathbb{E}[X]^T \mathbb{E}[X] \\
&= \mathbb{E}_{d^\pi} \left[\psi_\theta(S, A)^T \psi_\theta(S, A) ((q^\pi(S, A) - \Phi(S)))^2 \right] - \mathbb{E}[X]^T \mathbb{E}[X] \\
&= \mathbb{E}_{d^\pi} \left[\psi_\theta(S, A)^T \psi_\theta(S, A) \Phi(S)^2 \right] - 2\mathbb{E}_{d^\pi} \left[\psi_\theta(S, A)^T \psi_\theta(S, A) (q^\pi(S, A) \Phi(S)) \right] + \\
&\quad \mathbb{E}_{d^\pi} \left[\psi_\theta(S, A)^T \psi_\theta(S, A) q^\pi(S, A)^2 \right] - \mathbb{E}[X]^T \mathbb{E}[X]
\end{aligned}$$

573 The last term in above is the variance of the method with primary reward i.e.
574 $\mathbb{E}_{d^\pi}[\psi_\theta(S, A)^T \psi_\theta(S, A) q^\pi(S, A)^2] - \mathbb{E}[X]^T \mathbb{E}[X] = \mathbb{V}_{d^\pi}[X]$

$$\mathbb{V}_{d^\pi}[Y] - \mathbb{V}_{d^\pi}[X] = \mathbb{E}_{d^\pi}[\psi_\theta(S, A)^T \psi_\theta(S, A) \Phi(S)^2] - 2\mathbb{E}_{d^\pi}[\psi_\theta(S, A)^T \psi_\theta(S, A) (q^\pi(S, A) \Phi(S))]$$

575 Hence variance of Y will be higher than X if $\mathbb{E}_{d^\pi}[\psi_\theta(S, A)^T \psi_\theta(S, A) \Phi(S)^2] -$
576 $2\mathbb{E}_{d^\pi}[\psi_\theta(S, A)^T \psi_\theta(S, A) (q^\pi(S, A) \Phi(S))] > 0$

577 Now lets look at variance of $\hat{\Delta}(\theta, \tilde{r})$ and $\hat{\Delta}(\theta, r_p)$, lets use $X_i = \psi_\theta(S_i, A_i) q^\pi(S_i, A_i)$, similarly for Y_i and
578 Z_i . Hence we can write $\hat{\Delta}(\theta, \tilde{r}) = \sum_{i=0}^T Y_i$ and $\hat{\Delta}(\theta, r_p) = \sum_{i=0}^T X_i$. Hence, as we saw that it is possible
579 that $\mathbb{V}[Y] > \mathbb{V}[X]$, let's look at the variance of $\mathbb{V}[\sum_{i=0}^T Y_i] > \mathbb{V}_{d^\pi}[\sum_{i=0}^T X_i]$, for simplicity lets assume only
580 2-time steps, or states at time i, j .

$$\mathbb{V}[Y_i + Y_j] = \mathbb{V}[Y_i] + \mathbb{V}[Y_j] + 2\text{Cov}(Y_i, Y_j)$$

581 and

$$\mathbb{V}[X_i + X_j] = \mathbb{V}[X_i] + \mathbb{V}[X_j] + 2\text{Cov}(X_i, X_j)$$

582 As we can see that $\mathbb{V}[Y] > \mathbb{V}[X]$, hence we need to see the relationship between $\text{Cov}(Y_i, Y_j)$ and $\text{Cov}(X_i, X_j)$
583 to understand if $\mathbb{V}[Y_i + Y_j] > \mathbb{V}[X_i + X_j]$.

584 We can write $Y_i = X_i - Z_i$, hence

$$\begin{aligned}
\text{Cov}(Y_i, Y_j) &= \mathbb{E}[Y_i Y_j] - \mathbb{E}[Y_i] \mathbb{E}[Y_j] \\
&= \mathbb{E}[(X_i - Z_i)(X_j - Z_j)] - \mathbb{E}[X_i] \mathbb{E}[X_j] \\
&= \mathbb{E}[X_i X_j - X_i Z_j - X_j Z_i + Z_i Z_j] - \mathbb{E}[X_i] \mathbb{E}[X_j] \\
&= \mathbb{E}[X_i X_j] - \mathbb{E}[X_i Z_j] - \mathbb{E}[X_j Z_i] + \mathbb{E}[Z_i Z_j] - \mathbb{E}[X_i] \mathbb{E}[X_j] \\
&= \text{Cov}(X_i, X_j) - \mathbb{E}[X_i Z_j] - \mathbb{E}[X_j Z_i] + \mathbb{E}[Z_i Z_j] \\
&= \text{Cov}(X_i, X_j) - \mathbb{E}[X_i \mathbb{E}[Z_j | (S_i, A_i)]] - \mathbb{E}[X_j \mathbb{E}[Z_i | (S_j, A_j)]] + \mathbb{E}[Z_i \mathbb{E}[Z_j | S_i, A_i]]
\end{aligned}$$

585 And we have seen that $\mathbb{E}[Z_i] = 0 \forall i$, hence

$$= \text{Cov}(X_i, X_j) + 0$$

586 Hence, $\text{Cov}(X_i, X_j) = \text{Cov}(Y_i, Y_j)$, hence if $\mathbb{V}[X] < \mathbb{V}[Y]$ then $\mathbb{V}[X_i + X_j] < \mathbb{V}[Y_i + Y_j]$, and hence
587 $\mathbb{V}[\hat{\Delta}(\theta, r_p)] < \mathbb{V}[\hat{\Delta}(\theta, \tilde{r})]$.

588 □

589 **Examples:** Let's look at some example where the above condition can be true, first, lets consider that we
590 have a parameterizing where we separately update the policies for each state, i.e. a tabular representation then,
591 in that case, we can consider the variance of the update to the policy at each state separately i.e.

$$\mathbb{V}_\pi[Y|s] - \mathbb{V}_\pi[X|s] = \Phi(s)^2 \mathbb{E}_\pi[\psi_\theta(s, A)^T \psi_\theta(s, A)] - 2\Phi(s) \mathbb{E}_\pi[\psi_\theta(s, A)^T \psi_\theta(s, A) (q^\pi(s, A))]$$

592 Hence let's see under what conditions we might get that the variance of the potential-based method might be
593 more than the variance only on primary reward, i.e. for that to be true, the difference between the 2 terms should

594 be positive, i.e.

$$\begin{aligned} \Phi(s)^2 \mathbb{E}_\pi \left[\psi_\theta(s, A)^T \psi_\theta(s, A) \right] - 2\Phi(s) \mathbb{E}_\pi \left[\psi_\theta(s, A)^T \psi_\theta(s, A) (q^\pi(s, A)) \right] &> 0 \\ \Phi(s)^2 \mathbb{E}_\pi \left[\psi_\theta(s, A)^T \psi_\theta(s, A) \right] &> 2\Phi(s) \mathbb{E}_\pi \left[\psi_\theta(s, A)^T \psi_\theta(s, A) (q^\pi(s, A)) \right] \end{aligned}$$

595 Further, let's consider the case where $\Phi(s) \neq 0$, cause otherwise for those states, the variance would be the
596 same, and $\Phi(s) > 0$.

$$\begin{aligned} \cancel{\Phi(s)^2} \mathbb{E}_\pi \left[\psi_\theta(s, A)^T \psi_\theta(s, A) \right] &> \cancel{2\Phi(s)} \mathbb{E}_\pi \left[\psi_\theta(s, A)^T \psi_\theta(s, A) (q^\pi(s, A)) \right] \\ \Phi(s) \mathbb{E}_\pi \left[\psi_\theta(s, A)^T \psi_\theta(s, A) \right] &> 2 \mathbb{E}_\pi \left[\psi_\theta(s, A)^T \psi_\theta(s, A) (q^\pi(s, A)) \right] \end{aligned}$$

597 We can see that the above condition can easily be satisfied by choosing a potential function that might be overly
598 optimistic about the action-value of the state s , i.e. any $\Phi(s)$, s.t. $\Phi(s) > 2q^\pi(s, a) \forall a$ would lead to an increase
599 in variance. A good example of this could be using an optimal value function (as hinted by [42]) as a baseline
600 for a bad/mediocre policy initially.

601 **Property 2.** *There exists $r_\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ and $\gamma_\varphi \in [0, 1)$ such that $\Delta_{on}(\theta, \phi, \varphi) = \Delta_\gamma(\theta, r_p)$.*

602 *Proof.* Recall the definition of $\Delta_\gamma(\theta, r_p)$ from Section 3.1,

$$\Delta_\gamma(\theta, r_p) = \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \gamma^t \psi_\theta(S_t, A_t) \sum_{j=t}^T \gamma^{j-t} r_p(S_j, A_j) \right].$$

603 Using the law of total expectation,

$$\begin{aligned} \Delta_\gamma(\theta, r_p) &= \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \gamma^t \psi_\theta(S_t, A_t) \mathbb{E}_{\pi_\theta} \left[\sum_{j=t}^T \gamma^{j-t} r_p(S_j, A_j) \middle| S_t, A_t \right] \right] \\ &= \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \gamma^t \psi_\theta(S_t, A_t) q^{\pi_\theta}(S_t, A_t) \right] \\ &= \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \sum_{t=0}^T \gamma^t \Pr(S_t = s, A_t = a; \pi_\theta) \psi_\theta(s, a) q^{\pi_\theta}(s, a) \\ &= \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \psi_\theta(s, a) q^{\pi_\theta}(s, a) \sum_{t=0}^T \gamma^t \Pr(S_t = s, A_t = a; \pi_\theta) \\ &= \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \psi_\theta(s, a) q^{\pi_\theta}(s, a) d_\gamma^{\pi_\theta}(s, a). \end{aligned} \tag{11}$$

604 Notice from (9) and (10) that for any (s, a) pair, if $d_\gamma^{\pi_\theta}(s, a) > 0$, then it has to be that $\bar{d}^{\pi_\theta}(s, a) > 0$ as well
605 since $\gamma \geq 0$. Therefore, dividing and multiplying by $d^{\pi_\theta}(s, a)$,

$$\begin{aligned} \Delta_\gamma(\theta, r_p) &= \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \bar{d}^{\pi_\theta}(s, a) \psi_\theta(s, a) q^{\pi_\theta}(s, a) \frac{d_\gamma^{\pi_\theta}(s, a)}{\bar{d}^{\pi_\theta}(s, a)} \\ &= \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \sum_{t=0}^T \Pr(S_t = s, A_t = a; \pi_\theta) \psi_\theta(s, a) q^{\pi_\theta}(s, a) \frac{d_\gamma^{\pi_\theta}(s, a)}{\bar{d}^{\pi_\theta}(s, a)} \\ &= \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \psi_\theta(S_t, A_t) q^{\pi_\theta}(S_t, A_t) \frac{d_\gamma^{\pi_\theta}(S_t, A_t)}{\bar{d}^{\pi_\theta}(S_t, A_t)} \right]. \end{aligned}$$

606 Now, notice that if $\gamma_\varphi = 0$ and $r_\phi(s, a) = q^{\pi_\theta}(s, a) \frac{d_\gamma^{\pi_\theta}(s, a)}{\bar{d}^{\pi_\theta}(s, a)}$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$, then

$$\Delta_{on}(\theta, \phi, \varphi) = \Delta_\gamma(\theta, r_p).$$

607

□

608 **Property 3.** *There exists $r_\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ and $\gamma_\varphi \in [0, 1)$ such that $\Delta_{off}(\theta, \phi, \varphi) = \Delta_{off}(\theta, r_p)$.*

609 *Proof.* This proof follows a similar technique as the proof for Property 2. Recall the definition of $\Delta_{\text{off}}(\theta, r_p)$,

$$\begin{aligned}\Delta_{\text{off}}(\theta, r_p) &:= \mathbb{E}_\beta \left[\sum_{t=0}^T \gamma^t \psi_\theta(S_t, A_t) \sum_{j=t}^T \gamma^{j-t} \rho_j r_p(S_j, A_j) \right] \\ &:= \mathbb{E}_\beta \left[\sum_{t=0}^T \gamma^t \rho_t \psi_\theta(S_t, A_t) \sum_{j=t}^T \gamma^{j-t} \rho_{j-t} r_p(S_j, A_j) \right]\end{aligned}$$

610 Now using the law of total expectations,

$$\begin{aligned}\Delta_{\text{off}}(\theta, r_p) &= \mathbb{E}_\beta \left[\sum_{t=0}^T \gamma^t \rho_t \psi_\theta(S_t, A_t) \mathbb{E}_\beta \left[\sum_{j=t}^T \gamma^{j-t} \rho_{j-t} r_p(S_j, A_j) \middle| S_t, A_t \right] \right] \\ &= \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \gamma^t \psi_\theta(S_t, A_t) \mathbb{E}_{\pi_\theta} \left[\sum_{j=t}^T \gamma^{j-t} r_p(S_j, A_j) \middle| S_t, A_t \right] \right] \\ &= \mathbb{E}_{\pi_\theta} \left[\sum_{t=0}^T \gamma^t \psi_\theta(S_t, A_t) q^{\pi_\theta}(S_t, A_t) \right] \\ &= \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \psi_\theta(s, a) q^{\pi_\theta}(s, a) d_\gamma^{\pi_\theta}(s, a),\end{aligned}$$

611 where the last line follows similar to (11). Now, notice that for any (s, a) pair, the assumption that
612 $\pi_\theta(s, a)/\beta(s, a) < \infty$ for all $s \in \mathcal{S}, a \in \mathcal{A}$, implies $d_\gamma^{\pi_\theta}(s, a)/d_\gamma^\beta(s, a) < \infty$. Further, if $d_\gamma^\beta(s, a) > 0$
613 it has to be that $d^\beta(s, a) > 0$ as well. Therefore, $d_\gamma^{\pi_\theta}(s, a)/d^\beta(s, a) < \infty$ as well. Multiplying and dividing by
614 $d^\beta(s, a)$,

$$\begin{aligned}\Delta_{\text{off}}(\theta, r_p) &= \sum_{s \in \mathcal{S}, a \in \mathcal{A}} d^\beta(s, a) \psi_\theta(s, a) q^{\pi_\theta}(s, a) \frac{d_\gamma^{\pi_\theta}(s, a)}{d^\beta(s, a)} \\ &= \sum_{s \in \mathcal{S}, a \in \mathcal{A}} \sum_{t=0}^T \Pr(S_t = s, A_t = a; \beta) \psi_\theta(s, a) q^{\pi_\theta}(s, a) \frac{d_\gamma^{\pi_\theta}(s, a)}{d^\beta(s, a)} \\ &= \mathbb{E}_\beta \left[\sum_{t=0}^T \psi_\theta(S_t, A_t) q^{\pi_\theta}(S_t, A_t) \frac{d_\gamma^{\pi_\theta}(S_t, A_t)}{d^\beta(S_t, A_t)} \right].\end{aligned}$$

615 Now, notice that if $\gamma_\phi = 0$ and $r_\phi(s, a) = q^{\pi_\theta}(s, a) \frac{d_\gamma^{\pi_\theta}(s, a)}{d^\beta(s, a)}$ for all $s \in \mathcal{S}$ and $a \in \mathcal{A}$,

$$\Delta_{\text{off}}(\theta, \phi, \varphi) = \Delta_{\text{off}}(\theta, r_p).$$

616

□

617 B Algorithm

618 In this section we discuss the algorithm for the proposed method. As the proposed method does behavior
619 alignment reward function's implicit optimization, we name our method Barfi. Pseudo-code for Barfi is
620 presented in Algorithm 5. We will first build on some preliminaries to understand the concepts

621 B.1 Vector Jacobian Product

622 Lets assume the following, $x \in \mathbb{R}^d, y \in \mathbb{R}^m, f(x, y) \in \mathbb{R}$, then we know that $\partial f(x, y)/\partial x \in$
623 $\mathbb{R}^d, \partial f(x, y)/\partial y \in \mathbb{R}^m, \partial^2 f(x, y)/\partial y \partial x \in \mathbb{R}^{d \times m}$. Let us also assume that we have a vector $v \in \mathbb{R}^d$,
624 and when we need to calculate the following

$$\underbrace{\underbrace{v}_{\mathbb{R}^d} \underbrace{\frac{\partial^2 f(x, y)}{\partial y \partial x}}_{\mathbb{R}^{d \times m}}}_{\mathbb{R}^m} = \frac{\partial}{\partial y} \underbrace{\left\langle \underbrace{v}_{\mathbb{R}^d}, \underbrace{\frac{\partial f(x, y)}{\partial x}}_{\mathbb{R}^d} \right\rangle}_{\mathbb{R}^1}}_{\mathbb{R}^m}$$

625 As we can see the vector jacobian product can actually be broken down into differentiating a vector product but
626 shifting the place of multiplication, in which case we assume that the gradient simply passes through v w.r.t. y
627 and hence we don't ever have to deal with large multiplications. Also note, that the outer partial w.r.t can easily
628 be handled by `autodiff` packages. A simple pseudo-code is show in Algorithm 1.

Algorithm 1: Jacobian Vector Product

1 Input: $f(x, y) \in \mathbb{R}^1, x \in \mathbb{R}^d, y \in \mathbb{R}^m, v \in \mathbb{R}^d$
2 $f' \leftarrow \text{grad}(f(x, y), x)$
3 $jvp \leftarrow \text{grad}(f', y, \text{grad_outputs} = v)$
4 Return: jvp

629 B.2 Neumann Series Approximation for Hessian Inverse

630 Recall, that for a given real number $\beta \in \mathbb{R}$, such that $0 \leq \beta < 1$, we know that the geometric series of this has
631 a closed form solution, i.e.

$$\begin{aligned} s &= 1 + \beta^1 + \beta^2 + \beta^3 + \dots + \\ &= \frac{1}{1 - \beta} \end{aligned}$$

632 Similarly given we have a value α , such that $\beta = 1 - \alpha$, we can write α^{-1} as follows

$$\begin{aligned} \frac{1}{1 - \beta} &= 1 + \beta + \beta^2 + \beta^3 + \dots + \\ \frac{1}{1 - (1 - \alpha)} &= 1 + (1 - \alpha) + (1 - \alpha)^2 + (1 - \alpha)^3 + \dots + \\ \alpha^{-1} &= 1 + (1 - \alpha) + (1 - \alpha)^2 + (1 - \alpha)^3 + \dots + \\ \alpha^{-1} &= \sum_{i=0}^{\infty} (1 - \alpha)^i \end{aligned}$$

633 The same can be generalized for a matrix, i.e. given a matrix $\mathbf{A} \in \mathbb{R}^{d \times d}$, we can write \mathbf{A}^{-1} as follows

$$\mathbf{A}^{-1} = \sum_{i=0}^{\infty} (\mathbf{I} - \mathbf{A})^i$$

634 Note for the above to hold, matrix \mathbf{A} , where we represent $\text{eig}(\mathbf{A})$ as the eigen values of matrix \mathbf{A} , we should
635 have the following condition to hold, $0 < \text{eig}(\mathbf{A}) < 1$. Note here we would regularize \mathbf{A} to ensure that all
636 eigen values are positive, and then we can always scale the matrix \mathbf{A} , by its biggest eigen value to ensure that
637 the above condition holds. Lets say $\eta = 1/\max \text{eig}(\mathbf{A})$, then we can write the following

$$\begin{aligned} \mathbf{A}^{-1} &= \frac{\eta}{\eta} \mathbf{A}^{-1} \\ &= \eta (\eta \mathbf{A})^{-1} \\ &= \eta \sum_{i=0}^{\infty} (\mathbf{I} - \eta \mathbf{A})^i \end{aligned}$$

638 As $\eta \mathbf{A}$ would always satisfy the above condition.

639 B.3 Neumann Approximation for Hessian Vector Product

640 Given we have seen how we can approximate the Inverse of a matrix without relying $O(d^3)$ operations, through
641 Neumann approximation, lets look what needs to be done for our updates. Recall that the update ϕ, φ (7) and (8)
642 were,

$$\frac{\partial J(\theta(\phi, \varphi))}{\partial \phi} = - \underbrace{\frac{\partial J(\theta(\phi, \varphi))}{\partial \theta(\phi, \varphi)}}_v \left(\underbrace{\frac{\partial \Delta(\theta(\phi, \varphi), \phi, \varphi)}{\partial \theta(\phi, \varphi)}}_{\mathbf{H}} \right)^{-1} \underbrace{\frac{\partial \Delta(\theta(\phi, \varphi), \phi, \varphi)}{\partial \phi}}_{\mathbf{A}}$$

643 and

$$\frac{\partial (J(\theta(\phi, \varphi)) - \frac{1}{2} \|\gamma_\varphi\|^2)}{\partial \varphi} = - \underbrace{\frac{\partial J(\theta(\phi, \varphi))}{\partial \theta(\phi, \varphi)}}_v \left(\underbrace{\frac{\partial \Delta(\theta(\phi, \varphi), \phi, \varphi)}{\partial \theta(\phi, \varphi)}}_{\mathbf{H}} \right)^{-1} \underbrace{\frac{\partial \Delta(\theta(\phi, \varphi), \phi, \varphi)}{\partial \varphi}}_{\mathbf{B}} - \frac{\partial \gamma_\varphi}{\partial \varphi}.$$

644 Lets look closely at the update for ϕ and we can generalize things easily for the case of φ .

$$\frac{\partial J(\theta(\phi, \varphi))}{\partial \phi} = -v\mathbf{H}^{-1}\mathbf{A}$$

645 We first look at how we can approximate the value of $v\mathbf{H}^{-1}$ efficiently as we can always make use of the
 646 Jacobian Vector product later to get $(v\mathbf{H}^{-1})\mathbf{A}$, as $v\mathbf{H}^{-1}$ becomes a vector. Lets assume we wish to run the
 647 Neumann approximation upto n steps, i.e. we want to approximate \mathbf{H}^{-1} upto n order Neumann expansion,

$$\eta(\eta\mathbf{H}^{-1}) \approx \eta \sum_{i=0}^n (I - \eta\mathbf{H})^i \quad (12)$$

648 Over here we are assuming that the outer optimization for update (1) is for the function $J(\theta(\phi, \varphi))$ and the inner
 649 optimization which is represented by the update (4) is $f(\theta(\phi, \varphi), \phi, \varphi)$, i.e.

$$\Delta(\theta, r_p) = \frac{\partial J(\theta(\phi, \varphi))}{\partial \theta}$$

$$\Delta(\theta, \phi, \varphi) = \frac{\partial f(\theta(\phi, \varphi), \phi, \varphi)}{\partial \theta}$$

650 The most common form in which $f(\cdot; B)$ is usually defined is the following:

$$f(\theta, \phi, \varphi; B) := \frac{1}{|B|} \sum_{\tau \in B} \left[\sum_{t=0}^T \log(\pi_{\theta}(S_t^{\tau}, A_t^{\tau})) \sum_{j=t}^T \gamma_{\varphi}^{j-t} r_{\phi}(S_j^{\tau}, A_j^{\tau}) \right].$$

651 Similalry this can be defined for J , except making use of r_p and the actual γ ,

$$J(\theta; B) := \frac{1}{|B|} \sum_{\tau \in B} \left[\sum_{t=0}^T \log(\pi_{\theta}(S_t^{\tau}, A_t^{\tau})) \sum_{j=t}^T \gamma^{j-t} r_p(S_j^{\tau}, A_j^{\tau}) \right].$$

Algorithm 2: Vector Hessian Inverse Product for (7) i.e. $v\mathbf{H}^{-1}$

```

1 Input:  $\theta(\phi, \varphi), \phi, \varphi, J, f, n, \eta, \mathcal{D}_{\text{off}}, \mathcal{D}_{\text{on}}$ 
2  $v \leftarrow \text{grad}(J(\theta(\phi, \varphi); \mathcal{D}_{\text{on}}), \theta)$ 
3  $v' \leftarrow \eta \times \text{grad}(f(\theta(\phi, \varphi), \phi, \varphi; \mathcal{D}_{\text{off}}), \theta)$ 
4 Let:  $v_0 \leftarrow v, p_0 \leftarrow v$ 
5 for  $i \in [0, n)$  do
6    $v_{i+1} \leftarrow v_i - \text{grad}(v', \theta, \text{grad\_outputs} = v_i)$ 
7    $p_{i+1} \leftarrow p_i + v_{i+1}$ 
8 Return:  $\eta p_n$ ; // Approximation of  $v\mathbf{H}^{-1}$  as in (12)
```

652

653 Finally once we have $v\mathbf{H}^{-1}$ we can use the Vector Jacobian Product to calculate $(v\mathbf{H}^{-1})\mathbf{A}$ as follows:

Algorithm 3: Update for ϕ , i.e. (7) i.e. $v\mathbf{H}^{-1}\mathbf{A}$

```

1 Input:  $\theta(\phi, \varphi), \phi, \varphi, J, f, n, \eta, \mathcal{D}_{\text{off}}, \mathcal{D}_{\text{on}}$ 
2  $v \leftarrow \text{Algorithm 2}(\theta(\phi, \varphi), \phi, \varphi, J, f, n, \eta, \mathcal{D}_{\text{off}}, \mathcal{D}_{\text{on}})$ 
3  $v' \leftarrow \text{grad}(f(\theta(\phi, \varphi), \phi, \varphi; \mathcal{D}_{\text{off}}), \theta)$ 
4  $\Delta_{\phi} \leftarrow \text{grad}(v', \phi, \text{grad\_outputs} = v)$ 
5 Return  $\Delta_{\phi}$ 
```

654 We can similarly derive updates for φ . Note we are not including the different forms of regularizers over here to
 655 reduce clutter, but adding them is simple.

656 B.4 Pseduo Code

657 Line 8–10 and 21–23 represent the inner optimization process, and the outer optimization process if from lines
 658 16–17. Lines 8–10 is the initial step of updates to converge to the current values of ϕ, φ , and from there onwards
 659 after each update of outer optimization, we consequently update the policy in (21–23). The flow of the algorithm
 660 is show in Figure 6.

Algorithm 4: Update for ϕ , i.e. (8) i.e. $v\mathbf{H}^{-1}\mathbf{B}$

```
1 Input:  $\theta(\phi, \varphi), \phi, \varphi, J, f, n, \eta, \mathcal{D}_{\text{off}}, \mathcal{D}_{\text{on}}$ 
2  $v \leftarrow \text{Algorithm 2}(\theta(\phi, \varphi), \phi, \varphi, J, f, n, \eta, \mathcal{D}_{\text{off}}, \mathcal{D}_{\text{on}})$ 
3  $v' \leftarrow \text{grad}(f(\theta(\phi, \varphi), \phi, \varphi), \mathcal{D}_{\text{off}}), \theta)$ 
4  $\Delta_\varphi \leftarrow \text{grad}(v', \varphi, \text{grad\_outputs} = v)$ 
5 Return  $\Delta_\varphi$ 
```

Algorithm 5: Barfi: Behavior Alignment Reward Function's Implicit optimization

```
1 Input:  $J, f, \alpha_\theta, \alpha_\phi, \alpha_\varphi, \eta, n, \delta, \text{optim}, E, N_i, N_0$ ,
2 Initialize:  $\pi_\theta, r_\phi, \gamma_\varphi$ 
3 Initialize:  $\text{optim}_\theta \leftarrow \text{optim}(\alpha_\theta), \text{optim}_\phi \leftarrow \text{optim}(\alpha_\phi), \text{optim}_\varphi \leftarrow \text{optim}(\alpha_\varphi)$ 
4  $\mathcal{D}_{\text{off}} \leftarrow \{\}$ 
5 for  $e \in [1, N_0]$  do
6   Generate  $\tau_e$  using  $\pi_\theta$ 
7   Save  $\tau_e$  in  $\mathcal{D}_{\text{off}}$ 
8 for  $i \in [0, N_i + N_0]$  do
9   Sample a batch of trajectories  $B$  from  $\mathcal{D}_{\text{off}}$ 
10   $\theta \leftarrow \theta + \text{optim}_\theta(\text{grad}(f(\theta, \phi, \varphi; B), \theta))$ ; // Make use of  $B$  to perform
    update rule (4)
11 for  $e \in [N_0, E]$  do
    # Collect a batch of on-policy data
12   $\mathcal{D}_{\text{on}} \leftarrow \{\}$ 
13  for  $j \in [0, \delta)$  do
14    Generate trajectory  $\tau_{e+j}$  using  $\pi_\theta$  and save in  $\mathcal{D}_{\text{on}}$ 
15   $e \leftarrow e + \delta$ 
    # Update  $r_\phi$  and  $\gamma_\varphi$ 
16   $\Delta_\phi \leftarrow \text{Algorithm 3}(\theta(\phi, \varphi), \phi, \varphi, J, f, n, \eta, \mathcal{D}_{\text{off}}, \mathcal{D}_{\text{on}})$ 
17   $\Delta_\varphi \leftarrow \text{Algorithm 4}(\theta(\phi, \varphi), \phi, \varphi, J, f, n, \eta, \mathcal{D}_{\text{off}}, \mathcal{D}_{\text{on}})$ 
18   $\phi \leftarrow \phi + \text{optim}_\phi(\Delta_\phi)$ 
19   $\varphi \leftarrow \varphi + \text{optim}_\varphi(\Delta_\varphi)$ 
20   $\mathcal{D}_{\text{off}} \leftarrow \mathcal{D}_{\text{off}} \cup \mathcal{D}_{\text{on}}$ 
21  for  $i \in [0, N_i]$  do
22    Sample a batch of trajectories  $B$  from  $\mathcal{D}_{\text{off}}$ 
23     $\theta \leftarrow \theta + \text{optim}_\theta(\text{grad}(f(\theta, \phi, \varphi; B), \theta))$ ; // Make use of  $B$  to
        perform update rule (4)
```

661 As discussed in Section C, using regularizers in $\Delta(\theta, \phi, \varphi)$ smoothens the objective $J(\theta(\phi, \varphi))$ with respect to
662 ϕ and φ . This is helpful as gradual changes in r_ϕ and γ_φ can result in gradually changes in the fixed point for
663 the inner optimization. Therefore, for computational efficiency, we initialize the policy parameters from the
664 fixed-point of the previous inner-optimization procedure such that the inner-optimization process may start close
665 to the new fixed-point.

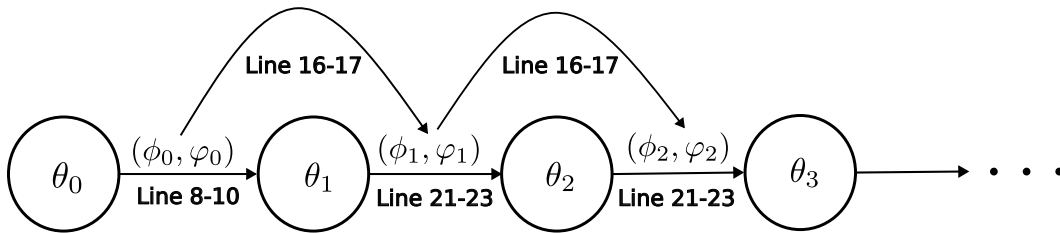


Figure 6: **Algorithm Flow:** The change in different parameters

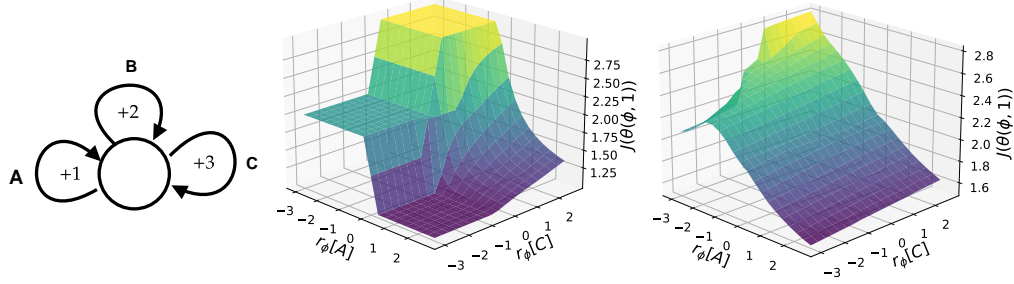


Figure 7: **(Left)** A bandit problem, where the data is collected from a policy β that samples action A mostly. **(Middle)** Each point on the 3D surface corresponds to the performance of $\theta(\phi, 1)$ returned by an `Alg` that uses the update rule $\Delta_{\text{off}}(\theta, \phi, 1)$ corresponding to the value of r_ϕ for actions A and C in the bottom axes; r_θ for action B is set to 0 to avoid another variable in a 3D plot. Notice that small perturbation in r_ϕ may lead to no or sudden changes in $J(\theta(\phi, 1))$. **(Right)** Performance of $\theta(\phi, 1)$ returned by an `Alg` that uses the update rule $\Delta_{\text{off}}(\theta, \phi, 1) - \theta$ that incorporates gradient of the L2 regularizer. Vector fields in Figure 1 were also obtained from this setup.

666 In lines 8–10, the inner optimization for the policy parameters θ are performed till (approximate) convergence.
 667 Note that only trajectories from past interactions are used and no new-trajectories are sampled for the inner
 668 optimization.

669 In Lines 13–14, a new batch \mathcal{D}_{on} of data is sampled using the policy returned by the inner-optimization process.
 670 This data is used to compute $\partial J(\theta(\phi, \varphi))/\partial \theta(\phi, \varphi)$. Existing data \mathcal{D}_{off} that was used in the inner-optimization
 671 process is then used to compute $\partial \theta(\phi, \varphi)/\partial \phi$ and $\partial \theta(\phi, \varphi)/\partial \varphi$. Using these in (7) and (8), the parameters for
 672 r_ϕ and γ_φ are updated in Lines 16 and 17, respectively.

673 Finally, the new data \mathcal{D}_{on} is merged into the existing data \mathcal{D}_{off} and the entire process continues.

674 C Smoothing the objective

675 To understand why $J(\theta(\phi, \varphi))$ might be ill-conditioned is to note that, often a small perturbation in the reward
 676 function doesn't necessarily lead to a change in the corresponding optimal policy. This can lead lack of gradient
 677 directions in the neighborhood of ϕ, φ for gradient methods to be effective. This issue can be addressed
 678 by employing common regularization techniques like L2 regularization of the policy parameters or entropy
 679 regularization for the policy. We discuss two ways to regularize the objective in the upcoming sections.

680 C.1 L2 Regularization

681 To understand how severely ill-conditioned $J(\theta(\phi, \varphi))$ can be, notice that a small perturbation in the reward
 682 function often does not change the corresponding optimal policies or the outcome of a policy optimization
 683 algorithm `Alg`. Therefore, if the parameters of the behavior alignment reward are perturbed from ϕ to ϕ' , it may
 684 often be that $J(\theta(\phi, \varphi)) = J(\theta(\phi', \varphi))$ and this limits any gradient based optimization for ϕ as $\partial J(\theta(\phi, \varphi))/\partial \phi$
 685 is 0. Similarly, minor perturbations in φ may result in no change in $J(\theta(\phi, \varphi))$ either.

686 Fortunately, there exists a remarkably simple solution: incorporate regularization for the *policy parameters* θ in
 687 objective for `Alg` in the inner-level optimization. For example, the optimal policy for the following regularized
 688 objective $\mathbb{E}_{\pi_\theta}[\sum_{t=0}^T \gamma_\varphi^t r_\phi(S_t, A_t)] - \frac{\lambda}{2} \|\theta\|^2$ varies smoothly to trade-off between the regularization value of θ
 689 and the magnitude of the performance characterized by (r_ϕ, γ_φ) , which changes with the values of r_ϕ and γ_φ .
 690 See Figure 7 for an example with L2 regularization.

691 C.2 Entropy Regularized

692 In Section C.1, smoothing of $J(\theta(\phi, \varphi))$ was done by using L2 regularization on the policy parameters θ in the
 693 inner-optimization process. However, alternate regularization methods can also be used. For example, in the
 694 following we present an alternate update rule for θ based on entropy regularization,

$$\Delta(\theta, \phi, \varphi) := \mathbb{E}_{\mathcal{D}} \left[\sum_{t=0}^T \psi_\theta(S_t, A_t) \sum_{j=t}^T \gamma_\varphi^{j-t} (r_\phi(S_j, A_j) - \lambda \ln \pi_\theta(S_j, A_j)) \right].$$

695 Notice that new update rule for ϕ and φ can be obtained from steps (5) to (8) with the following \mathbf{A} , \mathbf{B} , and \mathbf{H}
 696 instead, where for shorthand $\theta^* = \theta(\phi, \varphi)$,

$$\begin{aligned}\mathbf{A} &= \mathbb{E}_{\mathcal{D}} \left[\sum_{t=0}^T \psi_{\theta^*}(S_t, A_t) \left(\sum_{j=t}^T \gamma_{\varphi}^{j-t} \frac{r_{\phi}(S_j, A_j)}{\partial \phi} \right)^{\top} \right], \\ \mathbf{B} &= \mathbb{E}_{\mathcal{D}} \left[\sum_{t=0}^T \psi_{\theta^*}(S_t, A_t) \left(\sum_{j=t}^T \frac{\partial \gamma_{\varphi}^{j-t}}{\partial \varphi} (r_{\phi}(S_j, A_j) - \lambda \ln \pi_{\theta^*}(S_j, A_j)) \right) \right], \\ \mathbf{H} &= \mathbb{E}_{\mathcal{D}} \left[\sum_{t=0}^T \frac{\partial \psi_{\theta^*}(S_t, A_t)}{\partial \theta^*} \left(\sum_{j=t}^T \gamma_{\varphi}^{j-t} (r_{\phi}(S_j, A_j) - \lambda \ln \pi_{\theta^*}(S_j, A_j)) \right) - \lambda \psi_{\theta^*}(S_t, A_t) \left(\sum_{j=t}^T \gamma_{\varphi}^{j-t} \psi_{\theta^*}(S_j, A_j)^{\top} \right) \right].\end{aligned}$$

697

698 D Meta Learning via Implicit Gradient: Derivation

699 The general technique of implicit gradients [14, 34, 19] has been used in a vast range of applications, ranging
 700 from energy models [13, 36], differentiating through black-box solvers [61], few-shot learning [38, 49], model-
 701 based RL [50], differentiable convex optimization neural-networks layers [4, 2], to hyper-parameter optimization
 702 [37, 8, 12, 40]. In this work, we show how implicit gradients can also be useful to efficiently leverage auxiliary
 703 rewards r_a and overcome various sub-optimality.

704 Taking the derivative (6) $\Delta(\theta(\phi, \varphi), \phi, \varphi) = 0$ of the above point of convergence w.r.t to ϕ and φ we get the
 705 following updates.

706 Therefore, taking total derivative in (6) with respect to ϕ ,

$$\frac{d\Delta(\theta(\phi, \varphi), \phi, \varphi)}{d\phi} = \frac{\partial \Delta(\theta(\phi, \varphi), \phi, \varphi)}{\partial \phi} + \frac{\partial \Delta(\theta(\phi, \varphi), \phi, \varphi)}{\partial \theta(\phi, \varphi)} \frac{\partial \theta(\phi, \varphi)}{\partial \phi} = 0. \quad (13)$$

707 Lets try to understand why the above is true, considering the finite difference approach for this derivative,

$$\frac{d\Delta(\theta(\phi, \varphi), \phi, \varphi)}{d\phi} = \lim_{\|d\phi\| \rightarrow 0} \frac{\Delta(\theta(\phi + d\phi, \varphi), \phi + d\phi, \varphi) - \Delta(\theta(\phi, \varphi), \phi, \varphi)}{d\phi}$$

$$\begin{aligned}708 \quad \Delta(\theta(\phi + d\phi, \varphi), \phi + d\phi, \varphi) &= \Delta(\theta(\phi, \varphi), \phi, \varphi) = 0, \text{ as } \theta(\cdot, \cdot) \text{ defines convergence to fixed point} \\ &= \frac{0 - 0}{d\phi} = 0\end{aligned}$$

709 By re-arranging terms in (13) we obtain the term (b) in (5),

$$\frac{\partial \theta(\phi, \varphi)}{\partial \phi} = - \left(\frac{\partial \Delta(\theta(\phi, \varphi), \phi, \varphi)}{\partial \theta(\phi, \varphi)} \right)^{-1} \frac{\partial \Delta(\theta(\phi, \varphi), \phi, \varphi)}{\partial \phi}. \quad (14)$$

710 On combining (14) with (5) we obtain the desired gradient expression for ϕ ,

$$\frac{\partial J(\theta(\phi, \varphi))}{\partial \phi} = - \frac{\partial J(\theta(\phi, \varphi))}{\partial \theta(\phi, \varphi)} \underbrace{\left(\frac{\partial \Delta(\theta(\phi, \varphi), \phi, \varphi)}{\partial \theta(\phi, \varphi)} \right)^{-1}}_{\mathbf{H}} \underbrace{\frac{\partial \Delta(\theta(\phi, \varphi), \phi, \varphi)}{\partial \phi}}_{\mathbf{A}},$$

711 and following similar steps, it can be observed that the gradient expression for φ ,

$$\frac{\partial (J(\theta(\phi, \varphi)) - \frac{1}{2} \|\gamma_{\varphi}\|^2)}{\partial \varphi} = - \frac{\partial J(\theta(\phi, \varphi))}{\partial \theta(\phi, \varphi)} \underbrace{\left(\frac{\partial \Delta(\theta(\phi, \varphi), \phi, \varphi)}{\partial \theta(\phi, \varphi)} \right)^{-1}}_{\mathbf{H}} \underbrace{\frac{\partial \Delta(\theta(\phi, \varphi), \phi, \varphi)}{\partial \varphi}}_{\mathbf{B}} - \frac{\partial \gamma_{\varphi}}{\partial \varphi},$$

712 where using θ^* as a shorthand for $\theta(\phi, \varphi)$ the terms \mathbf{A} , \mathbf{B} and \mathbf{H} can be expressed as,

$$\begin{aligned}\mathbf{A} &= \mathbb{E}_{\mathcal{D}} \left[\sum_{t=0}^T \psi_{\theta^*}(S_t, A_t) \left(\sum_{j=t}^T \gamma_{\varphi}^{j-t} \frac{\partial r_{\phi}(S_j, A_j)}{\partial \phi} \right)^{\top} \right], \quad \mathbf{B} = \mathbb{E}_{\mathcal{D}} \left[\sum_{t=0}^T \psi_{\theta^*}(S_t, A_t) \left(\sum_{j=t}^T \frac{\partial \gamma_{\varphi}^{j-t}}{\partial \varphi} r_{\phi}(S_j, A_j) \right) \right], \\ \mathbf{H} &= \mathbb{E}_{\mathcal{D}} \left[\sum_{t=0}^T \frac{\partial \psi_{\theta^*}(S_t, A_t)}{\partial \theta^*} \left(\sum_{j=t}^T \gamma_{\varphi}^{j-t} r_{\phi}(S_j, A_j) \right) \right] - \lambda.\end{aligned}$$

These provide the necessary expressions for updating ϕ and φ in the outer loop. As **A** involves an outer product and **H** involves second derivatives, computing them *exactly* might not be practical when dealing with high-dimensions. Standard approximation techniques like conjugate-gradients or Neumann series can thus be used to make it more tractable [40]. In our experiments, we made use of the Neumann approximation to the Hessian Inverse vector product $(\mathbf{A}\mathbf{H})^{-1}$, which requires the same magnitude of resources as the baseline policy gradient methods that we build on top of.

Algorithm: Being based on implicit gradients, we call our method `Barfi`, shorthand for *behavior alignment reward function’s implicit* optimization. Overall, `Barfi` iteratively solves the bi-level optimization specified in (2) by alternating between using (4) till approximate converge of `Alg` to $\theta(\phi, \varphi)$ and then updating r_ϕ and γ_φ . Importantly, being based on (4) for sample efficiency, `Alg` leverages only the past samples and does *not* sample any new trajectories for the inner level optimization. Further, due to policy regularization which smoothens the objective as discussed in C, updates in r_ϕ and γ_φ changes the policy resulting from `Alg` gradually. Therefore, for compute efficiency, we start `Alg` from the policy obtained from the previous inner optimization, such that it is in proximity of the new fixed point. This allows `Barfi` to be both sample and compute efficient while solving the bi-level optimization iteratively online. Pseudo-code for `Barfi` and more details on the approximation techniques can be found in Appendix B.

E Details for the Empirical Results

E.1 Implementation Details

In this section we will briefly describe the implementation details around the different environments that were used.

GridWorld(GW): In the case of GridWorld we made use of the Fourier basis (of Order = 3) over the raw coordinates of agent position in the GridWorld. Details about this could be found in the `src/Utils/Basis.py` file.

MountainCar(MC): For this environment, to reduce the limitation because of the function approximator we used TileCoding [58], which offers a suitable representation for the MountainCar problem. We used 4 Tiles and Tilings of 5.

CartPole(CP): For CartPole also make use of Fourier Basis of (Order = 3), with linear function approximator on top of that.

MuJoCo(MJ): For this we made use of a neural network with 1 hidden layer of 32 nodes and ReLU activation as the function approximator over the raw observations. The output of the policy is continuous actions, hence we used a Gaussian representation, where the policy outputs the mean of the multivariate Gaussian and we used a fixed diagonal standard deviation, fixed to $\sigma = 0.1$.

General Details: All the outer returns are evaluated without any discounting, whereas all the inner optimizations were initialized with $\gamma_\varphi = 0.99$. Hence to do this we made φ a single bias unit, initialized to 4.6, and passed through a sigmoid (i.e., $\sigma(4.6) = 0.99$).

For GW, CP and MC r_ϕ is defined as below

$$r_\phi(s, a) = \phi_1(s) + \phi_2(s)r_p + \phi_3(s)r_a$$

Wherein ϕ_1, ϕ_2, ϕ_3 are scalar outputs of a 3-headed function, in this case simply a linear layer over the states inputs.

Whereas in the case of MJ, we have

$$r_\phi(s, a) = \phi_1 + r_p + \phi_3 r_a$$

Wherein ϕ_1 is initialized to zero and ϕ_3 is 1.0 act like bias units.

Gradient normalization was used for all the cases where Neural Nets were involved (i.e., MJ), and also for MJ we modified the Baseline (Reinforce) update to subtract the running average of the performance as a baseline to get acceptable performance for the baseline method.

E.2 Hyper-parameter Selection

As different make use of different function approximators hence the range of hyper-params can vary we talk about all the above over here. All the experiments were conducted on a personal computer with 32 GiB of memory and an Intel Core i7 CPU with 12 threads. Total runtime for all the experiments combined was less than a day.

Best-performing Parameters for different methods and environments are listed where

Hyper Parameter	Barfi Value	Reinforce Value	Actor Critic Value
α_θ	1×10^{-3}	1×10^{-3}	1×10^{-3}
α_ϕ	5×10^{-3}	—	—
α_φ	5×10^{-3}	—	—
optim	RMSprop	RMSprop	RMSprop
λ_θ	0.25	0.25	0.25
λ_ϕ	0.0625	—	—
λ_φ	4.0	—	—
Buffer	1000	—	—
Batch Size	1	1	1
η	0.0005	—	—
δ	3	—	—
n	5	—	—
N_0	150	—	—
N_i	15	—	—

Table 2: Hyper parameters for GridWorld

Hyper Parameter	Barfi Value	Reinforce Value	Actor Critic Value
α_θ	0.015625	0.125	0.03125
α_ϕ	0.0625	—	—
α_φ	0.0625	—	—
optim	RMSprop	RMSprop	RMSprop
λ_θ	0.0	0.0	0.25
λ_ϕ	0.0	—	—
λ_φ	0.25	—	—
Buffer	50	—	—
Batch Size	1	1	1
η	0.001	—	—
δ	3	—	—
n	5	—	—
N_0	50	—	—
N_i	15	—	—

Table 3: Hyper parameters for MountainCar

Hyper Parameter	Barfi Value	Reinforce Value	Actor Critic Value
α_θ	1×10^{-3}	1×10^{-3}	5×10^{-4}
α_ϕ	1×10^{-3}	—	—
α_φ	5×10^{-3}	—	—
optim	RMSprop	RMSprop	RMSprop
λ_θ	1.0	1.0	0.0
λ_ϕ	0.0	—	—
λ_φ	4.0	—	—
Buffer	10000	—	—
Batch Size	1	1	1
η	0.0005	—	—
δ	3	—	—
n	5	—	—
N_0	150	—	—
N_i	15	—	—

Table 4: Hyper parameters for CartPole

762 **Hyperparameter Sweep** : Here we list the details about how we swept the values for different hyper-params.
763 We used PyTorch [47] for all our implementations. We usually used an optimizer between RMSProp or Adam

Hyper Parameter	Barfi Value	Reinforce Value	Actor Critic Value
α_θ	7.5×10^{-5}	5×10^{-4}	2.5×10^{-4}
α_ϕ	2.5×10^{-3}	—	—
α_φ	0.0	—	—
optim	Adam	Adam	Adam
λ_ϕ	0.0625	—	—
λ_φ	0.0	—	—
Buffer	50	—	—
Batch Size	1	1	1
η	0.0005	—	—
δ	3	—	—
n	5	—	—
N_0	30	—	—
N_i	15	—	—

Table 5: Hyper parameters for MuJoco

with default parameters as provided in Pytorch. For $\alpha_\theta \in \{5 \times 10^{-3}, 2.5 \times 10^{-3}, 1 \times 10^{-3}, 5 \times 10^{-4}, 2.5 \times 10^{-4}, 1 \times 10^{-4}, 7.5 \times 10^{-5}\}$, we use similar ranges for $\alpha_\phi, \alpha_\varphi$ (which tend to be larger). For λ_θ and λ_ϕ , we swept from $[0, 0.25, 0.5, 1.0]$ and for λ_γ we swept from $[0, 0.25, 1.0, 4.0, 16.0]$. We simply list ranges for different values and later we present sensitivity curves showing that these values are usually robust for Barfi across different methods as we can see from the tables above. $\delta \in [1, 3, 5]$, $n \in [1, 3, 5]$, $N_i \in [1, 3, 6, 9, 12, 15]$, $\eta \in [1 \times 10^{-3}, 5 \times 10^{-4}, 1 \times 10^{-4}]$, $N_0 \in [30, 50, 100, 150]$, buffer $\in [25, 50, 100, 1000]$. α for Tilecoding was adopted from [58] and hence similar ranges were swept in that case. Most sweeps were done with around 10 seeds, and later the parameter ranges were reduced and performed with more seeds.

E.3 Compute

The computer is used for a cluster where the CPU class is Intel Xeon Gold 6240 CPU @2.60GHz. The total compute required for GW was around 3 CPU years², CP also required around 3 CPU years, and MC required around 4 CPU years. For MJ we needed around 5-6 CPU years. In total we utilized around 15-16 CPU years, where we needed around 1 GB of memory per thread.

F Extra Results & Ablations

Figure 8 and Figure 9 summarize the return based on r_ϕ and the γ learned by the agent across different domains and reward specification. We observe that Reinforce often optimizes the naive combination of reward for sure, but that doesn't really lead to a good performance on r_p , whereas Barfi does achieve appropriate return on r_ϕ , but is also able to successively decay γ as the learning progress across different domains. Particularly notice Figure 8 (a) Bottom, where Reinforce does optimize aux return a lot, but actually fails to solve the problem, as it simply learns to loop around the center state.

Ablations: Figure 10 represents the ablation of Barfi on GridWorld with the misspecified reward for its different params. We can see that usually having $\eta = 0.001, 0.0005$, $n = 5$ works for the approximation.

²1 CPU year := Compute equal to running a CPU thread for a year.

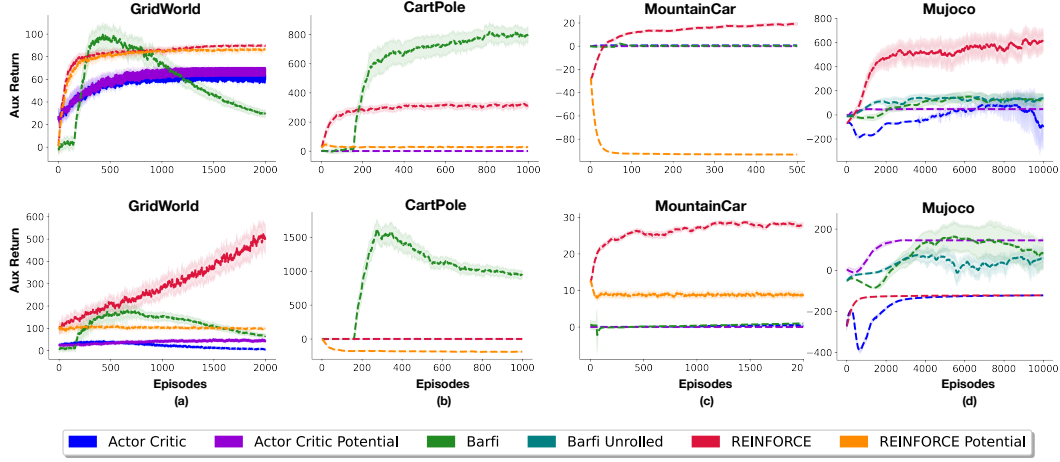


Figure 8: **Learned Reward Returns:** This figure illustrates the aux return collected by agent based on the learned r_ϕ , the curves are chosen based on best performing curves on r_p , and averaged over 20 runs (except 40 for GW). (a) Top – $r_{aux,GW}^1$, Bottom – $r_{aux,GW}^2$, (b) Top – $r_{aux,CP}^1$, Bottom – $r_{aux,CP}^2$, (c) Top – $r_{aux,MC}^2$, Bottom – $r_{aux,MC}^1$, (d) Top – $r_{aux,MJ}^1$, Bottom – $r_{aux,MJ}^2$.

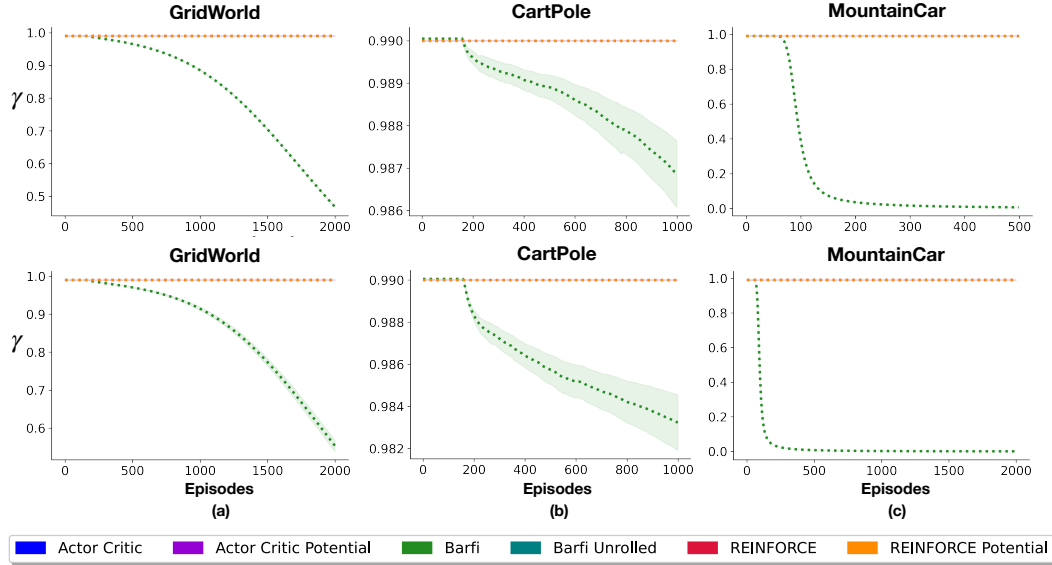


Figure 9: **Learned γ_ϕ :** This figure illustrates the learned γ_ϕ for **Barfi** and normal γ for other methods, the curves are chosen based on best-performing curves on r_p , and averaged over 20 runs (except 40 for GW). (a) Top – $r_{aux,GW}^1$, Bottom – $r_{aux,GW}^2$, (b) Top – $r_{aux,CP}^1$, Bottom – $r_{aux,CP}^2$, (c) Top – $r_{aux,MC}^2$, Bottom – $r_{aux,MC}^1$. MJ is not included as the γ was not learned in that case. We can observe that the agents start to learn to decay γ at the appropriate pace.

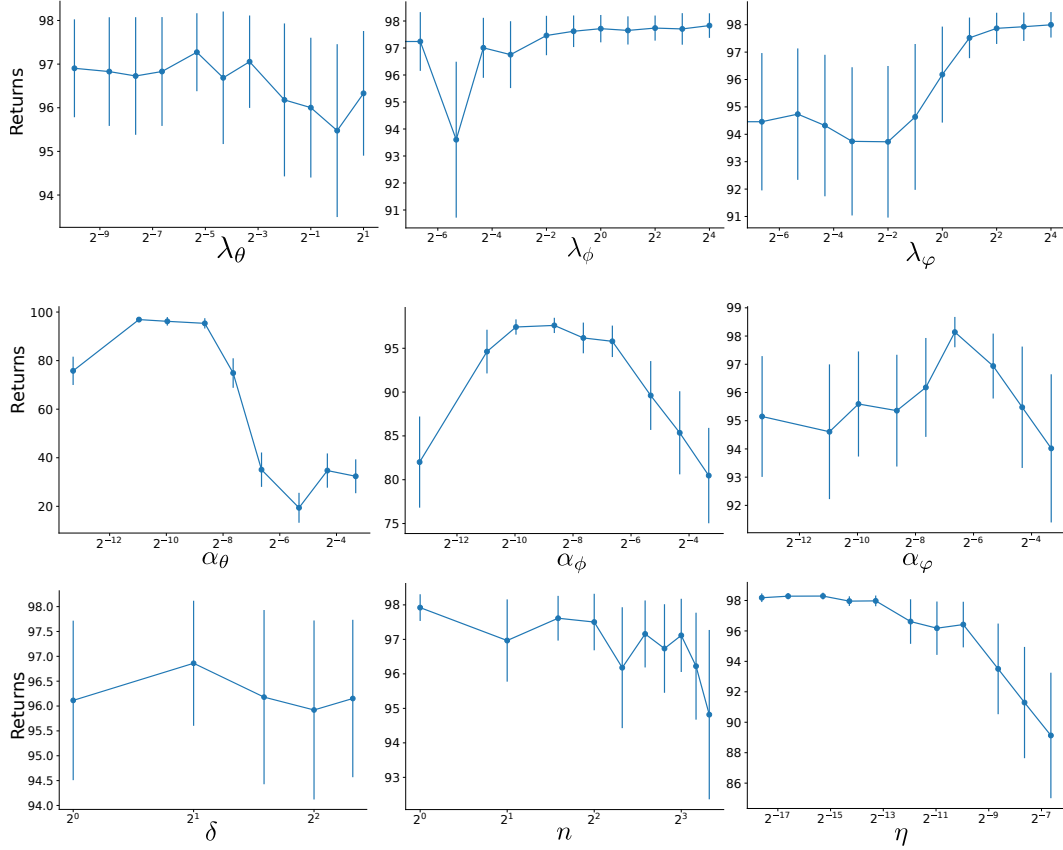


Figure 10: **Sensitivity Curves:** The set of graphs representing the sensitivity of different hyper-params keeping all the other params fixed. The sensitivity is for **Barfi** in GW with $r_{\text{aux, GW}}^2$, i.e., the misspecified reward. We choose the best-performing parameters and vary each parameter to see its influence. The curves are obtained for 50 runs (seeds) in each case, and error bars are standard errors. We can notice that α_θ and α_ϕ can have a large influence, and tend to stay around similar values. $\lambda_{\theta, \phi, \varphi}$ tends to help but doesn't really influence a lot in terms of its magnitude, except larger values of λ_φ seem to do better. Smaller values of η seems to work fine, hence something around 5×10^{-4} , 1×10^{-3} usually should suffice. n, δ can be chosen to around 5 and 3, and usually workout fine. We also defined $N_i = 5 \times \delta$ in this case.