

---

# What the hell is the title of this paper?

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

The abstract paragraph should be indented  $\frac{1}{2}$  inch (3 picas) on both the left- and right-hand margins. Use 10 point type, with a vertical spacing (leading) of 11 points. The word **Abstract** must be centered, bold, and in point size 12. Two line spaces precede the abstract. The abstract must be limited to one paragraph.

## 1 Introduction

Recent advancements have allowed reinforcement learning (RL) [29] to achieve impressive results in a wide variety of complex tasks, ranging from Atari [22], to the game of Go [28], to the control of sophisticated robotics systems [14, 20, 19]. The main limitation is that RL algorithms still require an enormous amount of experience samples before successfully learning such complicated tasks. One of the most promising solutions is transfer learning, which focuses on reusing past knowledge available to the agent in order to reduce the sample-complexity for learning new tasks. In the typical settings of transfer in RL [31], the agent is assumed to have already solved a set of *source tasks* generated from some unknown distribution. Then, given a *target task* drawn from the same distribution, or a slightly different one, the agent can rely on knowledge from the source tasks to speed up the learning process. This constitutes a significant advantage over plain RL, where the agent learns each new task from scratch independently of previous learning experience. Several algorithms have been proposed in the literature to transfer experience samples [18, 30], policies/options [10, 16], rewards [15], features [4], parameters [9, 13], and so on. We refer the reader to [31] for a thorough survey on transfer in RL.

Under the assumption that tasks follow a certain distribution, an intuitive choice for designing a transfer algorithm is to attempt at characterizing the uncertainty over the target task. Then, an ideal algorithm would leverage prior knowledge from the source tasks to interact with the target task in such a way that this uncertainty is reduced as quickly as possible. This simple intuition makes Bayesian methods appealing approaches for transfer in RL, and many previous works have been proposed in this direction. [33] assume tasks share similarities in their dynamics and rewards and propose a hierarchical Bayesian model for the distribution of these two elements. Similarly, [17] assume tasks are similar in their value functions and design a different hierarchical Bayesian model for transferring such information. More recently, [9], and its extension [13], consider tasks whose dynamics are governed by some hidden parameters, and propose efficient Bayesian models for quickly learning such parameters in new tasks. However, all these algorithms require specific, and sometimes restrictive, assumptions (e.g., on the distributions involved or the function approximators adopted), which might limit their practical applicability. The importance of having transfer algorithms that alleviate the need of strong assumptions and that can be easily adapted to different contexts motivates us to take a more general approach.

Similarly to [17], we assume tasks to share similarities in their value functions and use the given source tasks to learn a distribution over such functions. Then, we use this distribution as a prior for learning the target task and we propose an efficient variational approximation of the corresponding posterior. Leveraging on recent ideas from randomized value functions [23], we design a Thompson Sampling-based algorithm which efficiently explores the target task by repeatedly sampling from the posterior

probabilistic?  
distribu-  
tional?

and acting greedily w.r.t. (with respect to) the sampled value function. We show that our approach is very general, in the sense that it does not require any specific choice of function approximator or prior/posterior distribution models. We propose a finite-sample analysis that theoretically validates our approach, while empirically evaluating its performance on three domains with increasing level of difficulty.

## 2 Preliminaries

We consider a distribution  $\mathcal{D}$  over tasks, where each task is modeled as a discounted Markov Decision Process (MDP). We define an MDP as a tuple  $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, p_0, \gamma \rangle$ , where  $\mathcal{S}$  is the state-space,  $\mathcal{A}$  is a finite set of actions,  $\mathcal{P}(\cdot|s, a)$  is the distribution of the next state  $s'$  given that action  $a$  is taken in state  $s$ ,  $\mathcal{R} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  is the reward function,  $p_0$  is the initial-state distribution, and  $\gamma \in [0, 1)$  is the discount factor. We assume the reward function to be uniformly bounded by a constant  $R_{max} > 0$ . A deterministic policy  $\pi : \mathcal{S} \rightarrow \mathcal{A}$  is a mapping from states to actions. At the beginning of each episode of interaction, the initial state  $s_0$  is drawn from  $p_0$ . Then, the agent takes the action  $a_0 = \pi(s_0)$ , receives a reward  $\mathcal{R}(s_0, a_0)$ , transitions to the next state  $s_1 \sim \mathcal{P}(\cdot|s_0, a_0)$ , and the process is repeated. The goal is to find the policy maximizing the long-term return over a possibly infinite horizon:  $\max_{\pi} J(\pi) \triangleq \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t | \mathcal{M}, \pi]$ . To this end, we define the optimal value function  $Q^*(s, a)$  as the expected return obtained by taking action  $a$  in state  $s$  and following an optimal policy thereafter. Then, an optimal policy  $\pi^*$  is a policy that is greedy with respect to the optimal value function, i.e.,  $\pi^*(s) = \operatorname{argmax}_a Q^*(s, a)$  for all states  $s$ . It can be shown (e.g., [24]) that  $Q^*$  is the unique fixed-point of the optimal Bellman operator  $T$  defined by  $TQ(s, a) = \mathcal{R}(s, a) + \gamma \mathbb{E}_{\mathcal{P}}[\max_{a'} Q(s', a')]$  for any value function  $Q$ . From now on, we adopt the term  $Q$ -function to denote any plausible value function, i.e., any function  $Q : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$  uniformly bounded by  $\frac{R_{max}}{1-\gamma}$ .

Should we provide a more transfer-oriented definition?

We consider a parametric family of  $Q$ -functions,  $\mathcal{Q} = \{Q_{\mathbf{w}} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R} \mid \mathbf{w} \in \mathbb{R}^K\}$ , and we assume each function in  $\mathcal{Q}$  to be uniformly bounded by  $\frac{R_{max}}{1-\gamma}$ . When learning the optimal value function, a quantity of interest is how close a given function  $Q_{\mathbf{w}}$  is to the fixed-point of the Bellman operator. A possible measure is its Bellman error (or Bellman residual), defined by  $B(\mathbf{w}) \triangleq TQ_{\mathbf{w}} - Q_{\mathbf{w}}$ . Notice that  $Q_{\mathbf{w}}$  is optimal if, and only if,  $B(\mathbf{w})(s, a) = 0$  for all  $s, a$ . If we assume the existence of a distribution  $\nu$  over  $\mathcal{S} \times \mathcal{A}$ , an appealing objective is to directly minimize the squared Bellman error of  $Q_{\mathbf{w}}$  under  $\nu$ , denoted by  $\|B(\mathbf{w})\|_{\nu}^2$ . Unfortunately, it is well-known that an unbiased estimator of this quantity requires two independent samples of the next state  $s'$  for each  $s, a$  (e.g., [21]). In practice, the Bellman error is typically replaced by the TD error  $b(\mathbf{w})$ , which approximates the former using a single transition sample,  $b(\mathbf{w}) = r + \gamma \max_{a'} Q_{\mathbf{w}}(s', a') - Q_{\mathbf{w}}(s, a)$ . Finally, given a dataset  $D = \langle s_i, a_i, r_i, s'_i \rangle_{i=1}^N$  of  $N$  samples, the squared TD error is computed as  $\|B(\mathbf{w})\|_D^2 = \frac{1}{N} \sum_{i=1}^N (r_i + \gamma \max_{a'} Q_{\mathbf{w}}(s'_i, a') - Q_{\mathbf{w}}(s_i, a_i))^2 = \frac{1}{N} \sum_{i=1}^N b_i(\mathbf{w})^2$ . Whenever the distinction is clear from the context, with slight abuse of terminology, we refer to the squared Bellman error and squared TD error as Bellman error and TD error, respectively.

Should we define norms before?

## 3 Variational Transfer Learning

In this section, we describe our variational approach to transfer in RL. In Section 3.1, we start by introducing our algorithm from a high-level perspective, in such a way that it can be used for any choice of prior and posterior distributions. Then, in Sections 3.2 and 3.3, we propose practical implementations based on Gaussians and mixtures of Gaussians, respectively. We conclude with some considerations on how to optimize the proposed objective in Sec. 3.4.

### 3.1 Algorithm

We begin with a simple consideration: the distribution  $\mathcal{D}$  over tasks clearly induces a distribution over optimal  $Q$ -functions. Moreover, for any MDP, learning its optimal  $Q$ -function is sufficient for solving the problem. Thus, one can safely replace the distribution over tasks with the distribution over their optimal value functions. In our parametric settings, we reduce the latter to a distribution  $p(\mathbf{w})$  over weights.

Should we just make the assumption that  $Q$ -functions share knowledge in their weights?

---

**Algorithm 1** Variational Transfer

---

**Require:** Target task  $\tau$ , source  $Q$ -function weights  $\mathcal{W}_s$ , batch size  $M$

---

```
1: Estimate prior  $p(\mathbf{w})$  from  $\mathcal{W}_s$ 
2: Initialize variational parameters:  $\xi \leftarrow \operatorname{argmin}_{\xi} KL(q_{\xi} || p)$ 
3: Initialize replay buffer:  $D = \emptyset$ 
4: repeat
5:   Sample initial state:  $s_0 \sim p_0^{(\tau)}$ 
6:   while  $s_h$  is not terminal do
7:     Sample weights:  $\mathbf{w} \sim q_{\xi}(\mathbf{w})$ 
8:     Take action  $a_h = \operatorname{argmax}_a Q_{\mathbf{w}}(s_h, a)$ 
9:     Observe transition  $s_{h+1} \sim \mathcal{P}^{(\tau)}(\cdot | s_h, a_h)$  and collect reward  $r_h = \mathcal{R}^{(\tau)}(s_h, a_h)$ 
10:    Add sample to the replay buffer:  $D \leftarrow D \cup \langle s_h, a_h, r_h, s_{h+1} \rangle$ 
11:    Sample mini-batch  $D' = \langle s_i, a_i, r_i, s'_i \rangle_{i=1}^M$  from  $D$ 
12:    Estimate the gradient  $\nabla_{\xi} \mathcal{L}(\xi)$  using  $D'$ 
13:    Update  $\xi$  in the direction of  $-\nabla_{\xi} \mathcal{L}(\xi)$  using any stochastic optimizer (e.g., ADAM)
14:  end while
15: until forever
```

---

87 Assume, for the moment, that we know the distribution  $p(\mathbf{w})$  and consider a dataset  $D =$   
88  $\{(s_i, a_i, s'_i, r_i) \mid i = 1, 2, \dots, N\}$  of samples from some task  $\tau$  that we want to solve. Then, the  
89 posterior distribution over weights given such dataset can be computed by applying Bayes theorem  
90 as  $p(\mathbf{w}|D) \propto p(D|\mathbf{w})p(\mathbf{w})$ . Unfortunately, this cannot be directly used in practice since we do not  
91 have a model of the likelihood  $p(D|\mathbf{w})$ . In such case, it is very common to make strong assumptions  
92 on the MDPs or the  $Q$ -functions so as to get tractable posteriors. However, in our transfer settings  
93 all distributions involved depend on the family of tasks under consideration, and making such as-  
94 sumptions is likely to limit the applicability to specific problems. Thus, we take a different approach  
95 to derive a more general, but still well-grounded, solution. Recall that our final goal is move all  
96 probability mass over the weights minimizing some empirical loss measure, which in our case is the  
97 TD error  $\|B(\mathbf{w})\|_D^2$ . Then, given a prior  $p(\mathbf{w})$ , we know from PAC-Bayesian theory that the optimal  
98 Gibbs posterior  $q$  takes the form (e.g., [8]):

$$q(\mathbf{w}) = \frac{e^{-\Lambda \|B(\mathbf{w})\|_D^2} p(\mathbf{w})}{\int e^{-\Lambda \|B(\mathbf{w}')\|_D^2} p(d\mathbf{w}')} \quad (1)$$

99 for some parameter  $\Lambda > 0$ . Since  $\Lambda$  is typically chosen to increase with the number of samples  
100  $N$ , in the remaining we set it to  $\lambda^{-1}N$ , for some constant  $\lambda > 0$ . Notice that, whenever the term  
101  $e^{-\Lambda \|B(\mathbf{w})\|_D^2}$  can be interpreted as the actual likelihood of  $D$ ,  $q$  becomes a classic Bayesian posterior.  
102 Although we now have an appealing distribution, the integral at the denominator of (1) is intractable  
103 to compute even for simple  $Q$ -function models. Thus, we propose a variational approximation  $q_{\xi}$  by  
104 considering a simpler family of distributions parameterized by  $\xi \in \Xi$ . Then, our problem reduces to  
105 finding the variational parameters  $\xi$  such that  $q_{\xi}$  minimizes the Kullback-Leibler (KL) divergence  
106 w.r.t. the Gibbs posterior  $q$ . From the theory of variational inference (e.g., [5]), this can be shown to  
107 be equivalent to minimizing the well-known (negative) *evidence lower bound* (ELBO):

$$\min_{\xi \in \Xi} \mathcal{L}(\xi) = \mathbb{E}_{\mathbf{w} \sim q_{\xi}} [\|B(\mathbf{w})\|_D^2] - \frac{\lambda}{N} KL(q_{\xi}(\mathbf{w}) || p(\mathbf{w})) \quad (2)$$

108 Intuitively, the approximate posterior trades-off between placing probability mass over those weights  
109  $\mathbf{w}$  that have low expected TD error (first term), and staying close to the prior distribution (second  
110 term). Assuming that we are able to compute the gradients of (2) w.r.t. the variational parameters  $\xi$ ,  
111 our objective can be easily optimized using any stochastic optimization algorithm.

112 We now highlight our general transfer procedure in Alg. 1, while deferring a description of specific  
113 choices for the distributions involved to the next two sections. Given a set of weights  $\mathcal{W}_s$  from  
114 the source tasks' optimal  $Q$ -functions, we start by estimating the prior distribution (line 1) and we  
115 initialize the variational parameters by minimizing the KL divergence w.r.t. such distribution<sup>1</sup> (line

---

<sup>1</sup>If the prior and approximate posterior were in the same family of distributions we could simply set  $\xi$  to the prior parameters. However, we are not making this assumption at this point.

Maybe say  
how this is  
obtained

116 2). Then, at each time step of interaction, we re-sample the weights from the current approximate  
 117 posterior and act greedily w.r.t. the corresponding  $Q$ -function (lines 7,8). After collecting and storing  
 118 the new experience (lines 9-10), we draw a mini-batch of samples from the replay buffer (line 11), use  
 119 this to estimate the objective function gradient (line 12), and, finally, update the variational parameters  
 120 (line 13).

121 The key property of our approach is the weight resampling at line 7. This resembles the well-known  
 122 Thompson sampling approach adopted in multi-armed bandits [7] and closely relates to the recent  
 123 value function randomization [23]. In some sense, at each time we guess what is the task we are  
 124 trying to solve based on our current belief and we act as if such guess were actually true. This allows  
 125 an efficient adaptive exploration of the target task. Intuitively, during the first steps of interaction,  
 126 the agent is very uncertain about the current task and such uncertainty induces stochasticity in the  
 127 chosen actions, allowing rather informed exploration to take place. Consider, for instance, that actions  
 128 that are bad on average for all tasks are very unlikely to be sampled, while this cannot happen in  
 129 uninformed exploration strategies, like  $\epsilon$ -greedy, before learning takes place. As the learning process  
 130 goes on, the algorithm will quickly figure out which task is being solved, thus moving all probability  
 131 mass over the weights minimizing the TD error. From that point, sampling from the posterior is  
 132 approximately equivalent to deterministically taking such weights, and no more exploration will be  
 133 performed.

134 Finally, notice the generality of the proposed approach: as far as the objective  $\mathcal{L}$  is differentiable  
 135 in the variational parameters  $\xi$ , and its gradients can be efficiently computed, any approximator for  
 136 the  $Q$ -function and any prior/posterior distributions can be adopted. For the latter, we describe two  
 137 practical choices in the next two sections.

### 138 3.2 Gaussian Variational Transfer

139 We now restrict to a specific choice of the prior and posterior families that makes our algorithm  
 140 very efficient and easy to implement. We assume that optimal  $Q$ -functions (or better, their weights)  
 141 follow a multivariate Gaussian distribution. That is, we model the prior as  $p(\mathbf{w}) = \mathcal{N}(\boldsymbol{\mu}_p, \boldsymbol{\Sigma}_p)$   
 142 and we learn its parameters from the set of source weights using maximum likelihood estimation  
 143 (with small regularization to make sure the covariance is positive definite). Then, our variational  
 144 family is the set of all well-defined Gaussian distributions, i.e., the variational parameters are  
 145  $\Xi = \{(\boldsymbol{\mu}, \boldsymbol{\Sigma}) \mid \boldsymbol{\mu} \in \mathbb{R}^K, \boldsymbol{\Sigma} \in \mathbb{R}^{K \times K}, \boldsymbol{\Sigma} \succ 0\}$ . To prevent the covariance from becoming not pos-  
 146 itive definite, we consider its Cholesky decomposition  $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$  and learn the lower-triangular  
 147 Cholesky factor  $\mathbf{L}$  instead. In this case, deriving the gradient of the objective is very simple. Both the  
 148 KL between two multivariate Gaussians and its gradients have a simple closed-form expression. The  
 149 expected log-likelihood, on the other hand, can be easily differentiated by adopting the reparameteri-  
 150 zation trick (e.g., [12, 25]). Although these results are well-known in the literature, we report them in  
 151 App. ?? to have a more self-contained description.

### 152 3.3 Mixture of Gaussian Variational Transfer

153 Although the Gaussian assumption of the previous section is very appealing as it allows for a simple  
 154 and efficient way of computing the variational objective and its gradients, we believe that such  
 155 assumption almost never holds in practice. In fact, even for families of tasks in which the reward and  
 156 transition models are Gaussian, the  $Q$ -values might be far from it. Depending on the family of tasks  
 157 under consideration and, since we are learning a distribution over weights, on the chosen function  
 158 approximator, the prior might have arbitrarily complex shapes. When the information loss due to  
 159 the Gaussian approximation becomes too severe, the algorithm is likely to fail at capturing any  
 160 similarities between the tasks. We now propose a variant to successfully solve this problem, while  
 161 keeping the algorithm simple and efficient enough to be applied in practice.

162 Given the source tasks' weights  $\mathcal{W}_s$ , we model our estimated prior as a mixture of Gaussians with one  
 163 equally weighted isotropic Gaussian centered at each weight:  $p(\mathbf{w}) = \frac{1}{|\mathcal{W}_s|} \sum_{\mathbf{w}_s \in \mathcal{W}_s} \mathcal{N}(\mathbf{w}|\mathbf{w}_s, \sigma_p^2 \mathbf{I})$ .  
 164 This resembles a kernel density estimator [27] and, due to its nonparametric nature, it allows capturing  
 165 arbitrarily complex distribution. Consistently with the prior, we model our approximate posterior as a  
 166 mixture of Gaussians. However, we allow a different number of components (typically much less  
 167 than the prior's) and we adopt full covariances instead of only diagonals, so that our posterior has the  
 168 potential to match complex distributions with less components. Using  $C$  components, our posterior

169 is  $q_{\xi}(\mathbf{w}) = \frac{1}{C} \sum_{i=1}^C \mathcal{N}(\mathbf{w} | \boldsymbol{\mu}_i, \boldsymbol{\Sigma}_i)$ , with variational parameters  $\xi = (\boldsymbol{\mu}_1, \dots, \boldsymbol{\mu}_C, \boldsymbol{\Sigma}_1, \dots, \boldsymbol{\Sigma}_C)$ .  
 170 Once again, we learn Cholesky factors instead of full covariances.

171 Although this new model has the potential to capture much more complex distributions, it poses  
 172 a major complication: the KL divergence between two mixture of Gaussians has no closed-form  
 173 expression. To solve this issue, we rely on an upper bound to such quantity, so that negative ELBO  
 174 still upper bounds the KL between the approximate and true posterior. Among the many upper bounds  
 175 available in the literature, we adopt the one proposed in [11], which we report here for the sake of  
 176 completeness. We refer the reader to the original paper for the proof.

177 **Theorem 1** ([11]). *Let  $p = \sum_i c_i^{(p)} f_i^{(p)}$  and  $q = \sum_j c_j^{(q)} f_j^{(q)}$  be two mixture of Gaussian distri-*  
 178 *butions, where  $f_i^{(p)} = \mathcal{N}(\boldsymbol{\mu}_i^{(p)}, \boldsymbol{\Sigma}_i^{(p)})$  denotes the  $i$ -th component of  $p$ ,  $c_i^{(p)}$  denotes its weight, and*  
 179 *similarly for  $q$ . Introduce two vectors  $\chi^{(1)}$  and  $\chi^{(2)}$  such that  $c_i^{(p)} = \sum_j \chi_{j,i}^{(2)}$  and  $c_j^{(q)} = \sum_i \chi_{i,j}^{(1)}$ .*  
 180 *Then:*

$$KL(p||q) \leq KL(\chi^{(2)}||\chi^{(1)}) + \sum_{i,j} \chi_{j,i}^{(2)} KL(f_i^{(p)}||f_j^{(q)}) \quad (3)$$

181 Our new algorithm replaces the KL with the above-mentioned upper bound. Each time we require  
 182 its value, we have to recompute the parameters  $\chi^{(1)}$  and  $\chi^{(2)}$  that tighten the bound. As shown in  
 183 [11], this can be achieved by a simple fixed-point procedure. Finally, both terms in the objective are  
 184 now linear combinations of functions of the variational parameters of different components, and their  
 185 gradients are easily derived from the ones of the Gaussian case. We report the derivation in App. ??.

### 186 3.4 Optimizing the TD error

187 From Sections 3.2 and 3.3, we know that differentiating the negative ELBO  $\mathcal{L}$  w.r.t.  $\xi$  requires  
 188 differentiating  $\|B(\mathbf{w})\|_D^2$  w.r.t.  $\mathbf{w}$ . Unfortunately, the TD error is well-known to be non-differentiable  
 189 due to the presence of the max operator. This rarely represents a problem since typical value-based  
 190 algorithms are actually semi-gradient methods, i.e., they do not differentiate the targets (see, e.g.,  
 191 Chapter 11 of [29]). However, our transfer settings are rather different than common RL. In fact, our  
 192 algorithm is likely to start from  $Q$ -functions that are very close to an optimum, and the only thing that  
 193 needs to be done is to adapt the weights in a direction of lower error (i.e., higher likelihood) so as to  
 194 quickly converge to the solution of the task that is being solved. Unfortunately, this property cannot  
 195 be guaranteed for most semi-gradient algorithms. Even worse, many online RL algorithms combined  
 196 with complex function approximators (e.g., DQNs) are well-known to be unstable, especially when  
 197 approaching an optimum, and require many tricks and tuning to work well [26, 32]. This is obviously  
 198 an undesirable property in our case, as we only aim at adapting already good solutions. Thus, we  
 199 consider using a residual gradient algorithm [3]. In order to differentiate the targets, we replace  
 200 the optimal Bellman operator with the mellow Bellman operator introduced in [1], which adopts a  
 201 softened version of max called *mellowmax*:

$$\text{mm}_a Q_{\mathbf{w}}(s, a) = \frac{1}{\kappa} \log \frac{1}{|\mathcal{A}|} \sum_a e^{\kappa Q_{\mathbf{w}}(s, a)} \quad (4)$$

202 where  $\kappa$  is a hyperparameter and  $|\mathcal{A}|$  is the number of actions. The mellow Bellman operator, which  
 203 we denote as  $\tilde{T}$ , has several appealing properties that make it suitable for our settings: (i) it converges  
 204 to the maximum as  $\kappa \rightarrow \infty$ , (ii) it has a unique fixed point, and (iii) it is *differentiable*. Denoting  
 205 by  $\tilde{B}(\mathbf{w}) = \tilde{T}Q_{\mathbf{w}} - Q_{\mathbf{w}}$  the Bellman residual w.r.t. the mellow Bellman operator  $\tilde{T}$ , we have that  
 206 the corresponding TD error,  $\|\tilde{B}(\mathbf{w})\|_D^2$ , is now differentiable with respect to  $\mathbf{w}$ . Although residual  
 207 algorithms have guaranteed convergence, they are typically much slower than their semi-gradient  
 208 counterpart. [3] proposed to project the gradient in a direction that achieves higher learning speed,  
 209 while preserving convergence. This can be easily done by including a parameter  $\psi \in [0, 1]$  in the TD  
 210 error gradient such that:

$$\nabla_{\mathbf{w}} \|\tilde{B}(\mathbf{w})\|_D^2 = \frac{2}{N} \sum_{i=1}^N b_i(\mathbf{w}) \left( \gamma \psi \nabla_{\mathbf{w}} \text{mm}_{a'} Q_{\mathbf{w}}(s'_i, a') - \nabla_{\mathbf{w}} Q_{\mathbf{w}}(s_i, a_i) \right) \quad (5)$$

211 where  $b_i(\mathbf{w}) = r_i + \gamma \text{mm}_{a'} Q_{\mathbf{w}}(s'_i, a') - Q_{\mathbf{w}}(s_i, a_i)$ . Notice that  $\psi$  trades-off between the semi-  
 212 gradient ( $\psi = 0$ ) and the full residual gradient ( $\psi = 1$ ). A good criterion for choosing such

parameter is to start with values close to zero (to have faster learning) and move to higher values when approaching the optimum (to guarantee converge).

## 4 Theoretical Analysis

In this section, we theoretically analyze our variational transfer algorithm...

A first important question that we need to answer is whether replacing max with mellow-max in the Bellman operator constitutes a strong approximation or not. It has been proved [ ] that the mellow Bellman operator is a contraction under the  $L_\infty$ -norm and, thus, has a unique fixed-point. However, how such fixed-point differs from the one of the optimal Bellman operator remains an open question. Since mellow-max monotonically converges to max as  $\kappa \rightarrow \infty$ , it would be desirable if the corresponding operator also monotonically converged to the optimal one. We confirm that this property actually holds in the following theorem.

Cite MM

**Theorem 2.** *Let  $V$  be the fixed-point of the optimal Bellman operator  $T$ , and  $Q$  the corresponding action-value function. Define the action-gap function  $g(s)$  as the difference between the value of the best action and the second best action at each state  $s$ . Let  $\tilde{V}$  be the fixed-point of the mellow Bellman operator  $\tilde{T}$  with parameter  $\kappa > 0$  and denote by  $\beta > 0$  the inverse temperature of the induced Boltzmann distribution (as in [1]). Let  $\nu$  be a probability measure over the state-space. Then, for any  $p \geq 1$ :*

$$\|V - \tilde{V}\|_{\nu,p}^p \leq \frac{2R_{max}}{(1-\gamma)^2} \left\| 1 - \frac{1}{1 + |\mathcal{A}| e^{-\beta g}} \right\|_{\nu,p}^p \quad (6)$$

**Theorem 3.** *Let  $Q^*$  be the fixed-point of the optimal Bellman operator  $T$ . Define the action-gap function  $g(s)$  as the difference between the value of the best action and the second best action at each state  $s$ . Let  $\tilde{Q}$  be the fixed-point of the mellow Bellman operator  $\tilde{T}$  with parameter  $\kappa > 0$  and denote by  $\beta_\kappa > 0$  the inverse temperature of the induced Boltzmann distribution (as in [1]). Let  $\nu$  be a probability measure over the state-action space. Then, for any  $p \geq 1$ :*

$$\|Q^* - \tilde{Q}\|_{\nu,p}^p \leq \frac{2\gamma R_{max}}{(1-\gamma)^2} \left\| \frac{1}{1 + \frac{1}{|\mathcal{A}|} e^{\beta_\kappa g}} \right\|_{\nu,p}^p \quad (7)$$

## 5 Related Works

Our approach is mostly related to [17]. Although we both assume the tasks to share similarities in their value functions, [17] consider only linear approximators and adopt a hierarchical Bayesian model of the corresponding weights' distribution, which is assumed Gaussian. On the other hand, our variational approximation allows for more general distribution families and can be combined with non-linear approximators. Furthermore, [17] propose a Dirichlet process model for the case where weights cluster into different classes, which relates to our mixture formulation and proves again the importance of capturing more complicated task distributions. In [33], the authors propose a hierarchical Bayesian model for the distribution over tasks. Differently from our approach and [17], they consider a distribution over transition probabilities and rewards, rather than value functions. In the same spirit of our method, they consider a Thompson sampling-based procedure which, at each iteration, samples a new task from the posterior and solves it. However, [33] consider only finite MDPs, which poses a severe limitation on the algorithm's applicability. On the contrary, our approach can handle high-dimensional tasks. In [9], the authors consider a family of tasks whose dynamics are governed by some hidden parameters and use Gaussian processes (GPs) to model such dynamics across tasks. Recently, [13] extended this approach by replacing GPs with Bayesian neural networks, so as to obtain a more scalable approach.

In the RL community, our approach is related to value function randomization[23], which extends the well-known LSTD [6] by adopting Bayesian linear regression to model the uncertainty over the predicted value function weights, and use that to perform a form of Thompson sampling. Such approach was recently extended by [2], where the approximator is replaced by a Bayesian neural network, leading to an algorithm capable of solving much more complicated problems. Both these algorithms rely on the Gaussian assumption and, since they work in plain RL settings, have no

More comments about HiP-MDPs? Should we discuss more related works?



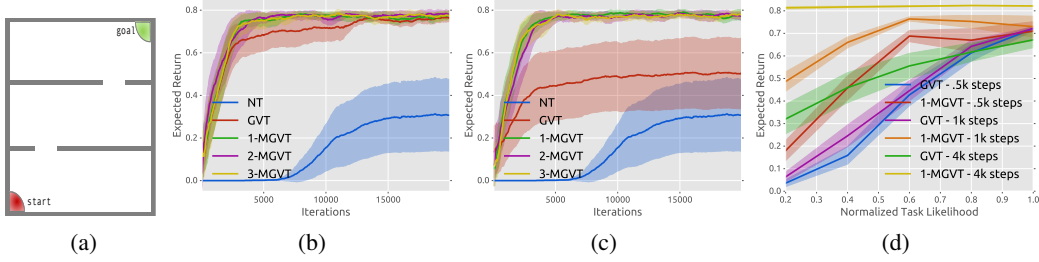


Figure 1: (a) the rooms environment, (b) transfer from 10 source tasks with both doors moving, (c) transfer from 10 source tasks with only one door moving, and (d) transfer performance as a function of how likely the target task is according to the prior.

informative prior available. On the other hand, our variational approximation allows more complex distributions (e.g., mixtures) to be adopted, while knowledge from the source tasks allows us to learn very informative priors.

## 6 Experiments

In this section, we provide an experimental evaluation of our approach in three different domains with increasing level of difficulty. In all experiments, we compare our Gaussian variational transfer algorithm (GVT) and the version using a  $c$ -component mixture of Gaussians ( $c$ -MGVT) to plain no-transfer RL (NT) with  $\epsilon$ -greedy exploration. To the best of our knowledge, no existing transfer algorithm is directly comparable to our approach from an experimental perspective.

Any better motivation?

### 6.1 The Rooms Problem

We consider an agent navigating in the rooms environment depicted in Fig. ?? . The agent starts in the bottom-left corner and must move from one room to another to reach the goal position in the top-right corner. The rooms are connected by small doors whose positions are unknown to the agent. The state-space is modeled as a  $10 \times 10$  continuous grid, while the action-space is the set of 4 movement directions (up, right, down, left). After each action, the agent moves by 1 in the chosen direction and the final position is corrupted by Gaussian noise with 0.2 standard deviation. In case the agent hits a wall, its position remains unchanged. The reward is 1 when reaching the goal (after which the process terminates) and 0 otherwise, while the discount factor is  $\gamma = 0.99$ . We consider a distribution over tasks in which doors have a fixed width of 1, while their positions are sampled uniformly in  $[0.5, 9.5]$ . Since the agent does not know where the doors are located in advance and receives only very sparse feedback, it must efficiently explore the environment to figure out (i) their positions, and (ii) how to reach the goal. While this might be a complicated problem for plain RL, our transfer algorithm should be able to quickly figure out the door positions. In fact, notice that, although different, the optimal  $Q$ -functions for all tasks share some similarities. For example, once the agent has passed the last door before the goal, the  $Q$ -values are exactly the same in all tasks. This does not hold for positions nearby the start state. However, it is clear that there should be a preference over actions up and right, rather than down and left (which are worse in all tasks). Thus, we expect our algorithm to efficiently explore any target task.

In order to prove that our guesses are correct, we generate a set of 50 source tasks for the three-room environment of Fig. ?? by sampling both door positions uniformly, and solve all of them by directly minimizing the TD error as presented in Sec. 3.4. In order to make sure that their solutions are accurate enough, we allow sampling the initial state uniformly in the environment and run until convergence. Then, we use our algorithms to transfer from 10 source tasks sampled from the previously generated set. The average return over the last 50 learning episodes as a function of the number of iterations is shown in Fig. ?? . Each curve is the result of 20 independent runs, each resampling the target and source tasks. 95% confidence intervals are shown. Further details on the parameters adopted in this experiment are given in App. ?? . As expected, the no-transfer (NT) algorithm fails at learning the task in so few iterations due to the limited exploration provided by an  $\epsilon$ -greedy policy. On the other hand, all our algorithms achieve a significant speed-up and are able to

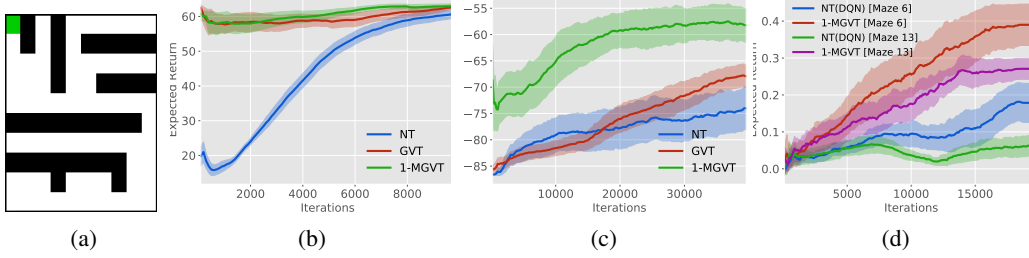


Figure 2: (a) Sample Maze (Maze 6), (b) transfer from 10 sources in Cartpole, (c) transfer from 10 source tasks in Mountain Car, and (d) transfer from 10 sources in Mazes 6 and 13—see App. ??.

297 converge to the optimal performance in few iterations, with GVT being slightly slower. Interestingly,  
 298 we notice that there is no advantage in adopting more than 1 component for the posterior in MGVT.  
 299 This is intuitive since, when the algorithm quickly figures out which task is being solved, it will move  
 300 all components in the same direction.

301 To better understand the differences between GVT and MGVT, we now consider transferring from a  
 302 slightly different distribution than the one from which target tasks are drawn. We generate again 50  
 303 source tasks but this time with the bottom door fixed in the center and the other one moving. Then, we  
 304 repeat the previous experiment, allowing both doors to move when sampling target tasks. The results  
 305 are shown in Fig. ?? . Interestingly, MGVT seems almost unaffected by this change, proving that is  
 306 has sufficient representation power to generalize to slightly different task distributions. The same  
 307 does not hold for GVT, which now is not able to solve many of the sampled target tasks, thus the very  
 308 high variance. This proves again that assuming Gaussian distributions can pose severe limitations in  
 309 our transfer settings.

310 Finally, we analyze the transfer performance as a function of how likely the target task is according to  
 311 the prior. We consider a two-room version of the environment of Fig. ?? . Differently from before, we  
 312 generate tasks by sampling the door position from a Gaussian with mean 5, and standard deviation  
 313 1.8, so that tasks where the door is near the sides are very unlikely. Fig. ?? shows the performance  
 314 reached by GVT and MGVT with 1 component at fixed iterations as a function of how likely the  
 315 target task is according to such distribution. As expected GVT achieves poor performance on very  
 316 unlikely tasks, even after many iterations. In fact, estimating a single Gaussian distribution definitely  
 317 entails some information loss, especially about the unlikely tasks. On the other hand, MGVT keeps  
 318 such information and, consequently, performs much better. Perhaps not surprisingly, MGVT reaches  
 319 the optimal performance in  $4k$  iterations no matter what task is being solved.

## 320 6.2 Classic Control

321 To evaluate the performance of our transfer algorithms in more common benchmarks, we chose to use  
 322 the well-known classic control environments: Cartpole and Mountain Car as defined in [29]. In the  
 323 case of Cartpole we generated 20 source tasks by uniformly sampling the cart mass, the pole mass and  
 324 length. Instead, for Mountain Car, we generated the 20 sources by uniformly sampling the base speed  
 325 of the car. We kept a discount factor of  $\gamma = 0.99$  for both environments. Moreover, to train the source  
 326 tasks we chose to use a Multilayer Perceptron (MLP) as approximator. We trained all source tasks  
 327 until convergence and to transfer, in both settings, we used GVT and 1-MGVT choosing randomly 10  
 328 of the trained sources. Analogously to Sec. 6.1, in Fig. ?? and ?? we show the expected return—as  
 329 the average of the last 50 episodes of learning—of transfer for 20 independent runs randomizing the  
 330 sources and target tasks. In Fig. ??, that corresponds to Cartpole, it is easily noticeable a near-optimal  
 331 jump-start, in both GVT and 1-MGVT, that during the following iterations converges to the optimal  
 332 performance. Furthermore, transfer in Mountain Car as shown in Fig. ??, a more complicated task  
 333 than Cartpole, causes GVT to struggle to solve the target task but 1-MGVT reaches convergence,  
 334 within the same number of iterations, to the optimal performances.

## 335 6.3 Maze Navigation

336 For this setting we propose a maze environment, such as those in Fig. ??, for a robotic agent to  
 337 navigate. We created a set of 20 mazes—see App ??—of size  $10 \times 10$  with continuous space within



which we ensured four groups: each group with the goal position in one of the corners of the maze. In order to simulate a robotic agent, we set an action space that allows to control the linear and angular motion by a fixed distance of 0.5 (*move forward*) and a fixed angle of  $\pm 22.5^\circ$  (*rotate counter- and clockwise*). Moreover, we enhanced the state space to include, in addition to the absolute position and orientation, a set of distance measurements to the nearest obstacles from the agent’s perspective with maximum range of 2 and a field of vision of 9 equally-spaced directions in  $[-90, 90^\circ]$  relative to its orientation. In this way, we provide the agent with information that a robot would normally collect from LIDARs, sonars and other sensors alike. Finally, in the same field of perception, we provide a boolean vector indicating whether the goal is within the range of observation. Regarding rewards, the agent gets 1 whenever it enters the goal area and 0 in any other situation and the MDP has a discount factor of  $\gamma = 0.99$ .

The main motivation for this set-up is that a robot is able to exploit locally-sensed information in addition to its estimation of the actual position—as opposed to the Rooms problem presented above—and it is such information that the robotic agent could learn to use and, also, exploit in different maze configurations leveraging its previous experience. Moreover, simple variations in the maze structure could, potentially, produce big differences in the value function with respect to the ones seen in previous tasks which makes the setting suitable to assess how robust our algorithm is to negative transfer.

The mazes were solved using Deep Q Networks (DQN) [ ] with an MLP with hidden layers  $32 \times 32$  and, in order to perform the knowledge transfer, the target maze would be excluded from the source tasks given to the algorithm. Finally, the 1-MGVT algorithm was used to solve the target tasks—Mazes 6 and 13. As before, 20 independent runs with randomized sources to transfer are averaged and shown in Fig. ??, in this case 5 source tasks were used for each trial. While these mazes pose a challenge for our algorithm, we can see that 1-MGVT is capable to still capture useful information from the sources and adapt to the target mazes which is demonstrated by the increased slope in the learning curve.

cite DQN?

actually, decide how to improve these results

## 7 Conclusion

## References

- [1] Kavosh Asadi and Michael L Littman. An alternative softmax operator for reinforcement learning. In *International Conference on Machine Learning*, pages 243–252, 2017.
- [2] Kamyar Azizzadenesheli, Emma Brunskill, and Animashree Anandkumar. Efficient exploration through bayesian deep q-networks. *arXiv preprint arXiv:1802.04412*, 2018.
- [3] Leemon Baird. Residual algorithms: Reinforcement learning with function approximation. In *Machine Learning Proceedings 1995*, pages 30–37. Elsevier, 1995.
- [4] André Barreto, Will Dabney, Rémi Munos, Jonathan J Hunt, Tom Schaul, Hado P van Hasselt, and David Silver. Successor features for transfer in reinforcement learning. In *Advances in neural information processing systems*, pages 4055–4065, 2017.
- [5] David M Blei, Alp Kucukelbir, and Jon D McAuliffe. Variational inference: A review for statisticians. *Journal of the American Statistical Association*, 112(518):859–877, 2017.
- [6] Justin A Boyan. Least-squares temporal difference learning.
- [7] Sébastien Bubeck, Nicolo Cesa-Bianchi, et al. Regret analysis of stochastic and nonstochastic multi-armed bandit problems. *Foundations and Trends® in Machine Learning*, 5(1):1–122, 2012.
- [8] Olivier Catoni. Pac-bayesian supervised classification: the thermodynamics of statistical learning. *arXiv preprint arXiv:0712.0248*, 2007.
- [9] Finale Doshi-Velez and George Konidaris. Hidden parameter markov decision processes: A semiparametric regression approach for discovering latent task parametrizations. In *IJCAI: proceedings of the conference*, volume 2016, page 1432. NIH Public Access, 2016.
- [10] Fernando Fernández and Manuela Veloso. Probabilistic policy reuse in a reinforcement learning agent. In *Proceedings of the fifth international joint conference on Autonomous agents and multiagent systems*, pages 720–727. ACM, 2006.

- [11] John R Hershey and Peder A Olsen. Approximating the kullback leibler divergence between gaussian mixture models. In *Acoustics, Speech and Signal Processing, 2007. ICASSP 2007. IEEE International Conference on*, volume 4, pages IV–317. IEEE, 2007.
- [12] Matthew D Hoffman, David M Blei, Chong Wang, and John Paisley. Stochastic variational inference. *The Journal of Machine Learning Research*, 14(1):1303–1347, 2013.
- [13] Taylor W Killian, Samuel Daulton, George Konidaris, and Finale Doshi-Velez. Robust and efficient transfer learning with hidden parameter markov decision processes. In *Advances in Neural Information Processing Systems*, pages 6250–6261, 2017.
- [14] Jens Kober and Jan R Peters. Policy search for motor primitives in robotics. In *Advances in neural information processing systems*, pages 849–856, 2009.
- [15] George Konidaris and Andrew Barto. Autonomous shaping: Knowledge transfer in reinforcement learning. In *Proceedings of the 23rd international conference on Machine learning*, pages 489–496. ACM, 2006.
- [16] George Konidaris and Andrew G Barto. Building portable options: Skill transfer in reinforcement learning.
- [17] Alessandro Lazaric and Mohammad Ghavamzadeh. Bayesian multi-task reinforcement learning. In *ICML-27th International Conference on Machine Learning*, pages 599–606. Omnipress, 2010.
- [18] Alessandro Lazaric, Marcello Restelli, and Andrea Bonarini. Transfer of samples in batch reinforcement learning. In *Proceedings of the 25th international conference on Machine learning*, pages 544–551. ACM, 2008.
- [19] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *The Journal of Machine Learning Research*, 17(1):1334–1373, 2016.
- [20] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. *arXiv preprint arXiv:1509.02971*, 2015.
- [21] Odalric-Ambrym Maillard, Rémi Munos, Alessandro Lazaric, and Mohammad Ghavamzadeh. Finite-sample analysis of bellman residual minimization. In *Proceedings of 2nd Asian Conference on Machine Learning*, pages 299–314, 2010.
- [22] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529, 2015.
- [23] Ian Osband, Benjamin Van Roy, and Zheng Wen. Generalization and exploration via randomized value functions. *arXiv preprint arXiv:1402.0635*, 2014.
- [24] Martin L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA, 1994.
- [25] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation and approximate inference in deep generative models. *arXiv preprint arXiv:1401.4082*, 2014.
- [26] Tom Schaul, John Quan, Ioannis Antonoglou, and David Silver. Prioritized experience replay. *arXiv preprint arXiv:1511.05952*, 2015.
- [27] David W Scott. *Multivariate density estimation: theory, practice, and visualization*. John Wiley & Sons, 2015.
- [28] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484–489, 2016.
- [29] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*, volume 1. MIT press Cambridge, 1998.
- [30] Matthew E Taylor, Nicholas K Jong, and Peter Stone. Transferring instances for model-based reinforcement learning. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 488–505. Springer, 2008.
- [31] Matthew E Taylor and Peter Stone. Transfer learning for reinforcement learning domains: A survey. *Journal of Machine Learning Research*, 10(Jul):1633–1685, 2009.

- 437 [32] Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning.  
438 2016.
- 439 [33] Aaron Wilson, Alan Fern, Soumya Ray, and Prasad Tadepalli. Multi-task reinforcement learning: a  
440 hierarchical bayesian approach. In *Proceedings of the 24th international conference on Machine learning*,  
441 pages 1015–1022. ACM, 2007.

## 442 A Proofs

443 **Theorem 2.** Let  $V$  be the fixed-point of the optimal Bellman operator  $T$ , and  $Q$  the corresponding  
 444 action-value function. Define the action-gap function  $g(s)$  as the difference between the value of  
 445 the best action and the second best action at each state  $s$ . Let  $\tilde{V}$  be the fixed-point of the mellow  
 446 Bellman operator  $\tilde{T}$  with parameter  $\kappa > 0$  and denote by  $\beta > 0$  the inverse temperature of the  
 447 induced Boltzmann distribution (as in [1]). Let  $\nu$  be a probability measure over the state-space. Then,  
 448 for any  $p \geq 1$ :

$$\|V - \tilde{V}\|_{\nu,p}^p \leq \frac{2R_{max}}{(1-\gamma)^2} \left\| 1 - \frac{1}{1 + |\mathcal{A}| e^{-\beta g}} \right\|_{\nu,p}^p \quad (6)$$

449 *Proof.* We begin by noticing that:

$$\begin{aligned} \|V - \tilde{V}\|_{\nu,p}^p &= \|TV - \tilde{T}\tilde{V}\|_{\nu,p}^p \\ &= \|TV - \tilde{T}V + \tilde{T}V - \tilde{T}\tilde{V}\|_{\nu,p}^p \\ &\leq \|TV - \tilde{T}V\|_{\nu,p}^p + \|\tilde{T}V - \tilde{T}\tilde{V}\|_{\nu,p}^p \\ &\leq \|TV - \tilde{T}V\|_{\nu,p}^p + \gamma \|V - \tilde{V}\|_{\nu,p}^p \end{aligned}$$

450 where the first inequality follows from Minkowsky's inequality and the second one from the contrac-  
 451 tion property of the mellow Bellman operator. This implies that:

$$\|V - \tilde{V}\|_{\nu,p}^p \leq \frac{1}{1-\gamma} \|TV - \tilde{T}V\|_{\nu,p}^p \quad (8)$$

452 Let us bound the norm on the right-hand side separately. In order to do that, we will bound the  
 453 function  $|TV(s) - \tilde{T}V(s)|$  point-wisely for any state  $s$ . By applying the definition of the optimal  
 454 and mellow Bellman operators, we obtain:

$$\begin{aligned} |TV(s) - \tilde{T}V(s)| &= \left| \max_a \{R(s, a) + \gamma \mathbb{E}[V(s')]\} - \min_a \{R(s, a) + \gamma \mathbb{E}[V(s')]\} \right| \\ &= \left| \max_a Q(s, a) - \min_a Q(s, a) \right| \end{aligned}$$

455 Recall that applying the mellow-max is equivalent to computing an expectation under a Boltzmann  
 456 distribution with inverse temperature  $\beta$  induced by  $\kappa$  []. Thus, we can write:

Cite MM

$$\begin{aligned} \left| \max_a Q(s, a) - \min_a Q(s, a) \right| &= \left| \sum_a \pi^*(a|s) Q(s, a) - \sum_a \pi_\beta(a|s) Q(s, a) \right| \\ &= \left| \sum_a Q(s, a) (\pi^*(a|s) - \pi_\beta(a|s)) \right| \\ &\leq \sum_a |Q(s, a)| |\pi^*(a|s) - \pi_\beta(a|s)| \\ &\leq \frac{R_{max}}{1-\gamma} \sum_a |\pi^*(a|s) - \pi_\beta(a|s)| \end{aligned} \quad (9)$$

457 where  $\pi^*$  is the optimal (deterministic) policy w.r.t.  $Q$  and  $\pi_\beta$  is the Boltzmann distribution induced  
 458 by  $Q$  with inverse temperature  $\beta$ :

$$\pi_\beta(a|s) = \frac{e^{\beta Q(s,a)}}{\sum_{a'} e^{\beta Q(s,a')}}$$

459 Denote by  $a_1(s)$  the optimal action for state  $s$  under  $Q$ . We can then write:

$$\begin{aligned}
\sum_a |\pi^*(a|s) - \pi_\beta(a|s)| &= |\pi^*(a_1(s)|s) - \pi_\beta(a_1(s)|s)| + \sum_{a \neq a_1(s)} |\pi^*(a|s) - \pi_\beta(a|s)| \\
&= |1 - \pi_\beta(a_1(s)|s)| + \sum_{a \neq a_1(s)} |\pi_\beta(a|s)| \\
&= 2 |1 - \pi_\beta(a_1(s)|s)|
\end{aligned} \tag{10}$$

460 Finally, let us bound this last term:

$$\begin{aligned}
|1 - \pi_\beta(a_1(s)|s)| &= \left| 1 - \frac{e^{\beta Q(s, a_1(s))}}{\sum_{a'} e^{\beta Q(s, a')}} \right| \\
&= \left| 1 - \frac{e^{\beta(Q(s, a_1(s)) - Q(s, a_2(s)))}}{\sum_{a'} e^{\beta(Q(s, a') - Q(s, a_2(s)))}} \right| \\
&= \left| 1 - \frac{e^{\beta g(s)}}{\sum_{a'} e^{\beta(Q(s, a') - Q(s, a_2(s)))}} \right| \\
&= \left| 1 - \frac{e^{\beta g(s)}}{e^{\beta g(s)} + \sum_{a' \neq a_1(s)} e^{\beta(Q(s, a') - Q(s, a_2(s)))}} \right| \\
&\leq \left| 1 - \frac{e^{\beta g(s)}}{e^{\beta g(s)} + |\mathcal{A}|} \right| \\
&= \left| 1 - \frac{1}{1 + |\mathcal{A}| e^{-\beta g(s)}} \right|
\end{aligned} \tag{11}$$

461 Combining Eq. (14), (15), and (16), we obtain:

$$\left| \max_a Q(s, a) - \min_a Q(s, a) \right| \leq \frac{2R_{max}}{1 - \gamma} \left| 1 - \frac{1}{1 + |\mathcal{A}| e^{-\beta g(s)}} \right|$$

462 Taking the norm and plugging this into Eq. (12) concludes the proof.  $\square$

463 **Theorem 3.** Let  $Q^*$  be the fixed-point of the optimal Bellman operator  $T$ . Define the action-gap  
464 function  $g(s)$  as the difference between the value of the best action and the second best action at  
465 each state  $s$ . Let  $\tilde{Q}$  be the fixed-point of the mellow Bellman operator  $\tilde{T}$  with parameter  $\kappa > 0$  and  
466 denote by  $\beta_\kappa > 0$  the inverse temperature of the induced Boltzmann distribution (as in [1]). Let  $\nu$  be  
467 a probability measure over the state-action space. Then, for any  $p \geq 1$ :

$$\left\| Q^* - \tilde{Q} \right\|_{\nu, p}^p \leq \frac{2\gamma R_{max}}{(1 - \gamma)^2} \left\| \frac{1}{1 + \frac{1}{|\mathcal{A}|} e^{\beta_\kappa g}} \right\|_{\nu, p}^p \tag{7}$$

468 *Proof.* We begin by noticing that:

$$\begin{aligned}
\left\| Q^* - \tilde{Q} \right\|_{\nu, p}^p &= \left\| TQ^* - \tilde{T}\tilde{Q} \right\|_{\nu, p}^p \\
&= \left\| TQ^* - \tilde{T}Q^* + \tilde{T}Q^* - \tilde{T}\tilde{Q} \right\|_{\nu, p}^p \\
&\leq \left\| TQ^* - \tilde{T}Q^* \right\|_{\nu, p}^p + \left\| \tilde{T}Q^* - \tilde{T}\tilde{Q} \right\|_{\nu, p}^p \\
&\leq \left\| TQ^* - \tilde{T}Q^* \right\|_{\nu, p}^p + \gamma \left\| Q^* - \tilde{Q} \right\|_{\nu, p}^p
\end{aligned}$$

469 where the first inequality follows from Minkowsky's inequality and the second one from the contrac-  
470 tion property of the mellow Bellman operator. This implies that:

$$\left\| Q^* - \tilde{Q} \right\|_{\nu, p}^p \leq \frac{1}{1 - \gamma} \left\| TQ^* - \tilde{T}Q^* \right\|_{\nu, p}^p \tag{12}$$

471 Let us bound the norm on the right-hand side separately. In order to do that, we will bound the  
 472 function  $|TQ^*(s, a) - \tilde{T}Q^*(s, a)|$  point-wisely for any state  $s, a$ . By applying the definition of the  
 473 optimal and mellow Bellman operators, we obtain:

$$\begin{aligned} |TQ^*(s, a) - \tilde{T}Q^*(s, a)| &= |R(s, a) + \gamma \mathbb{E} \left[ \max_{a'} Q^*(s', a') \right] - R(s, a) - \gamma \mathbb{E} \left[ \text{mm}_{a'} Q^*(s', a') \right]| \\ &= \gamma \left| \mathbb{E} \left[ \max_{a'} Q^*(s', a') \right] - \mathbb{E} \left[ \text{mm}_{a'} Q^*(s', a') \right] \right| \\ &\leq \gamma \mathbb{E} \left[ \left| \max_{a'} Q^*(s', a') - \text{mm}_{a'} Q^*(s', a') \right| \right] \end{aligned} \quad (13)$$

474 Thus, bounding this quantity reduces to bounding  $|\max_a Q^*(s, a) - \text{mm}_a Q^*(s, a)|$  point-wisely  
 475 for any  $s$ . Recall that applying the mellow-max is equivalent to computing an expectation under a  
 476 Boltzmann distribution with inverse temperature  $\beta_\kappa$  induced by  $\kappa$  []. Thus, we can write:

Cite MM

$$\begin{aligned} \left| \max_a Q^*(s, a) - \text{mm}_a Q^*(s, a) \right| &= \left| \sum_a \pi^*(a|s) Q^*(s, a) - \sum_a \pi_{\beta_\kappa}(a|s) Q^*(s, a) \right| \\ &= \left| \sum_a Q^*(s, a) (\pi^*(a|s) - \pi_{\beta_\kappa}(a|s)) \right| \\ &\leq \sum_a |Q^*(s, a)| |\pi^*(a|s) - \pi_{\beta_\kappa}(a|s)| \\ &\leq \frac{R_{max}}{1 - \gamma} \sum_a |\pi^*(a|s) - \pi_{\beta_\kappa}(a|s)| \end{aligned} \quad (14)$$

477 where  $\pi^*$  is the optimal (deterministic) policy w.r.t.  $Q^*$  and  $\pi_{\beta_\kappa}$  is the Boltzmann distribution induced  
 478 by  $Q^*$  with inverse temperature  $\beta_\kappa$ :

$$\pi_{\beta_\kappa}(a|s) = \frac{e^{\beta_\kappa Q^*(s, a)}}{\sum_{a'} e^{\beta_\kappa Q^*(s, a')}}.$$

479 Denote by  $a_1(s)$  the optimal action for state  $s$  under  $Q^*$ . We can then write:

$$\begin{aligned} \sum_a |\pi^*(a|s) - \pi_{\beta_\kappa}(a|s)| &= |\pi^*(a_1(s)|s) - \pi_{\beta_\kappa}(a_1(s)|s)| + \sum_{a \neq a_1(s)} |\pi^*(a|s) - \pi_{\beta_\kappa}(a|s)| \\ &= |1 - \pi_{\beta_\kappa}(a_1(s)|s)| + \sum_{a \neq a_1(s)} |\pi_{\beta_\kappa}(a|s)| \\ &= 2 |1 - \pi_{\beta_\kappa}(a_1(s)|s)| \end{aligned} \quad (15)$$

480 Finally, let us bound this last term:

$$\begin{aligned} |1 - \pi_{\beta_\kappa}(a_1(s)|s)| &= \left| 1 - \frac{e^{\beta_\kappa Q^*(s, a_1(s))}}{\sum_{a'} e^{\beta_\kappa Q^*(s, a')}} \right| \\ &= \left| 1 - \frac{e^{\beta_\kappa (Q^*(s, a_1(s)) - Q^*(s, a_2(s)))}}{\sum_{a'} e^{\beta_\kappa (Q^*(s, a') - Q^*(s, a_2(s)))}} \right| \\ &= \left| 1 - \frac{e^{\beta_\kappa g(s)}}{\sum_{a'} e^{\beta_\kappa (Q^*(s, a') - Q^*(s, a_2(s)))}} \right| \\ &= \left| 1 - \frac{e^{\beta_\kappa g(s)}}{e^{\beta_\kappa g(s)} + \sum_{a' \neq a_1(s)} e^{\beta_\kappa (Q^*(s, a') - Q^*(s, a_2(s)))}} \right| \\ &\leq \left| 1 - \frac{e^{\beta_\kappa g(s)}}{e^{\beta_\kappa g(s)} + |\mathcal{A}|} \right| \\ &= \left| \frac{1}{1 + \frac{1}{|\mathcal{A}|} e^{\beta_\kappa g(s)}} \right| \end{aligned} \quad (16)$$



Combining Eq. (14), (15), and (16), we obtain:

$$\left| \max_a Q(s, a) - \min_a Q(s, a) \right| \leq \frac{2R_{max}}{1-\gamma} \left| \frac{1}{1 + \frac{1}{|\mathcal{A}|} e^{\beta_{\kappa} g(s)}} \right|$$

Finally, using Eq. (13) we get:

$$\left| TQ^*(s, a) - \tilde{T}Q^*(s, a) \right| \leq \frac{2\gamma R_{max}}{1-\gamma} \left| \frac{1}{1 + \frac{1}{|\mathcal{A}|} e^{\beta_{\kappa} g(s)}} \right|$$

Taking the norm and plugging this into Eq. (12) concludes the proof.  $\square$

**Lemma 1.** Let  $p$  and  $\nu$  denote probability measures over  $Q$ -functions and state-action pairs, respectively. Assume  $Q^*$  is the unique fixed-point of the optimal Bellman operator  $T$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$  over the choice of a  $Q$ -function  $Q$ , the following holds:

$$\|Q - Q^*\|_{\nu}^2 \leq \frac{\mathbb{E}_p [\|B(Q)\|_{\nu}^2]}{(1-\gamma)\delta} \quad (17)$$

*Proof.* First notice that:

$$\begin{aligned} \|Q - Q^*\| &= \|Q + TQ - TQ - TQ^*\| \\ &\leq \|Q - TQ\| + \|TQ - TQ^*\| \\ &\leq \|Q - TQ\| + \gamma \|Q - Q^*\| \\ &= \|B(Q)\| + \gamma \|Q - Q^*\| \end{aligned}$$

which implies that:

$$\|Q - Q^*\| \leq \frac{1}{1-\gamma} \|B(Q)\|$$

Then we can write:

$$P(\|Q - Q^*\| > \epsilon) \leq P(\|B(Q)\| > \epsilon(1-\gamma)) \leq \frac{\mathbb{E}_p [\|B(Q)\|_{\nu}^2]}{(1-\gamma)\epsilon}$$

Settings the right-hand side equal to  $\delta$  and solving for  $\epsilon$  concludes the proof.  $\square$

**Corollary 1.** Let  $p$  and  $\nu$  denote probability measures over  $Q$ -functions and state-action pairs, respectively. Assume  $\tilde{Q}$  is the unique fixed-point of the mellow Bellman operator  $\tilde{T}$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$  over the choice of a  $Q$ -function  $Q$ , the following holds:

$$\|Q - \tilde{Q}\|_{\nu}^2 \leq \frac{\mathbb{E}_p [\|\tilde{B}(Q)\|_{\nu}^2]}{(1-\gamma)\delta} \quad (18)$$

**Lemma 2.** Assume  $Q$ -functions belong to a parametric space of functions bounded by  $\frac{R_{max}}{1-\gamma}$ . Let  $p$  and  $q$  be arbitrary distributions over the parameter space  $\mathcal{W}$ , and  $\nu$  be a probability measure over  $\mathcal{S} \times \mathcal{A}$ . Consider a dataset  $D$  of  $N$  samples and define  $v(\mathbf{w}) \triangleq \mathbb{E}_{\nu} [\text{Var}_{\mathcal{P}} [b(\mathbf{w})]]$ . Then, for any  $\delta > 0$ , with probability at least  $1 - \delta$ , the following two inequalities hold simultaneously:

$$\mathbb{E}_q [\|B(\mathbf{w})\|_{\nu}^2] \leq \mathbb{E}_q [\|B(\mathbf{w})\|_D^2] - \mathbb{E}_q [v(\mathbf{w})] + \frac{\lambda}{N} KL(q||p) + 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \quad (19)$$

498

$$\mathbb{E}_q [\|B(\mathbf{w})\|_D^2] \leq \mathbb{E}_q [\|B(\mathbf{w})\|_{\nu}^2] + \mathbb{E}_q [v(\mathbf{w})] + \frac{\lambda}{N} KL(q||p) + 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \quad (20)$$

499 *Proof.* From Hoeffding's inequality we have:

$$P \left( \left| \mathbb{E}_{\nu, \mathcal{P}} \left[ \|B(\mathbf{w})\|_D^2 \right] - \|B(\mathbf{w})\|_D^2 \right| > \epsilon \right) \leq 2 \exp \left( - \frac{2N\epsilon^2}{\left( 2 \frac{R_{max}}{1-\gamma} \right)^4} \right)$$

500 which implies that, for any  $\delta > 0$ , with probability at least  $1 - \delta$ :

$$\left| \mathbb{E}_{\nu, \mathcal{P}} \left[ \|B(\mathbf{w})\|_D^2 \right] - \|B(\mathbf{w})\|_D^2 \right| \leq 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}}$$

501 Under independence assumptions, the expected TD error can be re-written as:

$$\begin{aligned} \mathbb{E}_{\nu, \mathcal{P}} \left[ \|B(\mathbf{w})\|_D^2 \right] &= \mathbb{E}_{\nu, \mathcal{P}} \left[ \frac{1}{N} \sum_{i=1}^N (r_i + \gamma \min_{a'} Q_{\mathbf{w}}(s'_i, a') - Q_{\mathbf{w}}(s_i, a_i))^2 \right] \\ &= \mathbb{E}_{\nu, \mathcal{P}} \left[ (R(s, a) + \gamma \min_{a'} Q_{\mathbf{w}}(s', a') - Q_{\mathbf{w}}(s, a))^2 \right] \\ &= \mathbb{E}_{\nu} [\mathbb{E}_{\mathcal{P}} [b(\mathbf{w})^2]] \\ &= \mathbb{E}_{\nu} [Var_{\mathcal{P}} [b(\mathbf{w})] + \mathbb{E}_{\mathcal{P}} [b(\mathbf{w})]^2] \\ &= v(\mathbf{w}) + \|B(\mathbf{w})\|_{\nu}^2 \end{aligned}$$

502 where  $v(\mathbf{w}) \triangleq \mathbb{E}_{\nu} [Var_{\mathcal{P}} [b(\mathbf{w})]]$ . Thus:

$$\left| \|B(\mathbf{w})\|_{\nu}^2 + v(\mathbf{w}) - \|B(\mathbf{w})\|_D^2 \right| \leq 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \quad (21)$$

503 From the change of measure inequality [], we have that, for any measurable function  $f(\mathbf{w})$  and any  
504 two probability measures  $p$  and  $q$ :

$$\log \mathbb{E}_p \left[ e^{f(\mathbf{w})} \right] \geq \mathbb{E}_q [f(\mathbf{w})] - KL(q||p)$$

Find a reference for this

505 Thus, multiplying both sides of (21) by  $\lambda^{-1}N$  and applying the change of measure inequality with  
506  $f(\mathbf{w}) = \lambda^{-1}N \left| \|B(\mathbf{w})\|_{\nu}^2 + v(\mathbf{w}) - \|B(\mathbf{w})\|_D^2 \right|$ , we obtain:

$$\mathbb{E}_q [f(\mathbf{w})] - KL(q||p) \leq \log \mathbb{E}_p \left[ e^{f(\mathbf{w})} \right] \leq 4 \frac{R_{max}^2 \lambda^{-1}N}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}}$$

507 where the second inequality holds since the right-hand side of (21) does not depend on  $\mathbf{w}$ . Finally,  
508 we can explicitly write:

$$\mathbb{E}_q \left[ \left| \|B(\mathbf{w})\|_{\nu}^2 + v(\mathbf{w}) - \|B(\mathbf{w})\|_D^2 \right| \right] \leq \frac{\lambda}{N} KL(q||p) + 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}}$$

509 from which the lemma follows straightforwardly.  $\square$

510 **Lemma 3.** Let  $p$  be a prior distribution over the parameter space  $\mathcal{W}$ , and  $\nu$  be a probability measure  
511 over  $\mathcal{S} \times \mathcal{A}$ . Assume  $\hat{\xi}$  is the minimizer of  $ELBO(\xi) = \mathbb{E}_{q_{\xi}} \left[ \|B(\mathbf{w})\|_D^2 \right] + \frac{\lambda}{N} KL(q_{\xi}||p)$  for a  
512 dataset  $D$  of  $N$  samples. Define  $v(\mathbf{w}) \triangleq \mathbb{E}_{\nu} [Var_{\mathcal{P}} [b(\mathbf{w})]]$ . Then, for any  $\delta > 0$ , with probability at  
513 least  $1 - \delta$ :

$$\mathbb{E}_{q_{\hat{\xi}}} \left[ \|B(\mathbf{w})\|_{\nu}^2 \right] \leq \inf_{\xi \in \Xi} \left\{ \mathbb{E}_{q_{\xi}} \left[ \|B(\mathbf{w})\|_{\nu}^2 \right] + \mathbb{E}_{q_{\xi}} [v(\mathbf{w})] + 2 \frac{\lambda}{N} KL(q_{\xi}||p) \right\} + 2 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{N}}$$

514 *Proof.* Let us use Lemma 2 for the specific choice  $q = q_{\hat{\xi}}$ . From Eq. (19), we have:

$$\begin{aligned} \mathbb{E}_{q_{\hat{\xi}}} [\|B(\mathbf{w})\|_{\nu}^2] &\leq \mathbb{E}_{q_{\hat{\xi}}} [\|B(\mathbf{w})\|_D^2] - \mathbb{E}_{q_{\hat{\xi}}} [v(\mathbf{w})] + \frac{\lambda}{N} KL(q_{\hat{\xi}} \| p) + 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \\ &\leq \mathbb{E}_{q_{\hat{\xi}}} [\|B(\mathbf{w})\|_D^2] + \frac{\lambda}{N} KL(q_{\hat{\xi}} \| p) + 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \\ &= \inf_{\xi \in \Xi} \left\{ \mathbb{E}_{q_{\xi}} [\|B(\mathbf{w})\|_D^2] + \frac{\lambda}{N} KL(q_{\xi} \| p) \right\} + 4 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{2N}} \end{aligned}$$

515 where the second inequality holds since  $v(\mathbf{w}) > 0$ , while the equality holds from the definition of  $\hat{\xi}$ .

516 We can now use Eq. (20) to bound  $\mathbb{E}_{q_{\hat{\xi}}} [\|B(\mathbf{w})\|_D^2]$ , thus obtaining:

$$\mathbb{E}_{q_{\hat{\xi}}} [\|B(\mathbf{w})\|_{\nu}^2] \leq \inf_{\xi \in \Xi} \left\{ \mathbb{E}_{q_{\xi}} [\|B(\mathbf{w})\|_{\nu}^2] + \mathbb{E}_{q_{\xi}} [v(\mathbf{w})] + 2 \frac{\lambda}{N} KL(q_{\xi} \| p) \right\} + 2 \frac{R_{max}^2}{(1-\gamma)^2} \sqrt{\frac{\log \frac{2}{\delta}}{N}}$$

517 This concludes the proof.  $\square$

## 518 B Additional Details on the Algorithms

### 519 B.1 Gaussian Variational Transfer

520 Under Gaussian distributions, all quantities of interest for using Alg. 1 can be computed very easily.  
521 The KL divergence between the prior and approximate posterior can be computed in closed-form as:

$$KL(q_{\xi}(\mathbf{w}) \| p(\mathbf{w})) = \frac{1}{2} \left( \log \frac{|\Sigma_p|}{|\Sigma|} + \text{Tr}(\Sigma_p^{-1} \Sigma) + (\boldsymbol{\mu} - \boldsymbol{\mu}_p)^T \Sigma_p^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_p) - K \right) \quad (22)$$

522 for  $\xi = (\boldsymbol{\mu}, \mathbf{L})$  and  $\Sigma = \mathbf{L}\mathbf{L}^T$ . Its gradients with respect to the variational parameters are:

$$\nabla_{\boldsymbol{\mu}} KL(q_{\xi}(\mathbf{w}) \| p(\mathbf{w})) = \Sigma_p^{-1} (\boldsymbol{\mu} - \boldsymbol{\mu}_p) \quad (23)$$

523

$$\nabla_{\mathbf{L}} KL(q_{\xi}(\mathbf{w}) \| p(\mathbf{w})) = \Sigma_p^{-1} \mathbf{L} - (\mathbf{L}^{-1})^T \quad (24)$$

524 Finally, the gradients w.r.t. the expected likelihood term of the variational objective (2) can be  
525 computed using the reparameterization trick (e.g., [12, 25]):

$$\nabla_{\boldsymbol{\mu}} \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{L}\mathbf{L}^T)} [\|B(\mathbf{w})\|_D^2] = \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\nabla_{\mathbf{w}} \|B(\mathbf{w})\|_D^2] \quad \text{for } \mathbf{w} = \mathbf{L}\mathbf{v} + \boldsymbol{\mu} \quad (25)$$

526

$$\nabla_{\mathbf{L}} \mathbb{E}_{\mathbf{w} \sim \mathcal{N}(\boldsymbol{\mu}, \mathbf{L}\mathbf{L}^T)} [\|B(\mathbf{w})\|_D^2] = \mathbb{E}_{\mathbf{v} \sim \mathcal{N}(\mathbf{0}, \mathbf{I})} [\nabla_{\mathbf{w}} \|B(\mathbf{w})\|_D^2 \cdot \mathbf{v}^T] \quad \text{for } \mathbf{w} = \mathbf{L}\mathbf{v} + \boldsymbol{\mu} \quad (26)$$

### 527 B.2 Mixture of Gaussian Variational Transfer

## 528 C Additional Details on the Experiments

### 529 C.1 The Rooms Problem

### 530 C.2 Classic Control

### 531 C.3 Maze Navigation