
Exploiting Hierarchy for Learning and Transfer in KL-regularized RL

Dhruva Tirumala^{*1} Hyeonwoo Noh^{*2} Alexandre Galashov¹ Leonard Hasenclever¹ Arun Ahuja¹
Greg Wayne¹ Razvan Pascanu¹ Yee Whye Teh¹ Nicolas Heess¹

Abstract

As reinforcement learning agents are tasked with solving more challenging and diverse tasks, the ability to incorporate prior knowledge into the learning system and the ability to exploit reusable structure in solution space is likely to become increasingly important. The KL-regularized expected reward objective constitutes a convenient tool to this end. It introduces an additional component, a default or prior behavior, which can be learned alongside the policy and as such partially transforms the reinforcement learning problem into one of behavior modelling. In this work we consider the implications of this framework in case where both the policy and default behavior are augmented with latent variables. We discuss how the resulting hierarchical structures can be exploited to implement different inductive biases and how the resulting modular structures can be exploited for transfer. Empirically we find that they lead to faster learning and transfer on a range of continuous control tasks.

1. Introduction

Recent advances have greatly improved data efficiency, scalability, and stability of reinforcement learning (RL) algorithms leading to successful applications in a number of domains (Heess et al., 2017; Riedmiller et al., 2018; OpenAI et al., 2018; Mnih et al., 2015; OpenAI, 2018; Silver et al., 2016). Many problems, however, still remain challenging to solve or require large numbers of interactions with the environment; a situation that is likely to get worse as we attempt to tackle increasingly challenging and diverse problems.

A general avenue for reducing sample complexity in machine learning is the use of appropriate prior knowledge,

^{*}Equal contribution ¹DeepMind, London, UK ²Department of Computer Science and Engineering, POSTECH, Pohang, Korea. Correspondence to: Dhruva Tirumala <dhruvat@google.com>, Hyeonwoo Noh <shgusdngogo@postech.ac.kr>.

or inductive biases. In RL, curricula, learning and transfer across task distributions, and the use of demonstrations have been shown to be powerful tools in this regard. The success of these ideas, and more generally our ability build agents that have a chance of succeeding in lifelong learning scenarios, relies on the ability to transfer and generalize behaviours across tasks, and thus on mechanisms for extracting knowledge from existing solutions and for using such knowledge to shape the solutions to new problems.

Recently (Teh et al., 2017; Galashov et al., 2019) have developed a framework for knowledge representation and transfer in RL that allows to express and exploit inductive biases in reinforcement learning problems. The approach relies on a connection between the KL-regularized objective (Todorov, 2007; Kappen et al., 2012; Rawlik et al., 2012; Schulman et al., 2017a) and probabilistic models. One way to think of a policy is that, together with the environment dynamics, it defines a distribution over trajectories, and the framework uses probabilistic models of trajectories to express prior knowledge about the solution distribution of a RL problem. To this end it introduces a second component, a prior or default behaviour, to which the policy is encouraged to remain close in terms of the Kullback-Leibler (KL) divergence. This prior can simplify the learning problem, for instance by reducing the trajectory space that needs to be searched, or by steering the learning algorithm away from undesirable solutions. Rather than manually crafting specific constraints which may or may not be appropriate for a given task or task distribution, it is possible to learn the prior from data, and (Teh et al., 2017; Galashov et al., 2019) have shown that this can be effective both in multi-task and transfer learning scenarios. (Galashov et al., 2019) demonstrate how information asymmetry, i.e. restricting the prior’s access to certain parts of the state space can be used to selectively transfer or generalize certain aspects of a learned behavior across tasks or different parts of the state space.

In this paper we extend and generalize this framework and study priors and policies that are hierarchically structured and augmented with latent variables. The introduction of latent variables can increase model capacity, and allow for richer classes of distributions (e.g. non-Gaussian). Importantly, it enables a broader range of inductive biases: It allows us to express more specific constraints on the infor-

mation processing capabilities of the agent leading to priors that have more diverse generalization properties. In particular, the hierarchical model structure allows us to mirror more closely the natural structure of many problems. For instance, in motor control, one tends to observe a separation into rapidly varying motor commands and slower changing higher-level goals. Similarly, different types of information tend to be processed by different levels of the system (e.g. vision vs. proprioception).

The hierarchical formulation also gives rise to a modular structure that allows to selectively constrain, transfer and generalize certain aspects of the behavior, such as low-level skills or high-level goals. Exploiting the compositionality of probabilistic models, this modularity also facilitates the immediate reuse of some model components in a principled manner, such as the sharing of low-level skills across tasks. This allows imposing more flexible constraints and can reduce the number of model parameters leading to an increased statistical efficiency and faster learning.

The presence of latent variables introduces additional complexities and we present a general framework which addresses these in Appendix A. We elaborate on a particular but still general hierarchically structured variant in the main text, and develop suitable efficient off-policy algorithms. We provide empirical results on several tasks with physically simulated bodies and continuous action spaces and discrete grid worlds which demonstrate the benefits of the general framework and especially the advantages of the hierarchical structure, compared to the unstructured policies used in (Galashov et al., 2019).

2. RL as probabilistic modelling

In this section, we briefly review how the KL-regularized objective can connect RL and probabilistic model learning. We will denote states and actions at time t respectively with s_t and a_t . $r(s, a)$ is the instantaneous reward received in state s when taking action a . We will refer to the history up to time t as $x_t = (s_1, a_1, \dots, s_t)$ and the whole trajectory as $\tau = (s_1, a_1, s_2, a_2, \dots)$. The agent policy $\pi(a_t|x_t)$ denotes a distribution over next actions given history x_t , while $\pi_0(a_t|x_t)$ denotes a default or habitual policy.¹ The KL-regularized RL objective (Todorov, 2007; Kappen et al., 2012; Rawlik et al., 2012; Schulman et al., 2017a) takes the form:

$$\mathcal{L}(\pi, \pi_0) = \mathbb{E}_\tau \left[\sum_{t \geq 1} \gamma^t r(s_t, a_t) - \alpha \gamma^t \text{KL}(a_t|x_t) \right] \quad (1)$$

where γ is discount factor and α controls the relative contributions of both terms. $\mathbb{E}_\tau[\cdot]$ is taken w.r.t. the trajec-

¹ We generally work with history dependent policies since we will consider restricting access to state information from policies (for information asymmetry), which may render fully observed MDPs effectively partially observed.

tory distribution given by agent policy and system dynamics: $p(s_1) \prod_{t \geq 1} \pi(a_t|x_t) p(s_{t+1}|s_t, a_t)$. We use a convenient notation² for the KL divergence: $\text{KL}(a_t|x_t) = \mathbb{E}_{\pi(a_t|x_t)} [\log \frac{\pi(a_t|x_t)}{\pi_0(a_t|x_t)}]$.

When optimized with respect to π the objective can be seen to trade off expected reward with closeness (in terms of KL) between trajectories produced by executing π and π_0 (Appendix A). This is also evident from the optimal π in eq. (1)

$$\pi^*(a_t|x_t) = \pi_0(a_t|x_t) \exp \frac{1}{\alpha} (Q^*(x_t, a_t) - V^*(x_t)) \quad (2)$$

$$Q^*(x_t, a_t) = r(s_t, a_t) + \gamma \mathbb{E}_{s_{t+1}|s_t, a_t} [V^*(x_{t+1})] \quad (3)$$

$$V^*(x_t) = \alpha \log \int \pi_0(a_t|x_t) \exp \frac{1}{\alpha} Q^*(x_t, a_t) da_t, \quad (4)$$

where $Q^*(\cdot)$ and $V^*(\cdot)$ are optimal action value and value functions of eq. (1); See (e.g. Rawlik et al., 2012; Fox et al., 2016; Schulman et al., 2017a; Nachum et al., 2017) for derivations. We can thus think of π as a specialization of π_0 that is obtained by tilting π_0 towards high-value actions (as measured by the action value Q).

As discussed in (Galashov et al., 2019) the default behavior thus bears resemblance to a trajectory prior in a particular probabilistic model. Several recent works have considered optimizing eq. (1) when π_0 is of a fixed and simple form. For instance, when π_0 is chosen to be uniform entropy-regularized objective is recovered e.g. (Ziebart, 2010; Fox et al., 2016; Haarnoja et al., 2017; Schulman et al., 2017a; Hausman et al., 2018). More interestingly, in some cases π_0 can be used to inject detailed prior knowledge into the learning problem. In a transfer scenario π_0 can be a *learned* object, and the KL term plays effectively the role of a shaping reward.

π and π_0 can also be co-optimized. In this case the relative parametric forms of π_0 and π are of importance. The optimal π_0 in eq. (1) is

$$\pi_0^*(a_t|x_t) = \arg \max_{\pi_0} \mathbb{E}_{\pi(a_t|x_t)} [\log \pi_0(a_t|x_t)], \quad (5)$$

which maximizes only terms in eq. (1) depending on π_0 . Thus learning π_0 can be seen as supervised learning where π_0 is trained to match the history-conditional action sequences produced by π . It should be clear that $\pi_0 = \pi$ is the optimal solution when π_0 has sufficient information and capacity, and in this scenario the regularizing effect of π_0 is lost. When the information or the capacity of π_0 is limited then π_0 will be forced to generalize the behavior of π . For instance, (Teh et al., 2017) and (Galashov et al., 2019) consider a multitask scenario in which π is given task-identifying information, while π_0 is not. As a result, π_0 is forced to learn a marginal trajectory distribution over

²In the following, $\text{KL}(Y|X)$ always denotes $\mathbb{E}_{\pi(Y|X)} [\log \frac{\pi(Y|X)}{\pi_0(Y|X)}]$ for arbitrary variables X and Y .

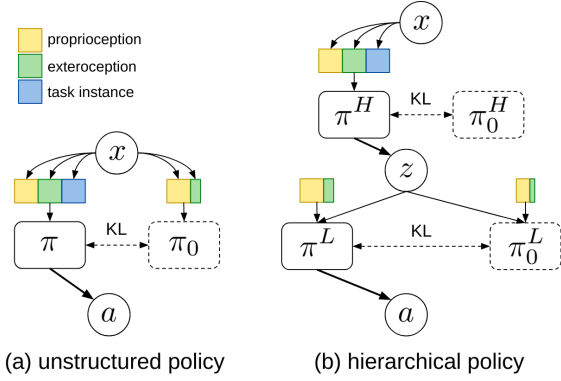


Figure 1. Diagram of the generic structure of the regularized KL-objective considered. (a) shows an unstructured policy, where information asymmetry (e.g. hiding task instance information) is exploited to induce a meaningful default policy π_0 (Galashov et al., 2019); (b) shows the scenario when we use structured policies composed from high-level π^H and low-level π^L policies that communicate through the latent action z . Note that now different forms of information asymmetry can be employed. See text for details.

tasks, which can be considered as a common default behaviour, and this behaviour is then shareable across tasks to regularize π . More generally, appropriate choices for the model classes of π_0 and π will allow us to influence both the learning dynamics as well as the final solutions to eq. (1) and thus provide us with a means of injecting prior knowledge into the learning problem.

3. Hierarchically structured policies

(Galashov et al., 2019) consider simple unstructured models and focus on the information that π_0 has access to. In this paper we focus on a complementary perspective. We explore how variations of the parametric forms of π_0 and π , via the introduction of latent variables, can give rise to richer and hierarchically structured models with different inductive biases and generalization properties. In this section we discuss a particular instantiation of this idea and discuss the general framework in Appendix A.

We consider learning multi-level representations of behaviour. Specifically, we consider a two level architecture in which the high-level decisions are concerned with task objectives but largely agnostic to details of actuation. The low-level control translates the high-level decisions into motor actions while being agnostic to task objectives. Successful learning of such abstractions can be useful to exploit repetitive structures within or across tasks. As two use cases we consider (a) multi-task control where different tasks require similar motor-skills; as well as (b) a scenario where we aim to solve similar tasks with different actuation systems.

We are interested in the role of the KL-regularized objective and the default policies for learning the desired multi-level

abstractions. In this view, the hierarchical structure and latent variables mirror the structure of the abstractions that we want to learn, and the structured default policy allows fine-grained control of its regularization effect in each level of abstractions. Besides, the latent variables allow us to work with richer (e.g. non-Gaussian) distribution classes, allow us to model temporal correlations, e.g. in the high-level goals, and they give rise to modular structure that enables parameter sharing.

Conceptually policies are divided into high-level and low-level components which interact via auxiliary latent variables. For concreteness, but without loss of generality, let z_t be a (continuous) latent variable for each time step t (we discuss alternative choices such as latent variables that are sampled infrequently in Appendices A and C). The agent policy is extended as $\pi(a_t, z_t|x_t) = \pi^H(z_t|x_t)\pi^L(a_t|z_t, x_t)$ and likewise for the default policy π_0 . z_t can be interpreted as a high-level or abstract action, taken according to the high-level (HL) controller π^H , and which is translated into low-level or motor action a_t by the low-level (LL) controller π^L . We extend the histories x_t and trajectories τ to appropriately include z_t 's. Note that as we elaborate in the appendix this allows for temporally correlated z_t 's including the case where z is only sampled once at the beginning of the episode. As will be discussed in Section 5, structuring a policy into HL and LL controllers has been studied e.g. (Heess et al., 2016; Hausman et al., 2018; Haarnoja et al., 2018a; Merel et al., 2019), but the concept of a default policy has not been widely explored in this context. We discuss the differences between these works and ours in detail in Appendix A.

In case z_t 's can take on many values or are continuous, the objective (1) becomes intractable as the marginal distributions $\pi(a_t|x_t)$ and $\pi_0(a_t|x_t)$ in the KL divergence cannot be computed in closed form. As discussed in more detail in Appendix A this problem can be addressed in different ways. For simplicity and concreteness we here assume that the latent variables in π and π_0 have the same dimension and semantics. We can then construct a lower bound for the objective by using the following upper bound for the KL:

$$\text{KL}(a_t|x_t) \leq \text{KL}(z_t|x_t) + \mathbb{E}_{\pi(z_t|x_t)}[\text{KL}(a_t|z_t, x_t)], \quad (6)$$

which is tractably approximated using Monte Carlo sampling. The derivation is in Appendix C.1. Note that:

$$\text{KL}(z_t|x_t) = \text{KL}(\pi^H(z_t|x_t) \parallel \pi_0^H(z_t|x_t)) \quad (7)$$

$$\text{KL}(a_t|z_t, x_t) = \text{KL}(\pi^L(a_t|z_t, x_t) \parallel \pi_0^L(a_t|z_t, x_t)). \quad (8)$$

The resulting lower bound for the objective is:

$$\mathcal{L}(\pi, \pi_0) \geq \mathbb{E}_{\tau} \left[\sum_{t \geq 1} \gamma^t r(s_t, a_t) - \alpha \gamma^t \text{KL}(z_t|x_t) - \alpha \gamma^t \text{KL}(a_t|z_t, x_t) \right], \quad (9)$$

where τ is a trajectory that appropriately includes z_t 's. Full derivation including discount terms is in Appendix C.2. In this paper we consider eq. (9) as a main objective function.

3.1. Structuring the default policy

Compared to e.g. (Galashov et al., 2019) who mostly use conditional Gaussian distributions to model the default policy, the presence of latent variables may allow more flexible trajectory models that can, in principle, capture richer, non-Gaussian state-conditional distributions as a marginal $\pi_0(a_t|x_t) = \int_{z_t} \pi_0^H(z_t|x_t)\pi_0^L(a_t|z_t, x_t)dz_t$, and thus to approximate π_0^* more accurately.³ We can also model different marginal distribution based on the parameterization of $\pi_0^H(z_t|x_t)$, which could be a useful inductive bias for generalizing the learned trajectory distribution.

In this work, we consider the following choices of of HL default policy: **Independent isotropic Gaussian.** $\pi_0^H(z_t|x_t) = \mathcal{N}(z_t|0, 1)$, i.e. the HL default policy assumes the abstract actions to be context independent. **AR(1) process.** $\pi_0^H(z_t|x_t) = \mathcal{N}(z_t|\alpha z_{t-1}, \sqrt{1-\alpha^2})$, i.e. the HL default policy is a first-order auto-regressive process with a fixed parameter $0 \leq \alpha < 1$ chosen to ensure a marginal distribution $\mathcal{N}(0, 1)$. This allows for more structured temporal dependence among the abstract actions. **Learned AR process.** Similar to the AR(1) process this default HL policy allows z_t to depend on z_{t-1} but now the high-level default policy is a Gaussian distribution with mean and variance that are learned functions of z_{t-1} with parameters ϕ : $\pi_0^H(z_t|x_t) = \mathcal{N}(z_t|\mu_\phi(z_{t-1}), \sigma_\phi^2(z_{t-1}))$. Note that the considered HL default policies are not conditioned on x_t . This is a form of information asymmetry for capturing task agnostic trajectory distribution as we will discuss in the following.

Regularizing via information asymmetries As discussed in (Galashov et al., 2019) restricting the information available to different policies is a powerful tool to force regularization and generalization. In our case we let this information asymmetry be reflected also in the separation between HL and LL controllers (see Figure 1). Specifically we introduce a separation of concerns between π^L and π^H by providing full information only to π^H while information provided to π^L is limited. In our experiments we vary the information provided to π^L ; it receives body-specific (proprioceptive) information as well as different amounts of environment-related (exteroceptive) information. The task is only known to π^H . Hiding task specific information from the LL controller makes it easier to transfer across tasks. It

³ Note that π_0^* is itself obtained by marginalizing over different contexts, e.g. different task-specific policies, cf. eq. (5)). Thus, even though the individual task-conditional policies may be well represented by conditional Gaussian distributions, their mixture may not be.

forces π^L to focus on learning task agnostic behaviour, and to rely on the abstract actions selected by π^H to solve the task. Similarly, we hide task specific information from both π_0^L and π_0^H , which in turn makes the marginal default policy π_0 to be agnostic to the task (see above). In the experiments we further consider transferring the HL controller across bodies, in situations where the abstract task is the same but the body changes. Here we additionally hide body-specific information from π^H , so that the HL controller is forced to learn body-agnostic behaviour.

Partial parameter sharing An advantage of the hierarchical structure of the policies is that it enables several options for partial parameter sharing, which when used in the appropriate context can make learning more statistically efficient. We explore the utility of sharing low-level controllers between agent and default policy, i.e. $\pi^L(a_t|z_t, x_t) = \pi_0^L(a_t|z_t, x_t)$, and study the associated trade-offs in different learning scenarios. In the case of sharing, the new lower bound is:

$$\mathcal{L}(\pi, \pi_0) \geq \mathbb{E}_\tau \left[\sum_{t \geq 1} \gamma^t r(s_t, a_t) - \alpha \gamma^t \text{KL}(z_t|x_t) \right] \quad (10)$$

This objective function is similar in spirit to current KL-regularized RL approaches discussed in Section 2, except that the KL divergence is between policies defined on abstract actions z_t as opposed to concrete actions a_t . The effect of this KL divergence is that it regularizes both the HL policies as well as the space of behaviours parameterised by the abstract actions. This special case of our framework also reveals a connection to (Goyal et al., 2019), where eq. (10) was motivated as an approximation of information bottleneck for learning a goal conditioned policy. Here, we provide a different perspective of eq. (10) in terms of probabilistic modeling and structured policies as a special case of eq. (9) (and more generally the framework in Appendix A), which suggests a broader range of algorithms that were not discussed in (Goyal et al., 2019); see Section 5 for detailed discussion.

4. Algorithm

We jointly optimize the default policy and the agent's policy, while the agent's policy is regularized by a *target default policy*, which is periodically updated to a new default policy. The objective for the default policies is similar to distillation (Parisotto et al., 2016; Rusu et al., 2016) or supervised learning, where agent's policies define the data distribution. Note that due to the particular way we lower bound the KL (see eq. (9)) the supervised step remains unproblematic despite the presence of the latent variables in π_0 .

To efficiently optimize the hierarchical policy with latent variables, we develop off-policy learning algorithm based on SVG(0) (Heess et al., 2015) with experience replay.

Algorithm 1 On-policy version of our algorithm

Flat (HL+LL) policy: $\pi_\theta(a_t|\epsilon_t, x_t)$, parameter θ
 Reparameterized HL policy: $z_t = f^H(x_t, \epsilon_t)$
 Default policy: $\pi_{0,\phi}^H(z_t|x_t)$, $\pi_{0,\phi}^L(a_t|z_t, x_t)$, parameter ϕ
 Q-function: $Q_\psi(a_t, z_t, x_t)$, parameter ψ
repeat
 for $t = 0, K, 2K, \dots, T$ **do**
 Rollout trajectory: $\tau_{t:t+K} = (s_t, a_t, r_t, \dots, r_{t+K})$
 Sample latent: $z_{t'} = f^H(x_{t'}, \epsilon_{t'})$, where $\epsilon_{t'} \sim \rho(\epsilon)$
 Compute KL: $\hat{\mathbf{K}}_{L_{t'}} = \text{KL}(z|x_{t'}) + \text{KL}(a|z_{t'}, x_{t'})$
 Bootstrap: $\hat{V} = \mathbb{E}_\pi Q(a, z_{t+K}, x_{t+K}) - \alpha \hat{\mathbf{K}}_{L_{t+K}}$
 Estimate Q target: $\hat{Q}_{t'} = \sum_{i=t'}^{t'+K-1} (r_i - \alpha \hat{\mathbf{K}}_{L_i}) + \hat{V}$
 Policy loss: $\hat{L}_\pi = \sum_{i=t}^{t+K-1} \mathbb{E}_\pi Q(a, z_i, x_i) - \alpha \hat{\mathbf{K}}_{L_i}$
 Q-value loss: $\hat{L}_Q = \sum_{i=t}^{t+K-1} \|\hat{Q}_i - Q(a, z_i, x_i)\|^2$
 Default policy loss: $\hat{L}_{\pi_0} = \sum_{i=t}^{t+K-1} \hat{\mathbf{K}}_{L_i}$
 $\theta \leftarrow \theta + \beta_\pi \nabla_\theta \hat{L}_\pi$ $\phi \leftarrow \phi + \beta_{\pi_0} \nabla_\phi \hat{L}_{\pi_0}$
 $\psi \leftarrow \psi - \beta_Q \nabla_\psi \hat{L}_Q$
end for
until

We follow a strategy similar to (Heess et al., 2016) and reparameterize $z_t \sim \pi^H(z_t|x_t)$ as $z_t = f^H(x_t, \epsilon_t)$, where $\epsilon_t \sim \rho(\epsilon_t)$ is a fixed noise distribution and $f^H(\cdot)$ is a deterministic function. In practice this means that the hierarchical policy can be treated as a flat policy $\pi(a_t|\epsilon_t, x_t) = \pi^L(a_t|f^H(x_t, \epsilon_t), x_t)$. We reparameterize action distribution of the flat policy $\pi(a_t|\epsilon_t, x_t)$ and optimize it by back-propagating the gradient from an action value function $Q(a_t, z_t, x_t)$. Note that the action value function is conditioned on z_t because it will be contained in history from the next time step and affect the policy accordingly. The action value function is optimized to match a target action value estimated by Retrace (Munos et al., 2016), which provides low variance estimate of action value from K-step windows of off-policy trajectories. Note that the off-policy learning with hierarchical policy requires extra care in terms of handling latent variables in action value function and applying Retrace. We describe these details in Appendix B. For illustration, Algorithm 1 provides the pseudo-code for a simple on-policy version of our algorithm. We implement our algorithm in a distributed setup similar to (Riedmiller et al., 2018) where multiple actors are used to collect trajectories and a single learner is used to optimize model parameters. Similarly to other KL-regularized RL approaches e.g. (Teh et al., 2017; Galashov et al., 2019), we additionally regularize the entropy of π^L to encourage exploration. More details about the learning algorithms are in Appendix B.

5. Related Work

Entropy regularized reinforcement learning (RL), also known as maximum entropy RL (Ziebart, 2010; Kappen

et al., 2012; Toussaint, 2009) is a special case of KL regularized RL. This framework connects probabilistic inference and sequential decision making problems. Recently, this idea has been adapted to deep reinforcement learning (Fox et al., 2016; Schulman et al., 2017a; Nachum et al., 2017; Haarnoja et al., 2017; Hausman et al., 2018; Haarnoja et al., 2018b). Another instance of KL regularized RL includes trust region based methods (Schulman et al., 2015; 2017b; Wang et al., 2017; Abdolmaleki et al., 2018). They use KL divergence between new policy and old policy as a trust region constraints for conservative policy update.

Introducing a parameterized default policy provides a convenient way to transfer knowledge or regularize the policy. Schmitt et al. (Schmitt et al., 2018) use a pretrained policy as the default policy; other works jointly learn the policy and default policy to capture reusable behaviour from experience (Teh et al., 2017; Czarnecki et al., 2018; Galashov et al., 2019; Grau-Moya et al., 2019). To retain the role of default policy as a regularizer, it has been explored to restrict its input (Galashov et al., 2019; Grau-Moya et al., 2019), parameteric form (Czarnecki et al., 2018) or to share it across different contexts (Teh et al., 2017; Ghosh et al., 2018).

Another closely related regularization for RL is using information bottleneck (Tishby & Polani, 2011; Still & Precup, 2012; Rubin et al., 2012; Ortega & Braun, 2013; Tiomkin & Tishby, 2017). Galashov et al. (Galashov et al., 2019) discussed the relation between information bottleneck and KL regularized RL. Strouse et al. (Strouse et al., 2018) learn to hide or reveal information for future use in multi-agent cooperation or competition. Goyal et al. (Goyal et al., 2019) consider identifying bottleneck states based on eq. (10) and using them for exploration during transfer. However, unlike (Goyal et al., 2019) motivating the objective simply as an approximation to an information bottleneck, we provide a new perspective that considers structuring policies with latent variables as a way of injecting inductive bias in the KL-regularized RL framework. This new perspective motivates the eq. (10) as one instance of a broader family of algorithms that were not discussed in (Goyal et al., 2019), and we demonstrate how different instances of the proposed family of algorithms can be useful in different learning scenarios.

The hierarchical RL literature (Dayan & Hinton, 1993; Parr & Russell, 1998; Sutton et al., 1999) has studied hierarchy extensively as a means to introduce inductive bias. Among various ways (Sutton et al., 1999; Bacon et al., 2017; Vezhnevets et al., 2017; Nachum et al., 2018; 2019; Xie et al., 2018; Lee et al., 2019), our approach resembles (Heess et al., 2016; Hausman et al., 2018; Haarnoja et al., 2018a; Merel et al., 2019), in that a HL controller modulates a LL controller through a continuous channel. For learning the LL

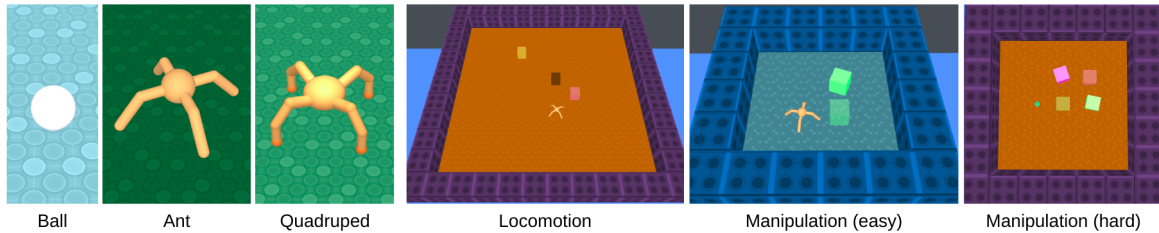


Figure 2. **Bodies and tasks for experiments.** **Left:** All considered bodies. **Right:** Example tasks.

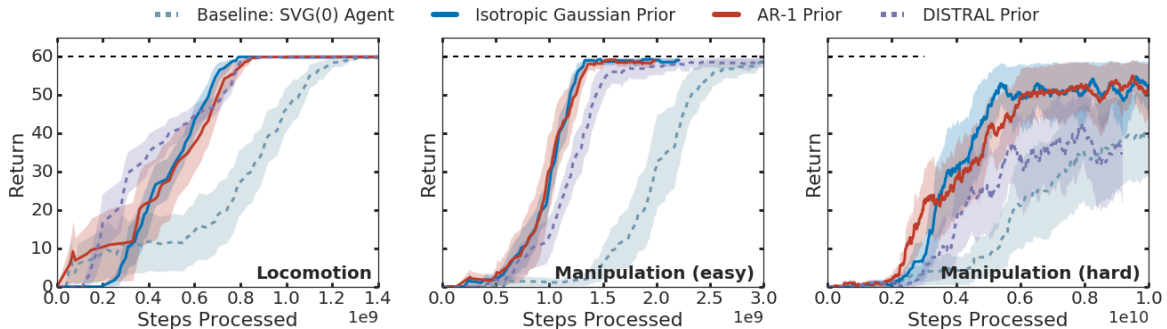


Figure 3. **Learning with structured policies.** **Left:** Locomotion with the Ant. **Center:** Manipulation (easy) with the Ant. **Right:** Manipulation (hard) with the Ball. The proposed models with hierarchy are denoted with the type of HL default policy: Isotropic Gaussian, AR-1, AR-Learned.

controller, imitation learning (Fox et al., 2017; Krishnan et al., 2017; Merel et al., 2019), unsupervised learning (Gregor et al., 2017; Eysenbach et al., 2019) and meta learning (Frans et al., 2018) have been employed. Similar to our approach, (Heess et al., 2016; Florensa et al., 2017; Hausman et al., 2018) use a pretraining task to learn a reusable LL controller. However, the concept of a default policy has not been widely explored in this context.

6. Experiments

We evaluate our method in several environments with continuous state and action spaces. We consider a set of structured, sparse reward tasks that can be executed by multiple bodies with different degrees of freedom. The tasks and bodies are illustrated in Figure 2.

We consider task distributions that are designed such that their solutions exhibit significant overlap in trajectory space so that transfer can reasonably be expected. They are further designed to contain instances of variable difficulty and hence provide a natural curriculum. Our tasks are as follows. **Locomotion:** reaching a specific target among 3 locations. **Locomotion with gap:** reaching the end of a corridor with a gap of variable length. Solving the task requires being able to jump over the gap. **Manipulation:** moving one of N boxes to one of K targets as indicated by the environment. (easy): $N=1, K=1$, (mid): $N=1, K=3$, (hard): $N=2, K=2$. **Manipulation (heavy):** manipulation (easy) with heavier box. **Manipulation (gather):** moving two boxes such that they are in contact with each other. **Combined task (And**

Or): moving a box to one target (And / Or) go to the other target in a single episode. See Appendix E for exact configurations.

We use three simulated bodies: Ball, Ant, and Quadruped. Ball and Ant have been used in (Heess et al., 2017; Xie et al., 2018; Galashov et al., 2019); we introduce the Quadruped as a variant of the Ant. The **Ball** has 2 actuators for moving forward or backward and turning left or right. The **Ant** has 8 actuators for moving its legs to walk and to interact with objects. The **Quadruped** is similar to the Ant, but with 12 actuators. Each body is characterized by a different set of proprioceptive features. Further details of the bodies are in Appendix E.

Throughout the experiments, we use 32 actors to collect trajectories and a single learner to optimize the model. We plot average episode return with respect to the number of steps processed by the learner. Note that the number of steps is different from the number of agent’s interaction with environment, because the collected trajectories are processed multiple times by a centralized learner to update model parameters. Hyperparameters, including KL cost and action entropy regularization cost, are optimized on a per-task basis. Details are provided in Appendices F and G.

6.1. Learning with structured policies

We study how the KL-regularized objective with different structures and parameterizations discussed in Section 3 affects learning. Our baselines are SVG-0 with entropy regularization and an unstructured KL regularized policy similar

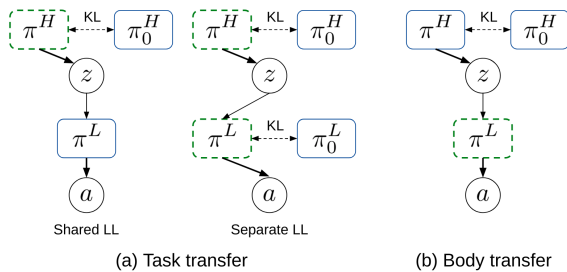


Figure 4. **Transfer learning scenarios.** The blue boxes denote modules that are transferred and fixed, and the green dotted boxes denote modules learned from scratch.

to (Galashov et al., 2019; Teh et al., 2017) (DISTRAL prior). As described in Section 3 we employ hierarchical structure with shared LL components (Shared LL) and separate LL components (Separate LL). Unless otherwise stated, we use Shared LL as our default hierarchical model. The HL controller receives full information while the LL controller (and hence the default policy) receives proprioceptive information plus the positions of the box(es) as indicated. The same information asymmetry is applied to the DISTRAL prior i.e. the default policy receives proprioception plus box positions as inputs. We explore multiple HL default policies: Isotropic Gaussian, AR(1) process, and learned AR prior.

Figure 3 illustrates the results of the experiment. Our main finding is that the KL regularized objective significantly speeds up learning, and that the hierarchical structure does as well or better than the flat, DISTRAL formulation. The gap increases for more challenging tasks (e.g. Manipulation (hard)).

7. Transfer Learning

The hierarchical structure introduces a modularity of the policy and default policy, which can be utilized for transfer learning. We consider two transfer scenarios (see Figure 4): 1) task transfer where we reuse the learned default policy to solve novel tasks, and 2) body transfer, where reusing the body agnostic HL policy and default policy transfers the goal directed behaviour to another body.

7.1. Task transfer

We consider transfer between task distributions whose solutions exhibit significant shared structure, e.g. because solution trajectories can be produced by a common set of skills or repetitive behaviour. If the default policy can capture and transfer this reusable structure it will facilitate learning similar tasks. Transfer then involves specializing the default behavior to the needs of the target task (e.g. by directing locomotion towards a goal).

For task transfer, we reuse pretrained goal agnostic com-

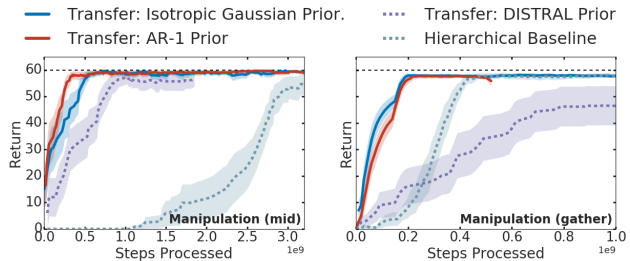


Figure 5. **Transfer learning with structured policies.** **Left:** From Manipulation (easy) to Manipulation (mid) with Ant. **Right:** From Manipulation (hard) to Manipulation (gather) with Ball.

ponents, including π_0^H and π_0^L . We learn a new π^H and either set the LL policy π^L identical to the pretrained LL default policy π_0^L (Shared LL), or allow π^L to diverge from π_0^L (Separate LL). In the case of Shared LL, similarly e.g. to (Heess et al., 2016; Hausman et al., 2018), the new HL policy π^H learns to control the LL policy π^L in the latent space. Unlike in previous work, however, we regularize π^H with π_0^H . Transferring the learned default policy to a new task is similar to (Galashov et al., 2019), but our approach is different in that we exploit the modular structure of π_0 and π ; It allows to re-use the pretrained π^L (Shared LL) or to initialize the new LL policy with the pretrained parameters (Separate LL with weight initialization).

We consider two baselines: (a) the Shared LL model learned from scratch (Hierarchical Agent); (b) a DISTRAL prior, i.e. we transfer a pretrained unstructured default policy to regularize the policy for the target task. The first baseline allows us to assess the benefit of transfer; the second baseline provides an indication whether the hierarchical policy structure is beneficial. Additionally, we compare different types of HL default policies. Specifics of the experiments including the information provided to HL and LL are provided in Appendix E.

Figure 5 shows the results with Shared LL. Transferring the pretrained default policy can bring clear benefits for related tasks. The hierarchical structure, which facilitates partial parameter sharing and transferring flexible non-Gaussian default policy, performs better than the DISTRAL prior regardless of type of HL default policy. To assess the benefit of the more flexible trajectory model we measure the KL divergence between π and π_0 in Manipulation (mid) for the DISTRAL prior and the hierarchical latent variable model with Gaussian prior. Values of 11.35 and 2.64 respectively are consistent with the idea that the latent variable default policy provides a better prior over the trajectory distribution and thus allows for a more favorable trade-off between KL and task reward. Figure 6 illustrates the benefits of the flexibility afforded by the full model in eq. (9) compared to eq. (10). In the presented transfer scenarios the adaptation of low level skills is required to learn the target task. Specifically, the transfer task requires the acquisition of a

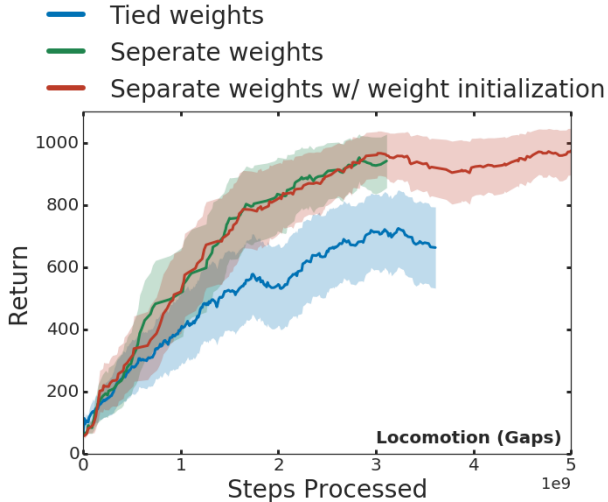


Figure 6. **Effect of partial parameter sharing.** From Locomotion (easy) to Locomotion (Gap) with Ant, AR-1.

new skill (jumping) in order to consistently solve the task. Allowing the π^L to diverge from π_0^L as in eq. (9) turns out to be critical in these scenarios.

7.2. Body transfer

We study whether our architecture can facilitate transfer of high-level, goal-directed behavior between bodies with different actuation systems. To this end we reuse the pre-trained body-agnostic components, HL policy π^H and the default policy π_0^H . We learn a new body-specific LL policy π^L , which is assumed to be shared with LL default policy π_0^L (see Figure 4b). The transferred HL policies provide goal-specific behaviour actuated on the latent space, which can then be instantiated by learning a new LL policy. During transfer, we optimize eq. (10) to exploit KL between the HL policies as a dense reward signal to guide learning new LL policy. Here, using the KL between the HL policies during transfer looks superficially similar to (Goyal et al., 2019), but their motivation and method are different. They rely on the *positive KL* as an exploration bonus, trying to get the agent to explore the new state. We rely on *negative KL*, encouraging the new policy to be close to the default policy. We discuss further differences with (Goyal et al., 2019) in Section 5.

Figure 7 illustrates result for transferring from Ant to Quadruped in Locomotion task with egocentric vision as observation input. We compare our approach with baseline learning the hierarchical policy from scratch and analyze the effects of the KL regularization. It shows that transferring the HL component and using KL regularization works best in this setting. We observed that this result is consistent in multiple settings with different tasks and bodies both in continuous control and discrete grid world environments. These

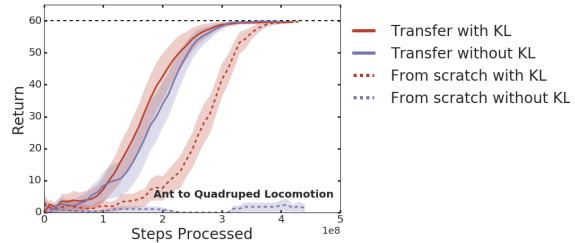


Figure 7. **Body transfer from vision, AR-1 Prior.** Ant to Quadruped, Locomotion. Task information is provided by egocentric vision.

additional experiments can be found in Appendix D.3.

8. Discussion

In this work we have studied the benefit of learned probabilistic models of default behaviors for multi-task and transfer scenarios in RL. In particular, we have outlined a generic framework for hierarchically structured models together with suitable off-policy RL algorithms. The hierarchical model structure allows to model richer distributions, can give rise to modular network structure and generally allows to express a broad range of inductive biases with different generalization properties. Here, we have studied a particular model variant and shown its empirical advantages. Looking forward we believe that as the importance of multitask and lifelong learning scenarios will grow, frameworks such as ours, that provide flexible mechanisms for introducing prior knowledge and reusing previous learned solutions will provide a fertile ground for advancing RL algorithms.

References

- Abdolmaleki, A., Springenberg, J. T., Tassa, Y., Munos, R., Heess, N., and Riedmiller, M. Maximum a posteriori policy optimisation. In *International Conference on Learning Representations*, 2018.
- Agakov, F. V. and Barber, D. An auxiliary variational method. In *Neural Information Processing, 11th International Conference, ICONIP 2004, Calcutta, India, November 22-25, 2004, Proceedings*, pp. 561–566, 2004.
- Bacon, P.-L., Harb, J., and Precup, D. The option-critic architecture. In *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- Czarnecki, W., Jayakumar, S., Jaderberg, M., Hasenclever, L., Teh, Y. W., Heess, N., Osindero, S., and Pascanu, R. Mix & match agent curricula for reinforcement learning. In *Proceedings of the 35th International Conference on Machine Learning*, 2018.
- Dayan, P. and Hinton, G. E. Feudal reinforcement learning.

- In *Advances in neural information processing systems*, pp. 271–278, 1993.
- Espeholt, L., Soyer, H., Munos, R., Simonyan, K., Mnih, V., Ward, T., Doron, Y., Firoiu, V., Harley, T., Dunning, I., et al. Impala: Scalable distributed deep-rl with importance weighted actor-learner architectures. In *Proceedings of the 35th International Conference on Machine Learning*, 2018.
- Eysenbach, B., Gupta, A., Ibarz, J., and Levine, S. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations*, 2019.
- Florensa, C., Duan, Y., and Abbeel, P. Stochastic neural networks for hierarchical reinforcement learning. In *International Conference on Learning Representations*, 2017.
- Fox, R., Pakman, A., and Tishby, N. Taming the noise in reinforcement learning via soft updates. In *Proceedings of the Thirty-Second Conference on Uncertainty in Artificial Intelligence*, pp. 202–211. AUAI Press, 2016.
- Fox, R., Krishnan, S., Stoica, I., and Goldberg, K. Multi-level discovery of deep options. *CoRR*, abs/1703.08294, 2017. URL <http://arxiv.org/abs/1703.08294>.
- Frans, K., Ho, J., Chen, X., Abbeel, P., and Schulman, J. Meta learning shared hierarchies. In *International Conference on Learning Representations*, 2018.
- Galashov, A., Jayakumar, S., Hasenclever, L., Tirumala, D., Schwarz, J., Desjardins, G., Czarnecki, W. M., Teh, Y. W., Pascanu, R., and Heess, N. Information asymmetry in KL-regularized RL. In *International Conference on Learning Representations*, 2019.
- Ghosh, D., Singh, A., Rajeswaran, A., Kumar, V., and Levine, S. Divide-and-conquer reinforcement learning. In *International Conference on Learning Representations*, 2018.
- Goyal, A., Islam, R., Strouse, D., Ahmed, Z., Larochelle, H., Botvinick, M., Levine, S., and Bengio, Y. Transfer and exploration via the information bottleneck. In *International Conference on Learning Representations*, 2019.
- Grau-Moya, J., Leibfried, F., and Vrancx, P. Soft q-learning with mutual-information regularization. In *International Conference on Learning Representations*, 2019.
- Gregor, K., Rezende, D. J., and Wierstra, D. Variational intrinsic control. In *International Conference on Learning Representations*, 2017.
- Haarnoja, T., Tang, H., Abbeel, P., and Levine, S. Reinforcement learning with deep energy-based policies. In *International Conference on Machine Learning*, pp. 1352–1361, 2017.
- Haarnoja, T., Hartikainen, K., Abbeel, P., and Levine, S. Latent space policies for hierarchical reinforcement learning. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 1851–1860, 2018a.
- Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 1861–1870, 2018b.
- Hausman, K., Springenberg, J. T., Wang, Z., Heess, N., and Riedmiller, M. Learning an embedding space for transferable robot skills. In *International Conference on Learning Representations*, 2018.
- Heess, N., Wayne, G., Silver, D., Lillicrap, T., Erez, T., and Tassa, Y. Learning continuous control policies by stochastic value gradients. In *Advances in Neural Information Processing Systems*, 2015.
- Heess, N., Wayne, G., Tassa, Y., Lillicrap, T., Riedmiller, M., and Silver, D. Learning and transfer of modulated locomotor controllers. *arXiv preprint arXiv:1610.05182*, 2016.
- Heess, N., Tirumala, D., Sriram, S., Lemmon, J., Merel, J., Wayne, G., Tassa, Y., Erez, T., Wang, Z., Eslami, A., Riedmiller, M., et al. Emergence of locomotion behaviours in rich environments. *arXiv preprint arXiv:1707.02286*, 2017.
- Johnson, M., Duvenaud, D. K., Wiltchko, A., Adams, R. P., and Datta, S. R. Composing graphical models with neural networks for structured representations and fast inference. In *Advances in Neural Information Processing Systems 29*, pp. 2946–2954. 2016.
- Kappen, H. J., Gómez, V., and Opper, M. Optimal control as a graphical model inference problem. *Machine learning*, 87(2):159–182, 2012.
- Krishnan, S., Fox, R., Stoica, I., and Goldberg, K. DDCO: discovery of deep continuous options for robot learning from demonstrations. *CoRR*, abs/1710.05421, 2017. URL <http://arxiv.org/abs/1710.05421>.
- Lee, Y., Sun, S.-H., Somasundaram, S., Hu, E., and Lim, J. J. Composing complex skills by learning transition policies with proximity reward induction. In *International Conference on Learning Representations*, 2019.

- Merel, J., Hasenclever, L., Galashov, A., Ahuja, A., Pham, V., Wayne, G., Teh, Y. W., and Heess, N. Neural probabilistic motor primitives for humanoid control. In *International Conference on Learning Representations*, 2019.
- Mnih, V., Kavukcuoglu, K., Silver, D., Rusu, A. A., Veness, J., Bellemare, M. G., Graves, A., Riedmiller, M., Fidjeland, A. K., Ostrovski, G., et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529, 2015.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. Asynchronous methods for deep reinforcement learning. In *International conference on machine learning*, pp. 1928–1937, 2016.
- Munos, R., Stepleton, T., Harutyunyan, A., and Bellemare, M. Safe and efficient off-policy reinforcement learning. In *Advances in Neural Information Processing Systems*, 2016.
- Nachum, O., Norouzi, M., Xu, K., and Schuurmans, D. Bridging the gap between value and policy based reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 2775–2785, 2017.
- Nachum, O., Gu, S. S., Lee, H., and Levine, S. Data-efficient hierarchical reinforcement learning. In *Advances in Neural Information Processing Systems 31*, pp. 3307–3317. 2018.
- Nachum, O., Gu, S., Lee, H., and Levine, S. Near-optimal representation learning for hierarchical reinforcement learning. In *International Conference on Learning Representations*, 2019.
- OpenAI. Openai five. <https://blog.openai.com/openai-five/>, 2018.
- OpenAI, Andrychowicz, M., Baker, B., Chociej, M., Jozefowicz, R., McGrew, B., Pachocki, J., Petron, A., Plappert, M., Powell, G., Ray, A., et al. Learning dexterous in-hand manipulation. *arXiv preprint arXiv:1808.00177*, 2018.
- Ortega, P. A. and Braun, D. A. Thermodynamics as a theory of decision-making with information-processing costs. *Proceedings of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 469(2153):20120683, 2013.
- Parisotto, E., Ba, J. L., and Salakhutdinov, R. Actor-mimic: Deep multitask and transfer reinforcement learning. In *International Conference on Learning Representations*, 2016.
- Parr, R. and Russell, S. J. Reinforcement learning with hierarchies of machines. In *Advances in neural information processing systems*, pp. 1043–1049, 1998.
- Rawlik, K., Toussaint, M., and Vijayakumar, S. On stochastic optimal control and reinforcement learning by approximate inference. In *Robotics: science and systems*, volume 13, pp. 3052–3056, 2012.
- Riedmiller, M., Hafner, R., Lampe, T., Neunert, M., Degrave, J., van de Wiele, T., Mnih, V., Heess, N., and Springenberg, J. T. Learning by playing solving sparse reward tasks from scratch. In *Proceedings of the 35th International Conference on Machine Learning*, pp. 4344–4353, 2018.
- Rubin, J., Shamir, O., and Tishby, N. Trading value and information in mdps. *Decision Making with Imperfect Decision Makers*, pp. 57–74, 2012.
- Rusu, A. A., Colmenarejo, S. G., Gulcehre, C., Desjardins, G., Kirkpatrick, J., Pascanu, R., Mnih, V., Kavukcuoglu, K., and Hadsell, R. Policy distillation. In *International Conference on Learning Representations*, 2016.
- Salimans, T., Kingma, D. P., and Welling, M. Markov Chain Monte Carlo and Variational Inference: Bridging the Gap. *ArXiv e-prints*, October 2014.
- Schmitt, S., Hudson, J. J., Zidek, A., Osindero, S., Doersch, C., Czarnecki, W. M., Leibo, J. Z., Kuttler, H., Zisserman, A., Simonyan, K., et al. Kickstarting deep reinforcement learning. *arXiv preprint arXiv:1803.03835*, 2018.
- Schulman, J., Levine, S., Moritz, P., Jordan, M., and Abbeel, P. Trust region policy optimization. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning-Volume 37*, pp. 1889–1897. JMLR. org, 2015.
- Schulman, J., Chen, X., and Abbeel, P. Equivalence between policy gradients and soft q-learning. *arXiv preprint arXiv:1704.06440*, 2017a.
- Schulman, J., Wolski, F., Dhariwal, P., Radford, A., and Klimov, O. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017b.
- Silver, D., Huang, A., Maddison, C. J., Guez, A., Sifre, L., Van Den Driessche, G., Schrittwieser, J., Antonoglou, I., Panneershelvam, V., Lanctot, M., et al. Mastering the game of go with deep neural networks and tree search. *nature*, 529(7587):484, 2016.
- Still, S. and Precup, D. An information-theoretic approach to curiosity-driven reinforcement learning. *Theory in Biosciences*, 131(3):139–148, 2012.

- Strouse, D., Kleiman-Weiner, M., Tenenbaum, J., Botvinick, M., and Schwab, D. J. Learning to share and hide intentions using information regularization. In *Advances in Neural Information Processing Systems*, pp. 10270–10281, 2018.
- Sutton, R. S., Precup, D., and Singh, S. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2): 181–211, 1999.
- Teh, Y., Bapst, V., Czarnecki, W. M., Quan, J., Kirkpatrick, J., Hadsell, R., Heess, N., and Pascanu, R. Distral: Robust multitask reinforcement learning. In *Advances in Neural Information Processing Systems*, pp. 4496–4506, 2017.
- Tiomkin, S. and Tishby, N. A unified bellman equation for causal information and value in markov decision processes. *arXiv preprint arXiv:1703.01585*, 2017.
- Tishby, N. and Polani, D. Information theory of decisions and actions. *Perception-Action Cycle*, pp. 601–636, 2011.
- Todorov, E. Linearly-solvable markov decision problems. In *Advances in Neural Information Processing Systems*, 2007.
- Toussaint, M. Robot trajectory optimization using approximate inference. In *Proceedings of the 26th annual international conference on machine learning*, pp. 1049–1056. ACM, 2009.
- Vezhnevets, A. S., Osindero, S., Schaul, T., Heess, N., Jaderberg, M., Silver, D., and Kavukcuoglu, K. Feudal networks for hierarchical reinforcement learning. In *International Conference on Machine Learning*, pp. 3540–3549, 2017.
- Wang, Z., Bapst, V., Heess, N., Mnih, V., Munos, R., Kavukcuoglu, K., and de Freitas, N. Sample efficient actor-critic with experience replay. In *International Conference on Learning Representations*, 2017.
- Xie, S., Galashov, A., Liu, S., Hou, S., Pascanu, R., Heess, N., and Teh, Y. W. Transferring task goals via hierarchical reinforcement learning, 2018. URL <https://openreview.net/forum?id=S1Y6TtJvG>.
- Ziebart, B. D. *Modeling Purposeful Adaptive Behavior with the Principle of Maximum Causal Entropy*. PhD thesis, Carnegie Mellon University, 2010.

Appendix

A. A general framework for RL as probabilistic modelling

In Sections 2 and 3 of the main text we have introduced the KL-regularized objective and explored a particular formulation that uses latent variables in the default policy and policy (Section 3 and experiments). The particular choice in Section 3 arises as a special case of a more general framework which we here outline briefly.

For both the default policy and for agent policy we can consider general directed latent variable models of the following form

$$\pi_0(\tau) = \int \pi_0(\tau|y)\pi_0(y)dy, \quad (11)$$

$$\pi(\tau) = \int \pi(\tau|z)\pi(z)dz \quad (12)$$

where both y and z can be time varying, e.g. $y = (y_1, \dots, y_T)$, and can be causally dependent on the trajectory prefix x_t , e.g. $y_t \sim p(\cdot|x_t)$ (and equivalently for z). The latent variables can further be continuous or discrete, and y_t or z_t can exhibit further structure (and thus include e.g. binary variables that model option termination). The general form of the objective presented in the main text

$$\mathcal{L}(\pi, \pi_0) = \mathbb{E}_\tau \left[\sum_{t \geq 1} \gamma^t r(s_t, a_t) - \alpha \gamma^t \text{KL}(a_t|x_t) \right],$$

remains valid regardless of the particular form of π_0 and π . This form can be convenient when $\pi_0(a_t|x_t)$ and $\pi(a_t|x_t)$ are tractable (e.g. when z or y have a small number of discrete states or decompose conveniently over time, e.g. as in (Fox et al., 2017; Krishnan et al., 2017)).

In general, however, latent variables in π_0 and π may introduce the need for additional approximations. In this case different models and algorithms can be instantiated based on a) the particular approximation chosen there, as well as b) choices for sharing of components between π_0 and π . A possible starting point when π_0 contains latent variables is the following lower bound to \mathcal{L} :

$$\mathcal{L} = \mathbb{E}_\pi \left[\sum_t r(s_t, a_t) \right] - \text{KL}[\pi(\tau)|\pi_0(\tau)] \quad (13)$$

$$\geq \mathbb{E}_\pi \left[\sum_t r(s_t, a_t) \right] + \mathbb{E}_f \left[\log \frac{\pi_0(\tau, y)}{f(y|\tau)} \right] + \text{H}[\pi(\tau)] \quad (14)$$

$$= \mathbb{E}_\pi \left[\sum_t r(s_t, a_t) \right] + \mathbb{E}_f \left[\log \pi_0(\tau|y) \right] - \text{KL}[f(y|\tau)|\pi_0(y)] + \text{H}[\pi(\tau)]. \quad (15)$$

If y_t are discrete and take on a small number of values we can compute $f(y|\tau)$ exactly (e.g. using the forward-backward algorithm as in (Fox et al., 2017; Krishnan et al.,

2017)); in other cases we can learn a parameterized approximation to the true posterior or can conceivably apply mixed inference schemes (e.g. Johnson et al., 2016).

Latent variables in the policy π can require an alternative approximation discussed e.g. in (Hausman et al., 2018):

$$\mathcal{L} \geq \mathbb{E}_\pi \left[\sum_t r(s_t, a_t) + \log \pi_0(\tau) + \log g(z|\tau) + \text{H}[\pi(\tau|z)] \right] + \text{H}[g(Z)], \quad (16)$$

where g is a learned approximation to the true posterior $\pi(z|\tau)$. (But see e.g. (Haarnoja et al., 2018a) who consider a parametric form for policies with latent variables for which the entropy term can be computed analytically and no approximation is needed.) This formulation bears interesting similarities with diversity inducing regularization schemes based on mutual information (e.g. Gregor et al., 2017; Florensa et al., 2017) but arises here as an approximation to trajectory entropy. This formulation also has interesting connections to auxiliary variable formulations in the approximate inference literature (Salimans et al., 2014; Agakov & Barber, 2004).

When both π_0 and π contain latent variables eqs. (15,16) can be combined. The model described in Section 3 in the main text then arises when the latent variable is “shared” between π_0 and π and we effectively use the policy itself as the inference network for π_0 : $f(y|\tau) = \prod_t \pi(y_t|x_t)$. In this case the objective simplifies to

$$\mathcal{L} \geq \mathbb{E}_\pi \left[\sum_t r(s_t, a_t) + \log \frac{\pi_0(\tau|z)\pi_0(z)}{\pi(\tau|z)\pi(z)} \right]. \quad (17)$$

When we further set $\pi_0(\tau|z) = \pi(\tau|z)$ we recover the model discussed in the main text of the paper.

As a proof-of-concept for a model without a shared latent space, with latent variables in π_0 but not π , we consider a simple humanoid with 28 degrees of freedom and 21 actuators and consider two different tasks: 1) a dense-reward walking task, in which the agent has to move forward, backward, left, or right at a fixed speed. The direction is randomly sampled at the beginning of an episode and changed to a different direction half-way through the episode and 2) a sparse reward go-to-target task, in which the agent has to move to a target whose location is supplied to the agent as a feature vector similar to those considered in (Galashov et al., 2019).

Figure 8 shows some exploratory results. In a first experiment we compare different prior architectures on the directional walking task. We let the prior marginalize over task

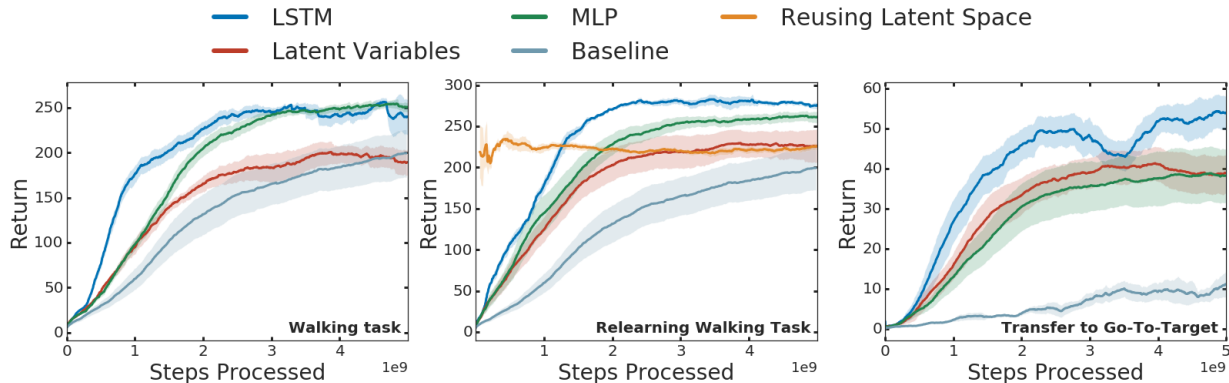


Figure 8. Results with a latent variable prior. Left: Walking task with the simple humanoid Center: Relearning the walking task with fixed priors. Right: Transfer to a go-to-target task.

condition. We include a feed-forward network, an LSTM, and a latent variable model with one latent variable per time step in the comparison. For the latent variable model we chose an inference network $f(z_t|z_{t-1}, \tau)$ so that eq. (15) decomposes over time. All priors studied in this comparison gave a modest speed-up in learning. While the latent variable prior works well, it does not work as well as the LSTM and MLP priors in this setup. In a first set of transfer experiments, we used the learned priors to learn the walking task again. Again, the learned priors led to a modest speed-up relative to learning from scratch.

We also experimented with parameter sharing for transfer as in the main text. We can freeze the conditional distribution $\pi_0(a|s, z)$ and learn a new policy $\pi(z|s)$, effectively using the learned latent space as an action space. In a second set of experiments, we study how well a prior learned on the walking task can transfer to the sparse go-to-target task. Here all learned priors led to a significant speed up relative to learning from scratch. Small return differences aside, all three different priors considered here solved the task with clear goal directed movements. On the other hand, the baseline only learned to go to very close-by targets. Reusing the latent space did not work well on this task. We speculate that the required turns are not easy to represent in the latent space resulting from the walking task.

B. Algorithm

We develop efficient off-policy learning algorithms for optimizing hierarchically structured model with default policies. Our off-policy learning algorithms are adapted from different existing algorithms based on the environments. Specifically, the algorithm for continuous control environments is based on SVG(0) (Heess et al., 2015) augmented with experience replay, and the algorithm for discrete action space environments is based on IMPALA (Espenholt et al., 2018),

which is a framework for off-policy actor critic algorithm. Note that the introduction of latent variables and default policies introduces additional challenges for off-policy learning algorithms in terms of gradient estimation of hierarchical policy, parameterization and estimation of value function, off-policy correction of the target value estimation, and appropriate incorporation of KL terms in the algorithm. We describe how we address such challenges in this section. Unless otherwise mentioned, we follow notations from the main paper.

B.1. Reparameterized latent for hierarchical policy

To estimate gradient for the hierarchical policy, we follow a strategy similar to (Heess et al., 2016) and reparameterize $z_t \sim \pi^H(z_t|x_t)$ as $z_t = f^H(x_t, \epsilon_t)$, where $\epsilon_t \sim \rho(\epsilon_t)$ is a fixed distribution. The $f^H(\cdot)$ is a deterministic function that outputs distribution parameters. In practice this means that the hierarchical policy can be treated as a flat policy $\pi(a_t|\epsilon_t, x_t) = \pi^L(a_t|f^H(x_t, \epsilon_t), x_t)$. We exploit the reparameterized flat policy to employ existing distributed learning algorithm with minimal modification.

B.2. Continuous control

In continuous control experiments, we employ distributed version of SVG(0) (Heess et al., 2015) augmented with experience replay and off-policy correction algorithm called Retrace (Munos et al., 2016). In the distributed setup, behaviour policies in multiple actors are used to collect off-policy trajectories and a single learner is used to optimize model parameters. The SVG(0) reparameterize a policy $p(a|s)$ and optimize it by backpropagating gradient from a learned action value function $Q(a, s)$ through a sampled action a .

To employ this algorithm, we reparameterize action from flat policy $a_t \sim \pi_\theta(a_t|\epsilon_t, x_t)$ with parameter θ as $a_t =$

$h_\theta(\epsilon_t, x_t, \xi_t)$, where $\xi_t \sim \rho(\xi_t)$ is a fixed distribution, and $h_\theta(\epsilon_t, x_t, \xi_t)$ is a deterministic function outputting a sample from the distribution $\pi_\theta(a_t|\epsilon_t, x_t)$. We also introduce the action value function $Q(a_t, z_t, x_t)$. Unlike policies without hierarchy, we estimate the action value depending on the sampled action z_t as well, so that it could capture the future returns depending on z_t . Given the flat policy and the action value function, SVG(0) (Heess et al., 2015) suggests to use following gradient estimate

$$\begin{aligned} & \nabla_\theta \mathbb{E}_{\pi_\theta(a|\epsilon_t, x_t)} Q(a, z_t, x_t) \\ &= \nabla_\theta \mathbb{E}_{\rho(\xi)} Q(h_\theta(\epsilon_t, x_t, \xi), z_t, x_t) \\ &= \mathbb{E}_{\rho(\xi)} \frac{\partial Q}{\partial h} \frac{\partial h}{\partial \theta} \approx \frac{1}{M} \sum_{i=1}^M \frac{\partial Q}{\partial h} \frac{\partial h}{\partial \theta} \Big|_{\xi=\xi_i}, \end{aligned} \quad (18)$$

which facilitates using backpropagation. Note that policy parameter θ could be learned through z_t as well, but we decide not to because it tends to make learning unstable.

To learn action value function $Q(a_t, z_t, x_t)$ and learn policy, we use off-policy trajectories from experience replay. We use Retrace (Munos et al., 2016) to estimate the action values from off-policy trajectories. The main idea behind Retrace is to use importance weighting to correct for the difference between the behavior policy μ and the online policy π , while cutting the importance weights to reduce variance. Specifically, we estimate corrected action value with

$$\hat{Q}_t^R = Q_t + \sum_{s \geq t} \gamma^{s-t} \left(\prod_{i=s}^t c_i \right) \delta_s Q, \quad (19)$$

where $\delta_s Q = r_s + \gamma(\hat{V}_{s+1} - \alpha \text{KL}_{s+1}) - Q_s$ and $Q_t = Q(a_t, z_t, x_t)$. $\hat{V}_s = \mathbb{E}_{\pi(a|\epsilon_t, x_t)} [Q(a, z_t, x_t)]$ is estimated bootstrap value, $\text{KL}_s = \text{KL}[\pi^H(z|x_s) \|\pi_0^H(z|x_s)] + \text{KL}[\pi^L(a|z_s, x_s) \|\pi_0^L(a|z_s, x_s)]$ and γ is discount. $c_i = \lambda \min\left(\frac{\pi(a_i|\epsilon_i, x_i)}{\mu(a_i|x_i)}\right)$ is truncated importance weight called *traces*.

There are, however, a few notable details that we adapt for our method. Firstly, we do not use the latent z_t sampled from behaviour policies in actors. This is possible because the latent does not affect the environment directly. Instead, we consider the behavior policy as $\mu(a|x)$, which does not depend on latents. This approach is useful since we do not need to consider the importance weight with respect to the HL policy, which might introduce additional variance in the estimator. Another detail is that the KL term at step s is not considered in $\delta_s Q$ because the KL at step s is not the result of action a_s . Instead, we introduce close form KL at step s as a loss to compensate for this. The pseudocode for the resulting algorithm is illustrated in Algorithm 2.

B.3. Discrete action space

For discrete control problems, we use distributed learning with V-trace (Espenholt et al., 2018) off-policy correction. Similarly to the distributed learning setup in continuous control, behaviours policies in multiple actors are used to collect trajectories and a single learner is used to optimize model parameters. The learning algorithm is almost identical to (Espenholt et al., 2018), but there are details that need to be considered mainly because of hierarchy with stochastic latent variable and temporal abstraction. Using negative KL as reward introduces another complication as well.

We consider optimizing objective with infrequent latent

$$\begin{aligned} \mathcal{L}(\pi, \pi_0) \geq & \\ & \mathbb{E}_\tau \left[\sum_{t \geq 1} \gamma^t r(s_t, a_t) - \alpha \gamma^t \mathbb{1}_p(t) \text{KL}(z_t|x_t) \right], \end{aligned} \quad (20)$$

where $\mathbb{1}_p(t)$ is the indicator function whose value is 1 if $t \bmod p \equiv 1$ with period p . This lower bound will be discussed later in Appendix C.2. This infrequent latent case is used for discrete action space experiment, by defining period p to be equal to the effective step size of the body.

We learn latent conditional value function $V_\psi(z_t, x_t)$ and reparameterized flat policy $\pi_\theta(a_t|\epsilon_t, x_t)$. V-trace target is computed as follows

$$v_s = V_\psi^s + \sum_{t \geq s} \gamma^{t-s} \left(\prod_{i=s}^{t-1} c_i \right) \delta_t V, \quad (21)$$

where $\delta_t V = \rho_t(r_t + \gamma(V_\psi^{t+1} - \text{KL}_{t+1}^p) - V_\psi^t)$, $\text{KL}_t^p = \mathbb{1}_p(t) \text{KL}[\pi_\theta^H(z|x_t) \|\pi_{0,\phi}^h(z|x_t)]$, and $V_\psi^t := V_\psi(z_t, x_t)$ is bootstrapped value at time step t . Importance weights are computed by $c_i := \min(\bar{c}, w_i)$ and $\rho_i := \min(\bar{\rho}, w_i)$, where $w_i = \frac{\pi_\theta(a_i|\epsilon_i, x_i)}{\mu(a_i|x_i)}$. \bar{c} and $\bar{\rho}$ are truncation coefficient identical to ones from original V-trace paper (Espenholt et al., 2018). Note that here we ignore latent sampled by behaviour policy and just consider states and actions from the trajectory. As discussed in Appendix B.2, we sample latent on-policy and this helps avoiding additional variance introduced with importance weight for HL policy.

Computed V-trace target is used for training both policy and value function with actor-critic algorithm. For training policy, we use policy gradient defined as

$$\sum_{t \geq 1} \rho_t \nabla_\theta \log \pi_\theta(a_t|z_t, x_t) \delta_t v \quad (22)$$

where $\delta_t v = \hat{r}_t + \gamma(v_{t+1} - \text{KL}_{t+1}^p) - V_\psi^t$ and $V_\psi^t = V_\psi(z_t, x_t)$. We optimize negative KL for time step t by adding an analytic loss function for HL policy $\pi_\theta^H(z|x_t)$ and default policy $\pi_\phi^H(z|x_t)$

$$\sum_{t \geq 1} \nabla_{\theta, \phi} \mathbb{1}_p(t) \text{KL}[\pi_\theta^H(z|x_t) \|\pi_{0,\phi}^H(z|x_t)] \quad (23)$$

For training value function, perform gradient descent over l_2 loss

$$\sum_{t \geq 1} (v_t - V_\psi^t) \nabla_\psi V_\psi(z_t, x_t). \quad (24)$$

Algorithm 2 SVG(0) (Heess et al., 2015) with experience replay for hierarchical policy

 Flat policy: $\pi_\theta(a_t|\epsilon_t, x_t)$ with parameter θ

 HL policy: $\pi_\theta^H(z_t|x_t)$, where latent is sampled by reparameterization $z_t = f_\theta^H(x_t, \epsilon_t)$

 Default policies: $\pi_{0,\phi}^H(z_t|x_t)$ and $\pi_{0,\phi}^L(a_t|z_t, x_t)$ with parameter ϕ

 Q-function: $Q_\psi(a_t, z_t, x_t)$ with parameter ψ

 Initialize target parameters $\theta' \leftarrow \theta$, $\phi' \leftarrow \phi$, $\psi' \leftarrow \psi$.

 Target update counter: $c \leftarrow 0$

 Target update period: P

 Replay buffer: \mathcal{B}
repeat
for $t = 0, K, 2K, \dots, T$ **do**

 Sample partial trajectory $\tau_{t:t+K}$ with action log likelihood $l_{t:t+K}$ from replay buffer \mathcal{B} :

$$\tau_{t:t+K} = (s_t, a_t, r_t, \dots, r_{t+K}),$$

$$l_{t:t+K} = (l_t, \dots, l_{t+K}) = (\log \mu(a_t|x_t), \dots, \log \mu(a_{t+K}|x_{t+K}))$$

 Sample latent: $\epsilon_{t'} \sim \rho(\epsilon)$, $z_{t'} = f_\theta^H(x_{t'}, \epsilon_{t'})$

 Compute KL: $\hat{\text{KL}}_{t'} = \text{KL} \left[\pi_\theta^H(z|x_{t'}) \parallel \pi_{0,\phi}^H(z|x_{t'}) \right] + \text{KL} \left[\pi_\theta^L(a|z_{t'}, x_{t'}) \parallel \pi_{0,\phi}^L(a|z_{t'}, x_{t'}) \right]$

Compute KL for Distillation:

$$\hat{\text{KL}}_{t'}^{\mathcal{D}} = \text{KL} \left[\pi_\theta^H(z|x_{t'}) \parallel \pi_{0,\phi}^H(z|x_{t'}) \right] + \text{KL} \left[\pi_\theta^L(a|z_{t'}, x_{t'}) \parallel \pi_{0,\phi}^L(a|z_{t'}, x_{t'}) \right]$$

 Compute action entropy: $\hat{\text{H}}_{t'} = \mathbb{E}_{\pi_\theta(a|\epsilon_{t'}, x_{t'})} [\log \pi_\theta(a|\epsilon_{t'}, x_{t'})]$

 Estimate bootstrap value: $\hat{V}_{t'} = \mathbb{E}_{\pi_\theta(a|\epsilon_{t'}, x_{t'})} [Q_{\psi'}(a, z_{t+K}, x_{t+K})] - \alpha \hat{\text{KL}}_{t+K}$

 Estimate traces (Munos et al., 2016): $\hat{c}_{t'} = \lambda \min \left(\frac{\pi_\theta(a_{t'}|\epsilon_{t'}, x_{t'})}{l_{t'}} \right)$

Apply Retrace to estimate Q targets (Munos et al., 2016):

$$\hat{Q}_{t'}^R = Q_{\psi'}(a_{t'}, z_{t'}, x_{t'}) + \sum_{s \geq t'} \gamma^{s-t'} \left(\prod_{i=s}^{t'} \hat{c}_i \right) \left(r_s + \gamma \left(\hat{V}_{s+1} - \alpha \hat{\text{KL}}_{s+1} \right) - Q_{\psi'}(a_s, z_s, x_s) \right)$$

 Policy loss: $\hat{L}_\pi = \sum_{i=t}^{t+K-1} \mathbb{E}_{\pi_\theta(a|\epsilon_i, x_i)} Q_{\psi'}(a, z_i, x_i) - \alpha \hat{\text{KL}}_i + \alpha_H \hat{\text{H}}_i$

 Q-value loss: $\hat{L}_Q = \sum_{i=t}^{t+K-1} \|\hat{Q}_i^R - Q_\psi(a, z_i, x_i)\|^2$

 Default policy loss: $\hat{L}_{\pi_0^H} = \sum_{i=t}^{t+K-1} \hat{\text{KL}}_i^{\mathcal{D}}$
 $\theta \leftarrow \theta + \beta_\pi \nabla_\theta \hat{L}_\pi$ $\phi \leftarrow \phi + \beta_{\pi_0^H} \nabla_\phi \hat{L}_{\pi_0^H}$ $\psi \leftarrow \psi - \beta_Q \nabla_\psi \hat{L}_Q$

 Increment counter $c \leftarrow c + 1$
if $c > P$ **then**

 Update target parameters $\theta' \leftarrow \theta$, $\phi' \leftarrow \phi$, $\psi' \leftarrow \psi$
 $c \leftarrow 0$
end if
end for
until

Additionally we include an action entropy bonus to encourage exploration (Mnih et al., 2016)

$$\sum_{t \geq 1} \nabla_\theta \text{H}[\pi_\theta(a|\epsilon_t, x_t)], \quad (25)$$

 where $\text{H}[\cdot]$ is close form entropy. We optimize the gradients from all four objectives jointly. Unlike continuous control, we do not maintain target parameters separately for the discrete action space experiments.

C. Derivations

This section includes derivations not described in the main paper.

C.1. Upper bound of KL divergence

The upper bound of KL divergence in eq. (6) of the main paper is derived as

$$\begin{aligned} \text{KL}(a_t|x_t) &\leq \text{KL}(a_t|x_t) + \mathbb{E}_{\pi(a_t|x_t)} [\text{KL}(z_t|a_t, x_t)] \\ &= \mathbb{E}_{\pi(a_t|x_t)} \left[\log \frac{\pi(a_t|x_t)}{\pi_0(a_t|x_t)} \right] \\ &\quad + \mathbb{E}_{\pi(a_t|x_t)} \left[\mathbb{E}_{\pi(z_t|a_t, x_t)} \left[\log \frac{\pi(z_t|a_t, x_t)}{\pi_0(z_t|a_t, x_t)} \right] \right] \\ &= \mathbb{E}_{\pi(a_t, z_t|x_t)} \left[\log \frac{\pi(a_t, z_t|x_t)}{\pi_0(a_t, z_t|x_t)} \right] = \text{KL}(z_t, a_t|x_t) \\ &= \mathbb{E}_{\pi(z_t|x_t)} \left[\log \frac{\pi(z_t|x_t)}{\pi_0(z_t|x_t)} \right] \\ &\quad + \mathbb{E}_{\pi(z_t|x_t)} \left[\mathbb{E}_{\pi(a_t|z_t, x_t)} \left[\log \frac{\pi(a_t|z_t, x_t)}{\pi_0(a_t|z_t, x_t)} \right] \right] \\ &= \text{KL}(z_t|x_t) + \mathbb{E}_{\pi(z_t|x_t)} [\text{KL}(a_t|z_t, x_t)] \quad (26) \end{aligned}$$

with the last expression being tractably approximated using Monte Carlo sampling. Note that:

$$\begin{aligned}\text{KL}(z_t|x_t) &= \text{KL}(\pi^H(z_t|x_t)||\pi_0^H(z_t|x_t)) \\ \text{KL}(a_t|z_t, x_t) &= \text{KL}(\pi^L(a_t|z_t, x_t)||\pi_0^L(a_t|z_t, x_t))\end{aligned}$$

C.2. Trajectory based derivation of lower bound

This section derives lower bound of eq. (1) in the main paper. We consider high level policy and default policy, whose latent are sampled every p steps. As in the lower bound in eq. (9) of the main paper, we consider a case where a latent is shared between the policy and the prior. We derive lower bound in trajectory level while considering both the period of high level action p and the discount γ . Here we introduce notation for trajectory including latent $\eta = (s_1, z_1, a_1, \dots)$, not including latent $\tau = (s_1, a_1, \dots)$, and including only the sequence of latent $\zeta = (z_1, z_{1+p}, \dots)$. In this section, we derive the following lower bound

$$\begin{aligned}\mathcal{L}(\pi, \pi_0) &= \mathbb{E}_\tau \left[\sum_{t \geq 1} \gamma^t r(s_t, a_t) - \alpha \gamma^t \text{KL}(a_t|x_t) \right] \\ &\geq \mathbb{E}_\eta \left[\sum_{t \geq 1} \gamma^t r(s_t, a_t) - \alpha \gamma^t \mathbb{1}_p(t) \text{KL}(z_t|x_t) \right. \\ &\quad \left. - \alpha \gamma^t \text{KL}(a_t|z_{p(t)}, x_t) \right],\end{aligned}\tag{27}$$

where $\mathbb{E}_\tau[\cdot]$ is taken with respect to the distribution over trajectories defined by the agent policy and system dynamics: $p(s_1) \prod_{t \geq 1} \pi(a_t|x_t) p(s_{t+1}|s_t, a_t)$. $\mathbb{E}_\eta[\cdot]$ is taken with respect to the distribution over trajectories of hierarchical policy including latent: $p(s_1) \prod_{t \geq 1} \pi(a_t|z_{p(t)}, x_t) \pi(z_t|x_t)^{\mathbb{1}_p(t)} p(s_{t+1}|s_t, a_t)$. $z_{p(t)}$ is the latest latent sample before time step t based on the period p and $\mathbb{1}_p(t)$ is the indicator function whose value is 1 if $t \bmod p \equiv 1$ with period p .

Note that these equations are composed of reward maximization and KL regularization term. We could show equality in reward maximization, and derive lower bound with respect to the KL regularization term. Therefore, we will present equality and inequality with respects to these two terms separately.

C.2.1. EQUALITY IN REWARD MAXIMIZATION

The equality of reward maximization term holds as follows

$$\begin{aligned}\mathbb{E}_{\pi_\tau} \left[\sum_{t \geq 1} \gamma^t r(s_t, a_t) \right] &= \int_\tau \pi(\tau) \left[\sum_{t \geq 1} \gamma^t r(s_t, a_t) \right] d\tau \\ &= \int_\tau \int_\xi \pi(\eta) d\xi \left[\sum_{t \geq 1} \gamma^t r(s_t, a_t) \right] d\tau \\ &= \int_\tau \int_\xi \pi(\eta) \left[\sum_{t \geq 1} \gamma^t r(s_t, a_t) \right] d\xi d\tau \\ &= \int_\eta \pi(\eta) \left[\sum_{t \geq 1} \gamma^t r(s_t, a_t) \right] d\eta \\ &= \mathbb{E}_{\pi_\eta} \left[\sum_{t \geq 1} \gamma^t r(s_t, a_t) \right].\end{aligned}\tag{28}$$

C.2.2. LOWER BOUND OF KL REGULARIZATION ON TRAJECTORY

We can show inequality in eq. (27) by deriving it from KL regularization term. We first derive lower bound of trajectory level KL regularization without considering discount.

$$\begin{aligned}\int_\tau \pi(\tau) \log \frac{\pi_0(\tau)}{\pi(\tau)} d\tau &= \int_\tau \pi(\tau) \left[\log \int_\xi \frac{\pi_0(\tau, \xi)}{\pi(\tau)} d\xi \right] d\tau \\ &= \int_\tau \pi(\tau) \left[\log \int_\xi \pi(\xi|\tau) \frac{\pi_0(\tau, \xi)}{\pi(\tau)\pi(\xi|\tau)} d\xi \right] d\tau \\ &\geq \int_\tau \pi(\tau) \left[\int_\xi \pi(\xi|\tau) \log \frac{\pi_0(\tau, \xi)}{\pi(\tau)\pi(\xi|\tau)} d\xi \right] d\tau \\ &= \int_\eta \pi(\eta) \log \frac{\pi_0(\eta)}{\pi(\eta)} d\eta,\end{aligned}\tag{29}$$

where the inequality holds by Jensen's inequality.

C.2.3. LOWER BOUND OF KL REGULARIZATION WITH DISCOUNT

To derive lower bound with discount, we first introduce the following equation to handle discount in the derivation.

$$\begin{aligned}\sum_{t \geq 1}^T \gamma^t a_t &= \sum_{t \geq 1}^{T-1} \left[(1-\gamma) \gamma^t \sum_{u \geq 1}^t a_u \right] \\ &\quad + \gamma^T \sum_{t \geq 1}^T a_t.\end{aligned}\tag{30}$$

This equality is useful to derive lower bound from summation without discounting ($\sum_{u \geq 1}^t a_u$) and then recombine it with the discounting terms.

To derive lower bound with respect to KL regularization with discount, we first rewrite KL regularization term as a trajectory based KL term. As α is a constant we will ignore it in the derivation, but it is straightforward to include it. We could turn KL regularization term as a trajectory based KL

$$\begin{aligned}\mathbb{E}_{\pi_\tau} \left[\sum_{t \geq 1}^T -\text{KL}(a_t|x_t) \right] &= - \int_\tau \pi(\tau) \left[\sum_{t \geq 1}^T \text{KL}(a_t|x_t) \right] d\tau \\ &= \int_\tau \pi(\tau) \left[\sum_{t \geq 1}^T \int_{a_t} \pi(a_t|x_t) \log \frac{\pi_0(a_t|x_t)}{\pi(a_t|x_t)} da_t \right] d\tau \\ &= \int_\tau \pi(\tau) \left[\sum_{t \geq 1}^T \log \frac{\pi_0(a_t|x_t)}{\pi(a_t|x_t)} \right] d\tau \\ &= \int_\tau \pi(\tau) \log \prod_{t \geq 1}^T \frac{\pi_0(a_t|x_t)}{\pi(a_t|x_t)} d\tau \\ &= \int_\tau \pi(\tau) \log \frac{\pi_0(\tau)}{\pi(\tau)} d\tau.\end{aligned}\tag{31}$$

We use eqs. (30) and (31) to rewrite discounted KL regular-

ization as a trajectory based equation.

$$\begin{aligned}
& \mathbb{E}_{\pi_\tau} \left[\sum_{t \geq 1}^T -\gamma^t \text{KL}(a_t | x_t) \right] \\
&= \mathbb{E}_{\pi_\tau} \left[-\sum_{t \geq 1}^{T-1} [(1-\gamma)\gamma^t \sum_{u \geq 1}^t \text{KL}(a_u | x_u)] \right. \\
&\quad \left. - \gamma^T \sum_{t \geq 1}^T \text{KL}(a_t | x_t) \right] \\
&= -\sum_{t \geq 1}^{T-1} [(1-\gamma)\gamma^t \mathbb{E}_{\pi_\tau} [\sum_{u \geq 1}^t \text{KL}(a_u | x_u)]] \\
&\quad - \gamma^T \mathbb{E}_{\pi_\tau} [\sum_{t \geq 1}^T \text{KL}(a_t | x_t)] \\
&= \sum_{t \geq 1}^{T-1} [(1-\gamma)\gamma^t \int_{\tau_t} \pi(\tau_t) \log \frac{\pi_0(\tau_t)}{\pi(\tau_t)} d\tau_t] \\
&\quad + \gamma^T \int_{\tau_T} \pi(\tau_T) \log \frac{\pi_0(\tau_T)}{\pi(\tau_T)} d\tau_T, \tag{32}
\end{aligned}$$

where τ_t is trajectory until time step t . We derive lower bound of this trajectory based representation using eq. (29).

$$\begin{aligned}
& \sum_{t \geq 1}^{T-1} [(1-\gamma)\gamma^t \int_{\tau_t} \pi(\tau_t) \log \frac{\pi_0(\tau_t)}{\pi(\tau_t)} d\tau_t] \\
&\quad + \gamma^T \int_{\tau_T} \pi(\tau_T) \log \frac{\pi_0(\tau_T)}{\pi(\tau_T)} d\tau_T \\
&\geq \sum_{t \geq 1}^{T-1} [(1-\gamma)\gamma^t \int_{\eta_t} \pi(\eta_t) \log \frac{\pi_0(\eta_t)}{\pi(\eta_t)} d\eta_t] \\
&\quad + \gamma^T \int_{\eta_T} \pi(\eta_T) \log \frac{\pi_0(\eta_T)}{\pi(\eta_T)} d\eta_T, \tag{33}
\end{aligned}$$

To rearrange this trajectory based representation with KL regularization formulation, we use following equality

$$\begin{aligned}
& \int_{\eta} \pi(\eta) \log \frac{\pi_0(\eta)}{\pi(\eta)} d\eta \\
&= \int_{\eta} \pi(\eta) \log \prod_{t \geq 1}^T \frac{\pi_0(a_t | z_{p(t)}, x_t) \pi_0(z_t | x_t)^{\mathbb{1}_{p(t)}}}{\pi(a_t | z_{p(t)}, x_t) \pi(z_t | x_t)^{\mathbb{1}_{p(t)}}} d\eta \\
&= \int_{\eta} \pi(\eta) \left[\sum_{t \geq 1}^T \log \frac{\pi_0(a_t | z_{p(t)}, x_t)}{\pi(a_t | z_{p(t)}, x_t)} \right. \\
&\quad \left. + \mathbb{1}_p(t) \log \frac{\pi_0(z_t | x_t)}{\pi(z_t | x_t)} \right] d\eta \\
&= \int_{\eta} \pi(\eta) \left[\sum_{t \geq 1}^T -\mathbb{1}_p(t) \text{KL}(z_t | x_t) \right. \\
&\quad \left. - \text{KL}(a_t | z_{p(t)}, x_t) \right] d\eta \\
&= \mathbb{E}_{\pi_\eta} \left[\sum_{t \geq 1}^T -\mathbb{1}_p(t) \text{KL}(z_t | x_t) - \text{KL}(a_t | z_{p(t)}, x_t) \right], \tag{34}
\end{aligned}$$

where $z_{p(t)}$ is the latest latent sample before time step t based on the period p .

We rearrange eq. (33) based on eqs. (30) and (34).

$$\begin{aligned}
& \sum_{t \geq 1}^{T-1} [(1-\gamma)\gamma^t \int_{\eta_t} \pi(\eta_t) \log \frac{\pi_0(\eta_t)}{\pi(\eta_t)} d\eta_t] \\
&\quad + \gamma^T \int_{\eta_T} \pi(\eta_T) \log \frac{\pi_0(\eta_T)}{\pi(\eta_T)} d\eta_T \\
&= -\sum_{t \geq 1}^{T-1} [(1-\gamma)\gamma^t \mathbb{E}_{\pi_\eta} [\sum_{u \geq 1}^t \mathbb{1}_p(u) \text{KL}(z_u, x_u) \\
&\quad + \text{KL}(a_u | z_{p(u)}, x_u)]] \\
&\quad - \gamma^T \mathbb{E}_{\pi_\eta} [\sum_{t \geq 1}^T \mathbb{1}_p(t) \text{KL}(z_t | x_t) + \text{KL}(a_t | z_{p(t)}, x_t)] \\
&= -\mathbb{E}_{\pi_\eta} [\sum_{t \geq 1}^{T-1} [(1-\gamma)\gamma^t \sum_{u \geq 1}^t \mathbb{1}_p(u) \text{KL}(z_u | x_u) \\
&\quad + \text{KL}(a_u | z_{p(u)}, x_u)]] \\
&\quad + \gamma^T \sum_{t \geq 1}^T \mathbb{1}_p(t) \text{KL}(z_t | x_t) + \text{KL}(a_t | z_{p(t)}, x_t)] \\
&= \mathbb{E}_{\pi_\eta} [\sum_{t \geq 1}^T -\mathbb{1}_p(t) \gamma^t \text{KL}(z_t | x_t) - \gamma^t \text{KL}(a_t | z_{p(t)}, x_t)]. \tag{35}
\end{aligned}$$

By combining all results, we obtain the following inequality.

$$\begin{aligned}
\mathcal{L}(\pi, \pi_0) &= \mathbb{E}_\tau [\sum_{t \geq 1} \gamma^t r(s_t, a_t) - \alpha \gamma^t \text{KL}(a_t | x_t)] \\
&\geq \mathbb{E}_\eta [\sum_{t \geq 1} \gamma^t r(s_t, a_t) - \alpha \gamma^t \mathbb{1}_p(t) \text{KL}(z_t | x_t) \\
&\quad - \alpha \gamma^t \text{KL}(a_t | z_{p(t)}, x_t)]. \tag{36}
\end{aligned}$$

D. Additional Experimental Results

D.1. Alternative training regimes

In the main paper, we present results based on learning speed with respect to the number of time steps processed by learner in distributed learning setup. Note that the number of time steps processed by the learner does not necessarily correspond to the number of collected trajectory time steps because of the use of experience replay, which allows to learning to proceed regardless of the amount of collected trajectories. We also experimented with two alternative training regimes to ensure that the speedup results reported are consistent. In Figure 9a, we compare the learning curves for our method against the SVG-0 and DISTRAL baselines in a quasi on-policy training regime similar to that of (Espeholt et al., 2018). In Figure 9b, we perform a similar comparison in the original replay based off-policy training regime but with a single actor generating the data. In both cases, our method learns faster than both baselines.

D.2. Information asymmetry in task transfer

Figure 10 illustrates the necessity of information asymmetry and KL regularization during transfer. Here we train the agent on the Move box Or Go to Target task with different information given to the LL and then transfer to Move Box and GTT. Therefore the distributions of the trajectories during training and those required for transfer should be similar.

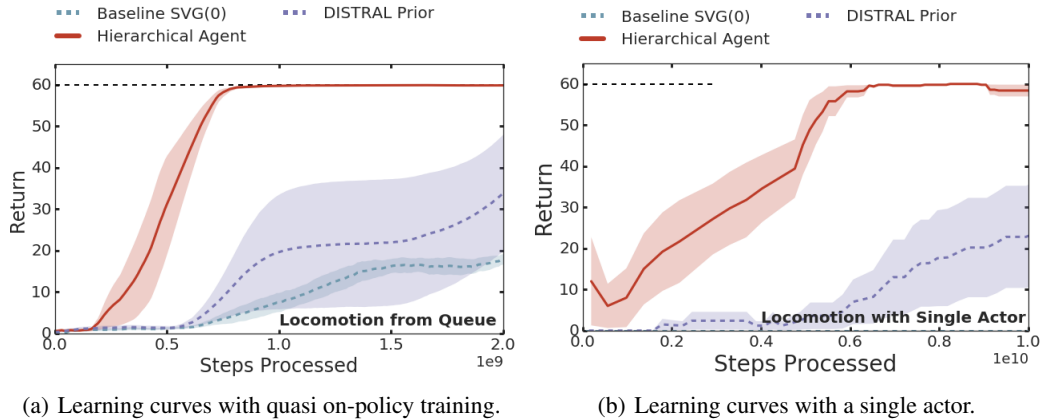


Figure 9. Experiments with alternative training regimes. (a) Locomotion with Ant. AR-1 process. (b) Locomotion with Ant. AR-1 process.

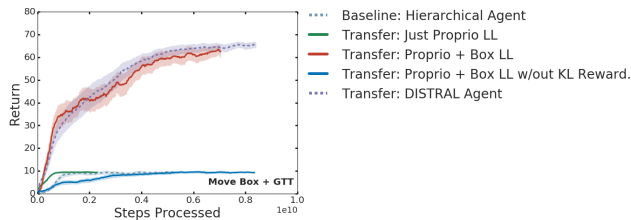


Figure 10. Analysis on task transfer results. Transfer from Combined task (OR) to Combined task (AND) with Ant and AR-1 process.

As the figure shows, we observe successful transfer only when the box position is given to the LL controller and KL regularization is used during transfer. Failed approaches usually converge to suboptimal policies, where the agent succeeds on the go to target task, but cannot move the box appropriately. In this transfer scenario, giving box position as input to LL controller is one way to specify inductive bias, which turns out to be useful to move box appropriately in target task. Interestingly, the unstructured DISTRAL prior performs comparably to our method in this experiment. We hypothesize that for tasks that take longer to learn, the benefit of the immediate parameter transfer in our approach is not as strong since this also leads to a fixed lower level behavior that cannot be adapted to the task. In this sense, the DISTRAL baseline is expected to be asymptotically stronger.

D.3. Additional body transfer experiments

We explore body transfer setup both in discrete and continuous environments. We compare performance to learning the hierarchical policy from scratch and analyze the effects of the KL regularization. The experimental setup in the continuous case is the same as before, and Figure 11 provides results for different types of bodies and tasks. Generally

transferring the HL component and relying on both the task reward and the KL term as a dense shaping reward signal for LL controller works best in these settings.

In the discrete case, we construct a **discrete go to target** task in a 2D grid world. An agent and goal are randomly placed in an 8×8 grid and the agent is rewarded for reaching the goal. The body agnostic task observation is the global x, y coordinates of the agent and goal. The different bodies in this case must take different numbers of actions to achieve an actual step in the grid. For instance the 4-step body needs to take 4 consecutive actions in the same direction to move forward by one step in that direction. We assume that the latent z_t is sampled every n steps, where n is the number of actions required to take a step. Details for models in which latent variables are sampled with a period > 1 are provided in Appendix B. The environment is described in Appendix E.

Figure 12 illustrates the result for transferring behavior from the 1-step to the 8-step body with AR-learned prior. (We were only able to solve the challenging 8-step version through body transfer with a KL reward.) In Figure 12, we visualize the negative KL divergence (KL reward) along the agent’s movement in every location of the grid. The size and the colour of arrows denotes the expected KL reward. This illustrates that the KL reward forms a vector field that guides the agent toward the goal, which provides a dense reward signal when transferring to a new body. This observation explains the gain from KL regularization, which can lead to faster learning and improve asymptotic performance.

E. Environments

E.1. Discrete control

We construct a **discrete go to target** task in 2D grid world. An agent and goal are randomly placed in an 8×8 grid

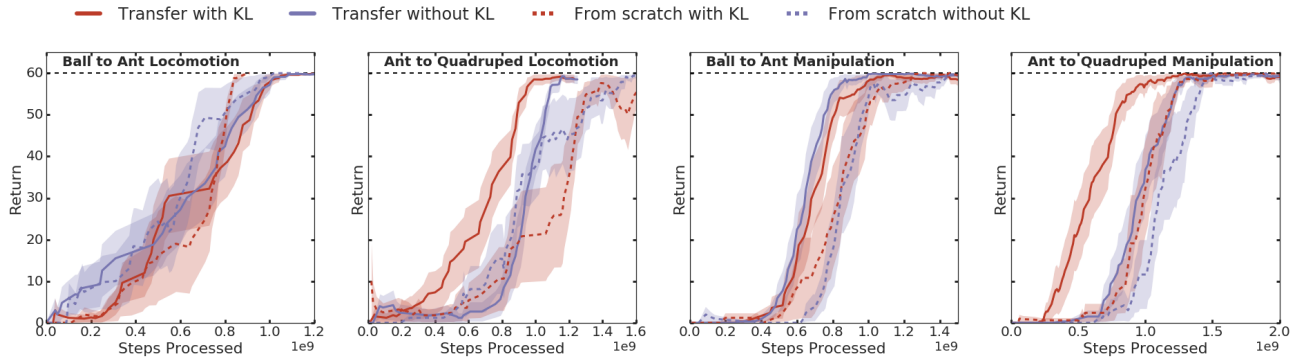


Figure 11. **Body transfer with the AR-1 Prior.** Column 1: Ball to Ant, Locomotion. Column 2: Ball to Ant, Manipulation (easy). Column 3: Ant to Quadruped, Locomotion. Column 4: Ant to Quadruped, Manipulation(easy).

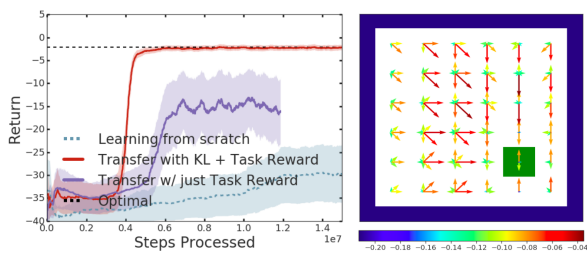


Figure 12. **Body transfer in 2D grid world, AR-1 prior.** Left Transfer from 1-step to 8-step body. Right KL reward visualization. The size and the color of arrows denotes KL reward (negative KL) for corresponding agents’ movement.

and the agent is given a reward of 1.0 for reaching the goal. The episode terminates when the agent reaches the goal or after 400 time steps. Additionally, the agent receives a penalty of -0.1 for every time step and a penalty of -0.2 if it collides with the walls. The body agnostic tasks observation is global x, y coordinate of agent and goal.

We consider a body that moves in any of 4 directions in the grid (up, down, left, right). We define different bodies based on their effective step size. The effective step size is the number of consecutive movements required to make a single displacement in the grid. Specifically, a body has 2 dimensional internal coordinate $[-n + 1, n - 1]^2$ with effective step size n . Agent’s action primarily affect the internal coordinate and it brings displacement to the external coordinate only if a value exceed its minimum or maximum. In this case, agent move 1 step in external coordinate and internal coordinate for the corresponding dimension is reset to 0. We denote different bodies with their step size (e.g. 1-step). We assume that the latent z_t is sampled every n steps, where n is the effective step size.

E.2. Continuous control

In this section, we describe detailed configuration of the continuous control tasks and bodies.

E.2.1. TASKS

Locomotion (Go to 1 of 3 targets) On a fixed 8×8 area, an agent and 3 targets are randomly placed at the beginning of episodes. In each episode, one of the 3 targets are randomly selected, and the agent should reach the selected target within 400 time steps. When the agent reach the selected target, the agent receives a reward of 60 and the episode terminates. Egocentric coordinates of 3 targets and an onehot vector representing one of the 3 targets are provided as task observations.

Locomotion (Gap) The task consists of a corridor with one gap in the middle. The length of the gap is chosen randomly at the start of each episode uniformly between 0.3 and 2.5. In order to successfully solve the task, the ant needs to jump across the gap and reach the end of the corridor. The ant receives a reward proportional to its velocity at each timestep. The observations given to the agent include proprioceptive information regarding the joint positions and velocities of the walker as well as the position and length of the gap for the current episode.

Manipulation (Move 1 of N boxes to 1 of K targets) On a fixed 3×3 area, an agent, N boxes and K targets are randomly placed at the beginning of episodes. In each episode, one of N boxes and one of the K targets are randomly selected, and the agent should move the box to the selected target within 400 time steps. When the box is placed on the selected target, the agent receives a reward of 60 and the episode terminates. Egocentric coordinates of K targets and 6 corner of the N boxes are provided as task observations with an onehot vector representing one of the K targets.

Manipulation (gather) On a fixed 3×3 area, an agent and 2 boxes are randomly locate at the beginning of episodes. In each episode, the agent should gather all the boxes within 400 time steps so that the boxes are contacted to each other. The agent receives a reward of 60 when it successfully gathers boxes. Egocentric coordinates of 6 corners of boxes are provided as task observations.

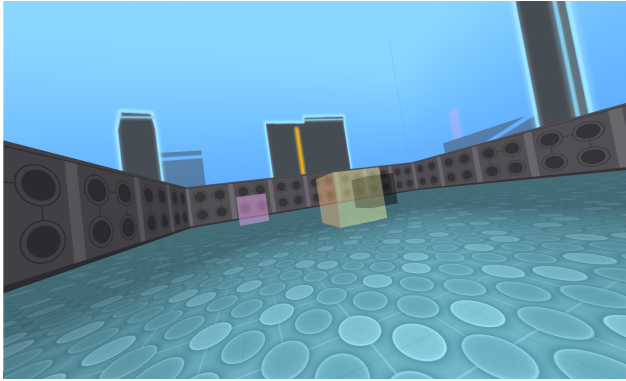


Figure 13. Vision input for Go to target from vision.

Combined task (Or) This task is a combination of the go to 1 of 2 targets task and the move box to 1 of 2 targets task. On each episode one of the two tasks is randomly sampled so that with probability 0.5 the agent must go to one of the 2 targets and with probability 0.5 it must push a box to a specific target. The agent receives a reward of 60 for successfully completing the corresponding task. Egocentric coordinates of both targets as well as the corners of the box are provided as task observations along with an encoding describing the specific task instance.

Combined task (And) This task contains the move box to 1 of 2 targets and go to 1 of 2 targets as two sub-tasks. In each episode, the walker must push the box to a target and then go to another target. A reward of 10 is awarded for each sub task that is completed and a bonus reward of 50 is awarded for completing both tasks. Egocentric coordinates of both targets as well as the corners of the box are provided as task observations.

Manipulation (vision) This task is identical to the Manipulation, but the task observation is egocentric vision. Instead of egocentric target coordinates, the agent observes a 64×64 image captured by the agent’s egocentric camera. The agent needs to recognize visual patterns of the targets and figure out the correct target (always colored black in this case). Figure 13 illustrates the egocentric visual observation before being rescaled to 64×64 .

E.2.2. BODIES

We use three different bodies: Ball, Ant, and Quadruped. Ball and Ant have been used in several previous works (Heess et al., 2017; Xie et al., 2018; Galashov et al., 2019), and we introduce Quadruped as a more complex variant of the Ant. The **Ball** is a body with 2 actuators for moving forward or backward, turning left, and turning right. The **Ant** is a body with 4 legs and 8 actuators, which moves its legs to walk and to interact with objects. The **Quadruped** is similar to Ant, but with 12 actuators. Each body has different proprioception (proprio) as a body

Task	Walker	LL information
Go to 1 of K Targets	Ant	Proprioception
Move box to target	Ant	Proprioception + Box
2 Boxes and 2 Targets	Ball	Proprioception + Boxes + Targets
Move Box or Go to Target (I)	Ant	Proprioception
Move Box or Go to Target (II)	Ant	Proprioception + Box
Gather Boxes	Ball	Proprioception + Boxes + Targets

Table 1. Information provided to the lower level controller for each task.

specific observation.

E.2.3. DETAILS OF INPUT TO THE LOWER LEVEL

1 illustrates the information provided to the lower level for each of the tasks for the speedup and transfer settings considered in the main text. In these cases, the HL received full information. For the tasks where the body was transferred, the LL was only given proprioceptive information while the HL received all other information relevant to the task.

F. Experimental Settings

F.1. General settings

Throughout the experiments, we use 32 actors to collect trajectories and a single learner to optimize the model. We plot average episode return with respect to the number of steps processed by the learner. Note that the number of steps is different from the number of agent’s interaction with environment, because the collected trajectories are processed multiple times by a centralized learner to update model parameters. When learning from scratch we report results as the mean and standard deviation of average returns for 5 random seeds. For the transfer learning experiments, we use 5 seed for the initial training, and then transfer all pretrained models to a new task, and train two new HL or LL controllers (two random seeds) per model on the transfer task. Thus, in total, 10 different runs are used to estimate mean and standard deviations of the average returns. Hyperparameters, including KL cost and action entropy regularization cost, are optimized on a per-task basis. More details are provided in Appendix G.

F.2. DISTRAL with parameter sharing

We introduced two baselines as variants of DISTRAL prior sharing parameters between the agent policy and the default

policy. **DISTRAL prior 2 cols** uses a 2-column architecture (see (Teh et al., 2017)), where the default policy network is reused in combination with another network (column) to output the final policy distribution. In this configuration, the default policy network does not access to task information, but another network (column) access to the full information. In **DISTRAL shared prior**, the policy network is reused to output the default policy distribution based on an additional branch on top of it. As default policy is constructed on policy network, information asymmetry is not used in this baseline.

G. Hyper parameters

Fully connected neural networks are used as function approximators for both the actor and the critic in the agent. In case of tasks with boxes, a separate common single layer MLP was used as a *box encoder torso*. Separate processing networks were implemented for *proprioception* for the baselines. For all tasks, multiple values were swept for actor networks and torso sizes with 5 random seeds each. ELU activations were used throughout. We use separate optimizers and learning rates for the actor and critic networks. For the hierarchical networks, fully connected MLP networks were used for the higher level and lower level policy cores. The relative contribution of the KL regularization to the reward was controlled by a *posterior entropy cost* which we denote α .

Below we provide the default hyperparameters used across tasks followed by the best parameters from the baselines and the hierarchical networks for each task.

G.1. Default parameters

Actor learning rate, β_{pi} = $1e-4$.

Critic learning rate, β_{pi} = $1e-4$.

DISTRAL baseline default policy learning rate, β_{pi} = $5e-4$.

Target network update period = 100.

DISTRAL policy target network update period = 100.

Baseline Actor network: MLP with sizes (200, 100).

Baseline Critic network: MLP with sizes (400, 300).

DISTRAL baseline default policy network: MLP with sizes (200, 100).

HL policy network: MLP with sizes (200, 10).

LL policy network: MLP with sizes (200, 100).

Box encoder network: MLP with sizes (50).

Batch size: 512.

Unroll length: 10.

Entropy bonus, λ = $1e-4$.

Posterior Entropy cost for HL, α = $1e-3$.

Posterior Entropy cost for DISTRAL default policy, α = 0.01

Distillation cost for DISTRAL default policy, α = 0.01

Number of actors: 32

G.2. Per-task parameters

Ant: Move Box to Target

Entropy bonus, λ = $1e-3$.

Policy network for Gaussian Prior: MLP with sizes (200, 100).

Action Entropy cost for Isotropic Gaussian Prior: = 0.

Action Entropy cost for AR-1 Prior: = 0.

AR Parameter: = 0.9

Posterior Entropy cost for AR-1 Prior, α = $1e-4$.

Ant: Go to 1 of 3 Targets

Action Entropy cost for DISTRAL default policy = 0.

AR Parameter: = 0.95

Distillation cost for DISTRAL default policy, α = 1.0

Ball: 2 Boxes to 2 Targets

Distillation cost for DISTRAL default policy, α = 0.1

Box encoder network: MLP with sizes (100, 20).

HL policy network for AR-Learned Prior: MLP with sizes (200, 4).

Ant: Move Box to 1 of 3 Targets

AR Parameter for baseline: = 0.9

Policy network for AR-1 Transfer Prior: MLP with sizes (200, 10).

Policy network for Gaussian Transfer Prior: MLP with sizes (200, 4).

Policy network for AR-Learned Transfer Prior: MLP with sizes (200, 4).

Posterior Entropy cost for AR-Learned Transfer Prior, α = $1e-2$.

Action Entropy cost for AR-Learned Transfer Prior: = $1e-4$.

Ball: Gather boxes

Policy network for Gaussian Transfer Prior: MLP with sizes (200, 4).

Policy network for AR-Learned Transfer Prior: MLP with

sizes (200, 4).

DISTRAL Actor learning rate, $\beta_{pi} = 5e-4$.

DISTRAL default policy distillation cost = 0.01.

Ant: Move Box and Go to Target

Policy network for AR-1 Prior: MLP with sizes (200, 10).

Ball to Ant: Go to Target

Unless otherwise specified, for all the body transfer tasks an *action entropy cost* of $1e-4$ worked best across tasks.

HL Policy network for agent from scratch: MLP with sizes (100, 4)

Ball to Ant: Move Box to Target

HL Policy network for agent from scratch: MLP with sizes (100, 4)

Posterior entropy cost for transfer agent: $1e-5$

Ant to Quadruped: Move Box to Target

Posterior entropy cost for transfer agent: $1e-2$

Action entropy cost for transfer agent: $1e-5$

Ant to Quadruped: Go To Target:

HL Policy network for agent from scratch: MLP with sizes (20)

Ant to Quadruped: Go To Target From Vision:

HL Policy network: Residual Network with an embedding size of (256,)