

1 **Reviewers # 5, 8:** Thank you for the appreciation! The best known lower bound is $\Omega(\sqrt{DSA/T})$, based on [26] for
2 the SO-OMDP setting. Our upper bound $\tilde{O}(D\sqrt{\Gamma SA/T})$ in Theorem 3.1 matches the $\tilde{O}(DS\sqrt{A/T})$ bound by [26]
3 for SO-OMDP. The best known upper bound for SO-OMDP is $\tilde{O}(c\sqrt{\Gamma SA/T})$ by [21], where $c \leq D$ is called the *span*,
4 a refined version of D in the SO-OMDP setting. The notion of span is inapplicable to MO-OMDP.

5 **Reviewer # 6:** Thank your for the comments! For a better appreciation on our contributions, we clarify as follows:

Justifying the objective function g_{MO} (5.1), (2.2), (2.3). We start by addressing (5.1). KPI stands for Key Performance Index. For *Target Set Objectives*, specifying $U = \{w : w_k \geq \rho_k \forall 1 \leq k \leq K\} = \prod_{k=1}^K [\rho_k, \infty)$ is sufficient for ensuring $\bar{V}_{1:T,k} \geq \rho_k$ whenever possible, thanks to the $\min_{u \in U}$ operator in (1). To see this, consider setting $L_1 = \dots = L_K = 0, L_0 = 1$. We claim that $g_{\text{MO}}(\bar{V}_{1:T}) = -(1/2K) \sum_{k=1}^K \max\{\rho_k - \bar{V}_{1:T,k}, 0\}^2$. Indeed,

$$g_{\text{MO}}(\bar{V}_{1:T}) = -\frac{1}{2K} \min_{u \in \prod_{k=1}^K [\rho_k, \infty)} \left\{ \sum_{k=1}^K (\bar{V}_{1:T,k} - u_k)^2 \right\} = -\frac{1}{2K} \sum_{k=1}^K \min_{u_k \in [\rho_k, \infty)} \{(\bar{V}_{1:T,k} - u_k)^2\}.$$

6 For the k th summand, if $\bar{V}_{1:T,k} \geq \rho_k$, the argmin is $\bar{V}_{1:T,k}$ and the summand = 0. Else, we have $\bar{V}_{1:T,k} < \rho_k$, the
7 argmin is ρ_k and the summand = $(\rho_k - \bar{V}_{1:T,k})^2$. Thus, the claim is proved.

8 Maximizing $g_{\text{MO}}(\bar{V}_{1:T})$ is equivalent to minimizing $(1/2K) \sum_{k=1}^K \max\{\rho_k - \bar{V}_{1:T,k}, 0\}^2$. If the KPI ρ is achievable,
9 then the optimal policy would generate $\bar{V}_{1:T}$ for which $\bar{V}_{1:T,k} \geq \rho_k$ for all k , yielding objective value $g_{\text{MO}}(\bar{V}_{1:T}) = 0$.
10 Otherwise, the shortfall of $\bar{V}_{1:T}$ compared to ρ is minimized in the mean squared error sense.

11 (2.2): Any maximizer $\bar{V}_{1:T}^*$ of $g_{\text{MO}}(\bar{V}_{1:T}) = -(1/2K) \sum_{k=1}^K \max\{1 - \bar{V}_{1:T,k}, 0\}^2$ is Pareto-optimal. To see this, first
12 observe that the $\bar{V}_{1:T}$ generated by any policy satisfies $\bar{V}_{1:T} \in [0, 1]^K$, since $V(s, a) \in [0, 1]$ always. Suppose the
13 contrary that there is a $\tilde{V}_{1:T}$, where $\tilde{V}_{1:T,k} \geq \bar{V}_{1:T,k} \forall k$, and $\tilde{V}_{1:T,1} > \bar{V}_{1:T,1}$. These mean that $0 \leq 1 - \tilde{V}_{1:T,k} \leq$
14 $1 - \bar{V}_{1:T,k} \forall k$, and $0 \leq 1 - \tilde{V}_{1:T,1} < 1 - \bar{V}_{1:T,1}$. Consequently, $g_{\text{MO}}(\tilde{V}_{1:T}) > g_{\text{MO}}(\bar{V}_{1:T}^*)$, contradicting the maximality
15 of $\bar{V}_{1:T}^*$ on g_{MO} . Thus, $\bar{V}_{1:T}^*$ is Pareto-optimal. Altogether, g_{MO} with suitably chosen ρ, U captures Pareto-optimality.
16 Moreover, g_{MO} captures the *State Space Exploration* problem, which goes beyond Pareto-optimality.

17 (2.3): Capturing Pareto-optimality allows us to model many real world problems. Our framework allows any smooth
18 concave g and not just g_{MO} (App. D), which captures other applications such as Maximum Entropy Exploration [23].

19 The design and analysis of GTP: (5.2), (2.1), (2.4). We start by addressing (5.2). For instance (1b), we claim that
20 $\text{opt}(\mathcal{P}_{\mathcal{M}}) = 0$. In addition, the solution x^* , defined as $x^*(s^1, \mathbf{r}1) = x^*(s^2, \mathbf{1}1) = 1/2$ and $x^*(s, a) = 0$ for all
21 other s, a , is optimal to $(\mathcal{P}_{\mathcal{M}})$. Indeed, x^* is feasible to $(\mathcal{P}_{\mathcal{M}})$ (recall $p(s^1|s^1, \mathbf{r}1) = p(s^2|s^2, \mathbf{1}1) = 1$), and that
22 $\sum_{s,a} v(s, a)x^*(s, a) = \binom{0}{1} * x^*(s^2, \mathbf{1}1) + \binom{0}{1} * x^*(s^1, \mathbf{r}1) = \binom{1}{1/2}$, and that $g_{\text{MO}}(\sum_{s,a} v(s, a)x^*(s, a)) = 0$.

23 The bad policy in Line 139, which causes $\bar{V}_{1:T} \approx (1/6, 1/6)^\top$, incurs $\text{Reg}(T) = 0 - (-(1/6 - 1/2)^2) = \Omega(1)$. The
24 $\Omega(1)$ regret is caused by the $\Theta(T)$ *implicit switching cost*, where the agent switches between s^1, s^2 (hence visits s^0) for
25 $\Theta(T)$ times in T time steps. In an MO-OMDP instance, the implicit switching cost occurs when the agent switches
26 form a recurrent class to another, and visits a state that does not contribute to the objective (like s^0) during the switch.

27 (2.1): GTP (see Lines 175-177) consists of the maintenance of distance measure Ψ in Line 13 in Algo 1, and the first
28 criterion $\Psi < Q$ in Line 9 in Algo 1. GTP keeps the implicit switching cost bounded, while balances the contributions
29 by $\{\bar{V}_{1:T,k}\}_{k=1}^K$. As said in Lines 188-193 for Fig 1b, GTP ensures the agent only switches between s^1, s^2 for $O(\sqrt{T})$
30 times in T steps, and $|\bar{V}_{1:T,k} - 0.5| = O(1/\sqrt{T})$ for $k = 1, 2$, thus $\text{Reg}(T) = O(1/\sqrt{T})$. GTP reduces the implicit
31 switching cost from $\Theta(T)$ to $O(\sqrt{T})$, by looping at each s^1, s^2 for $\Theta(\sqrt{t})$ times before switching (cf. Lines 190-191).

32 (2.4): Lemma 4.1 follows from concentration inequalities, which are not our contributions. GTP is new, and its
33 design and analysis are our novel contributions. Compared to UCRL2 for SO-OMDPs, analysing TFW-UCRL2 for
34 MO-OMDPs requires crucial effort on bounding two costs: (i) the implicit switch cost (see Lemma 4.3) due to GTP.
35 (ii) there is a *delay cost* caused by GTP on the gradient updates. In Fig 1b, the delay cost is the $O(1/\sqrt{T})$ error on
36 $|\bar{V}_{1:T,k} - 0.5|$. The delay cost, included by eqn. (11), is discussed in Lines 244-248 and bounded by Proposition 4.2.
37 These switch and delay costs are not present in UCRL2, and their analyses certainly do not follow from the literature.

38 **Reviewer # 8:** In fact, the regret under $Q = \bar{L}/\sqrt{K}$ (where $\bar{L} = L_0 + \max_k |L_k|$) is quite close to the optimal regret
39 by tuning Q . We chose $Q = \bar{L}/\sqrt{K}$ to optimize the dependence on \bar{L}, K in the regret order bound in Theorem
40 3.1. The regret could be improved by tuning Q online, or by optimizing Q in the actual regret bound. ρ, L_0, \dots, L_K
41 parameterize the objective function g_{MO} , which is assumed to be fixed, while Q parameterizes the algo, so we only
42 consider tuning Q . If accepted, we will conduct the suggested empirical comparisons with [26, 28, 34] in their settings.