
DARLA: Improving Zero-Shot Transfer in Reinforcement Learning

Irina Higgins^{*1} Arka Pal^{*1} Andrei Rusu¹ Loic Matthey¹ Christopher Burgess¹ Alexander Pritzel¹
 Matthew Botvinick¹ Charles Blundell¹ Alexander Lerchner¹

Abstract

Domain adaptation is an important open problem in deep reinforcement learning (RL). In many scenarios of interest data is hard to obtain, so agents may learn a *source* policy in a setting where data is readily available, with the hope that it generalises well to the *target* domain. We propose a new multi-stage RL agent, DARLA (Disentangled Representation Learning Agent), which learns to see before learning to act. DARLA’s vision is based on learning a disentangled representation of the observed environment. Once DARLA can see, it is able to acquire source policies that are robust to many domain shifts - even with no access to the target domain. DARLA significantly outperforms conventional baselines in zero-shot domain adaptation scenarios, an effect that holds across a variety of RL environments (Jaco arm, DeepMind Lab) and base RL algorithms (DQN, A3C and EC).

1. Introduction

Autonomous agents can learn how to maximise future expected rewards by choosing how to act based on incoming sensory observations via reinforcement learning (RL). Early RL approaches did not scale well to environments with large state spaces and high-dimensional raw observations (Sutton & Barto, 1998). A commonly used workaround was to embed the observations in a lower-dimensional space, typically via hand-crafted and/or privileged-information features. Recently, the advent of deep learning and its successful combination with RL has enabled end-to-end learning of such embeddings directly from raw inputs, sparking success in a wide variety of previously challenging RL domains (Mnih et al., 2015; 2016; Jaderberg et al., 2017). Despite the seemingly universal

efficacy of deep RL, however, fundamental issues remain. These include data inefficiency, the reactive nature and general brittleness of learnt policies to changes in input data distribution, and lack of model interpretability (Garnelo et al., 2016; Lake et al., 2016). This paper focuses on one of these outstanding issues: the ability of RL agents to deal with changes to the input distribution, a form of transfer learning known as *domain adaptation* (Bengio et al., 2013). In domain adaptation scenarios, an agent trained on a particular input distribution with a specified reward structure (termed the *source domain*) is placed in a setting where the input distribution is modified but the reward structure remains largely intact (the *target domain*). We aim to develop an agent that can learn a robust policy using observations and rewards obtained exclusively within the source domain. Here, a policy is considered as robust if it generalises with minimal drop in performance to the target domain without extra fine-tuning.

Past attempts to build RL agents with strong domain adaptation performance highlighted the importance of learning good internal representations of raw observations (Finn et al., 2015; Raffin et al., 2017; Pan & Yang, 2009; Barreto et al., 2016; Littman et al., 2001). Typically, these approaches tried to align the source and target domain representations by utilising observation and reward signals from both domains (Tzeng et al., 2016; Daftry et al., 2016; Parisotto et al., 2015; Guez et al., 2012; Talvitie & Singh, 2007; Niekum et al., 2013; Gupta et al., 2017; Finn et al., 2017; Rajendran et al., 2017). In many scenarios, such as robotics, this reliance on target domain information can be problematic, as the data may be expensive or difficult to obtain (Finn et al., 2017; Rusu et al., 2016). Furthermore, the target domain may simply not be known in advance. On the other hand, policies learnt exclusively on the source domain using existing deep RL approaches that have few constraints on the nature of the learnt representations often overfit to the source input distribution, resulting in poor domain adaptation performance (Lake et al., 2016; Rusu et al., 2016).

We propose tackling both of these issues by focusing instead on learning representations which capture an underlying low-dimensional factorised representation of the world and are therefore not task or domain specific. Many nat-

^{*}Equal contribution ¹DeepMind, 6 Pancras Square, Kings Cross, London, N1C 4AG, UK. Correspondence to: Irina Higgins <irinah@google.com>, Arka Pal <arkap@google.com>.

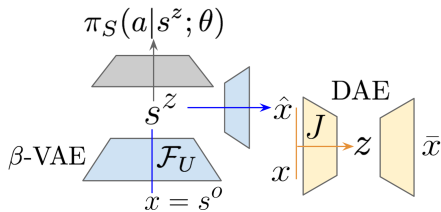


Figure 1. Schematic representation of DARLA. Yellow represents the denoising autoencoder part of the model, blue represents the β -VAE part of the model, and grey represents the policy learning part of the model.

uralistic domains such as video game environments, simulations and our own world are well described in terms of such a structure. Examples of such factors of variation are object properties like colour, scale, or position; other examples correspond to general environmental factors, such as geometry and lighting. We think of these factors as a set of high-level parameters that can be used by a world graphics engine to generate a particular natural visual scene (Kulkarni et al., 2015). Learning how to project raw observations into such a factorised description of the world is addressed by the large body of literature on disentangled representation learning (Schmidhuber, 1992; Desjardins et al., 2012; Cohen & Welling, 2014; 2015; Kulkarni et al., 2015; Hinton et al., 2011; Rippel & Adams, 2013; Reed et al., 2014; Yang et al., 2015; Goroshin et al., 2015; Kulkarni et al., 2015; Cheung et al., 2015; Whitney et al., 2016; Karaletos et al., 2016; Chen et al., 2016; Higgins et al., 2017). Disentangled representations are defined as interpretable, factorised latent representations where either a single latent or a group of latent units are sensitive to changes in single ground truth factors of variation used to generate the visual world, while being invariant to changes in other factors (Bengio et al., 2013). The theoretical utility of disentangled representations for supervised and reinforcement learning has been described before (Bengio et al., 2013; Higgins et al., 2017; Ridgeway, 2016); however, to our knowledge, it has not been empirically validated to date.

We demonstrate how disentangled representations can improve the robustness of RL algorithms in domain adaptation scenarios by introducing *DARLA* (Disentangled Representation Learning Agent), a new RL agent capable of learning a robust policy on the source domain that achieves significantly better out-of-the-box performance in domain adaptation scenarios compared to various baselines. *DARLA* relies on learning a latent state representation that is shared between the source and target domains, by learning a disentangled representation of the environment’s generative factors. Crucially, *DARLA* does not require target domain data to form its representations. Our approach utilises a three stage pipeline: 1) learning to see, 2) learning to act, 3) transfer. During the first stage,

DARLA develops its vision, learning to parse the world in terms of basic visual concepts, such as objects, positions, colours, etc. by utilising a stream of raw unlabelled observations – not unlike human babies in their first few months of life (Leat et al., 2009; Candy et al., 2009). In the second stage, the agent utilises this disentangled visual representation to learn a robust source policy. In stage three, we demonstrate that the *DARLA* source policy is more robust to domain shifts, leading to a significantly smaller drop in performance in the target domain even when no further policy finetuning is allowed (median 270.3% improvement). These effects hold consistently across a number of different RL environments (DeepMind Lab and Jaco/MuJoCo: Beattie et al., 2016; Todorov et al., 2012) and algorithms (DQN, A3C and Episodic Control: Mnih et al., 2015; 2016; Blundell et al., 2016).

2. Framework

2.1. Domain adaptation in Reinforcement Learning

We now formalise domain adaptation scenarios in a reinforcement learning (RL) setting. We denote the source and target domains as D_S and D_T , respectively. Each domain corresponds to an MDP defined as a tuple $D_S \equiv (\mathcal{S}_S, \mathcal{A}_S, \mathcal{T}_S, R_S)$ or $D_T \equiv (\mathcal{S}_T, \mathcal{A}_T, \mathcal{T}_T, R_T)$ (we assume a shared fixed discount factor γ), each with its own state space \mathcal{S} , action space \mathcal{A} , transition function \mathcal{T} and reward function R .¹ In domain adaptation scenarios the states \mathcal{S} of the source and the target domains can be quite different, while the action spaces \mathcal{A} are shared and the transitions \mathcal{T} and reward functions R have structural similarity. For example, consider a domain adaptation scenario for the Jaco robotic arm, where the MuJoCo (Todorov et al., 2012) simulation of the arm is the source domain, and the real world setting is the target domain. The state spaces (raw pixels) of the source and the target domains differ significantly due to the perceptual-reality gap (Rusu et al., 2016); that is to say, $\mathcal{S}_S \neq \mathcal{S}_T$. Both domains, however, share action spaces ($\mathcal{A}_S = \mathcal{A}_T$), since the policy learns to control the same set of actuators within the arm. Finally, the source and target domain transition and reward functions share structural similarity ($\mathcal{T}_S \approx \mathcal{T}_T$ and $R_S \approx R_T$), since in both domains transitions between states are governed by the physics of the world and the performance on the task depends on the relative position of the arm’s end effectors (i.e. fingertips) with respect to an object of interest.

2.2. DARLA

In order to describe our proposed *DARLA* framework, we assume that there exists a set \mathcal{M} of MDPs that is the set

¹For further background on the notation relating to the RL paradigm, see Section A.1 in the Supplementary Materials.

of all *natural world* MDPs, and each MDP D_i is sampled from \mathcal{M} . We define \mathcal{M} in terms of the state space $\hat{\mathcal{S}}$ that contains all possible conjunctions of high-level factors of variation necessary to generate any naturalistic observation in any $D_i \in \mathcal{M}$. A *natural world* MDP D_i is then one whose state space \mathcal{S} corresponds to some subset of $\hat{\mathcal{S}}$. In simple terms, we assume that there exists *some* shared underlying structure between the MDPs D_i sampled from \mathcal{M} . We contend that this is a reasonable assumption that permits inclusion of many interesting problems, including being able to characterise our own reality (Lake et al., 2016).

We now introduce notation for two state space variables that may in principle be used interchangeably within the source and target domain MDPs D_S and D_T – the agent observation state space \mathcal{S}^o , and the agent’s internal latent state space \mathcal{S}^z .² \mathcal{S}_i^o in D_i consists of raw (pixel) observations s_i^o generated by the true world simulator from a sampled set of data generative factors \hat{s}_i , i.e. $s_i^o \sim \mathbf{Sim}(\hat{s}_i)$. \hat{s}_i is sampled by some distribution or process \mathcal{G}_i on $\hat{\mathcal{S}}$, $\hat{s}_i \sim \mathcal{G}_i(\hat{\mathcal{S}})$.

Using the newly introduced notation, domain adaptation scenarios can be described as having different sampling processes \mathcal{G}_S and \mathcal{G}_T such that $\hat{s}_S \sim \mathcal{G}_S(\hat{\mathcal{S}})$ and $\hat{s}_T \sim \mathcal{G}_T(\hat{\mathcal{S}})$ for the source and target domains respectively, and then using these to generate different agent observation states $s_S^o \sim \mathbf{Sim}(\hat{s}_S)$ and $s_T^o \sim \mathbf{Sim}(\hat{s}_T)$. Intuitively, consider a source domain where oranges appear in blue rooms and apples appear in red rooms, and a target domain where the object/room conjunctions are reversed and oranges appear in red rooms and apples appear in blue rooms. While the true data generative factors of variation $\hat{\mathcal{S}}$ remain the same – room colour (blue or red) and object type (apples and oranges) – the particular source and target distributions \mathcal{G}_S and \mathcal{G}_T differ.

Typically deep RL agents (e.g. Mnih et al., 2015; 2016) operating in an MDP $D_i \in \mathcal{M}$ learn an end-to-end mapping from raw (pixel) observations $s_i^o \in \mathcal{S}_i^o$ to actions $a_i \in \mathcal{A}_i$ (either directly or via a value function $Q_i(s_i^o, a_i)$ from which actions can be derived). In the process of doing so, the agent implicitly learns a function $\mathcal{F} : \mathcal{S}_i^o \rightarrow \mathcal{S}_i^z$ that maps the typically high-dimensional raw observations s_i^o to typically low-dimensional latent states s_i^z ; followed by a policy function $\pi_i : \mathcal{S}_i^z \rightarrow \mathcal{A}_i$ that maps the latent states s_i^z to actions $a_i \in \mathcal{A}_i$. In the context of domain adaptation, if the agent learns a naive latent state mapping function $\mathcal{F}_S : \mathcal{S}_S^o \rightarrow \mathcal{S}_S^z$ on the source domain using reward signals to shape the representation learning, it is likely that \mathcal{F}_S will overfit to the source domain and will not generalise well to the target domain. Returning to our

²Note that we do not assume these to be Markovian i.e. it is not necessarily the case that $p(s_{t+1}^o | s_t^o) = p(s_{t+1}^o | s_t^o, s_{t-1}^o, \dots, s_1^o)$, and similarly for s^z . Note the index t here corresponds to time.

intuitive example, imagine an agent that has learnt a policy to pick up oranges and avoid apples on the source domain. Such a source policy π_S is likely to be based on an *entangled* latent state space \mathcal{S}_S^z of object/room conjunctions: oranges/blue \rightarrow good, apples/red \rightarrow bad, since this is arguably the most efficient representation for maximising expected rewards on the source task in the absence of extra supervision signals suggesting otherwise. A source policy $\pi_S(a | s_S^z; \theta)$ based on such an entangled latent representation s_S^z will not generalise well to the target domain without further fine-tuning, since $\mathcal{F}_S(s_S^o) \neq \mathcal{F}_S(s_T^o)$ and therefore crucially $\mathcal{S}_S^z \neq \mathcal{S}_T^z$.

On the other hand, since both $\hat{s}_S \sim \mathcal{G}_S(\hat{\mathcal{S}})$ and $\hat{s}_T \sim \mathcal{G}_T(\hat{\mathcal{S}})$ are sampled from the same *natural world* state space $\hat{\mathcal{S}}$ for the source and target domains respectively, it should be possible to learn a latent state mapping function $\hat{\mathcal{F}} : \mathcal{S}^o \rightarrow \mathcal{S}^z$, which projects the agent observation state space \mathcal{S}^o to a latent state space \mathcal{S}^z expressed in terms of factorised data generative factors that are representative of the natural world i.e. $\mathcal{S}^z \approx \hat{\mathcal{S}}$. Consider again our intuitive example, where $\hat{\mathcal{F}}$ maps agent observations (s_S^o : orange in a blue room) to a factorised or *disentangled* representation expressed in terms of the data generative factors (s_S^z : room type = blue; object type = orange). Such a *disentangled* latent state mapping function should then directly generalise to both the source and the target domains, so that $\hat{\mathcal{F}}(s_S^o) = \hat{\mathcal{F}}(s_T^o) = s^z$. Since \mathcal{S}^z is a *disentangled* representation of object and room attributes, the source policy π_S can learn a decision boundary that ignores the irrelevant room attributes: oranges \rightarrow good, apples \rightarrow bad. Such a policy would then generalise well to the target domain out of the box, since $\pi_S(a | \hat{\mathcal{F}}(s_S^o); \theta) = \pi_T(a | \hat{\mathcal{F}}(s_T^o); \theta) = \pi_T(a | s^z; \theta)$. Hence, DARLA is based on the idea that a good quality $\hat{\mathcal{F}}$ learnt exclusively on the source domain $D_S \in \mathcal{M}$ will zero-shot-generalise to all target domains $D_i \in \mathcal{M}$, and therefore the source policy $\pi(a | \mathcal{S}_S^z; \theta)$ will also generalise to all target domains $D_i \in \mathcal{M}$ out of the box.

Next we describe each of the stages of the DARLA pipeline that allow it to learn source policies π_S that are robust to domain adaptation scenarios, despite being trained with no knowledge of the target domains (see Fig. 1 for a graphical representation of these steps):

1) *Learn to see* (unsupervised learning of \mathcal{F}_U) – the task of inferring a factorised set of generative factors $\mathcal{S}_S^z = \hat{\mathcal{S}}$ from observations \mathcal{S}^o is the goal of the extensive disentangled factor learning literature (e.g. Chen et al., 2016; Higgins et al., 2017). Hence, in stage one we learn a mapping $\mathcal{F}_U : \mathcal{S}_U^o \rightarrow \mathcal{S}_U^z$, where $\mathcal{S}_U^z \approx \mathcal{S}_S^z$ (U stands for ‘unsupervised’) using an unsupervised model for learning disentangled factors that utilises observations collected by an agent with a random policy π_U from a visual pre-training

MDP $D_U \in \mathcal{M}$. Note that we require sufficient variability of factors and their conjunctions in D_U in order to have $S_U^z \approx S_S^z$;

2) *Learn to act* (reinforcement learning of π_S in the source domain D_S utilising previously learned \mathcal{F}_U) – an agent that has learnt to see the world in stage one in terms of the natural data generative factors is now exposed to a source domain $D_S \in \mathcal{M}$. The agent is tasked with learning the source policy $\pi_S(a|s_S^z; \theta)$, where $s_S^z = \mathcal{F}_U(s_S^o) \approx s_S^z$, via a standard reinforcement learning algorithm. Crucially, we do not allow \mathcal{F}_U to be modified (e.g. by gradient updates) during this phase;

3) *Transfer* (to a target domain D_T) – in the final step, we test how well the policy π_S learnt on the source domain generalises to the target domain $D_T \in \mathcal{M}$ in a zero-shot domain adaptation setting, i.e. the agent is evaluated on the target domain without retraining. We compare the performance of policies learnt with a disentangled latent state S_S^z to various baselines where the latent state mapping function \mathcal{F}_U projects agent observations s^o to entangled latent state representations s^z .

2.3. Learning disentangled representations

In order to learn \mathcal{F}_U , DARLA utilises β -VAE (Higgins et al., 2017), a state-of-the-art unsupervised model for automated discovery of factorised latent representations from raw image data. β -VAE is a modification of the variational autoencoder framework (Kingma & Welling, 2014; Rezende et al., 2014) that controls the nature of the learnt latent representations by introducing an adjustable hyperparameter β to balance reconstruction accuracy with latent channel capacity and independence constraints. It maximises the objective:

$$\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}, \beta) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \beta D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \quad (1)$$

where ϕ, θ parametrise the distributions of the encoder and the decoder respectively. Well-chosen values of β - usually larger than one ($\beta > 1$) - typically result in more disentangled latent representations \mathbf{z} by limiting the capacity of the latent information channel, and hence encouraging a more efficient factorised encoding through the increased pressure to match the isotropic unit Gaussian prior $p(\mathbf{z})$ (Higgins et al., 2017).

2.3.1. PERCEPTUAL SIMILARITY LOSS

The cost of increasing β is that crucial information about the scene may be discarded in the latent representation \mathbf{z} , particularly if that information takes up a small proportion of the observations \mathbf{x} in pixel space. We encountered this issue in some of our tasks, as discussed in Section 3.1. The shortcomings of calculating the log-likelihood term

$\mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})]$ on a per-pixel basis are known and have been addressed in the past by calculating the reconstruction cost in an abstract, high-level feature space given by *another* neural network model, such as a GAN (Goodfellow et al., 2014) or a pre-trained AlexNet (Krizhevsky et al., 2012; Larsen et al., 2016; Dosovitskiy & Brox, 2016; Warde-Farley & Bengio, 2017). In practice we found that pre-training a denoising autoencoder (Vincent et al., 2010) on data from the visual pre-training MDP $D_U \in \mathcal{M}$ worked best as the reconstruction targets for β -VAE to match (see Fig. 1 for model architecture and Sec. A.3.1 in Supplementary Materials for implementation details). The new β -VAE_{DAE} model was trained according to Eq. 2:

$$\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}, \beta) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \|J(\hat{\mathbf{x}}) - J(\mathbf{x})\|_2^2 - \beta D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z})) \quad (2)$$

where $\hat{\mathbf{x}} \sim p_\theta(\mathbf{x}|\mathbf{z})$ and $J : \mathbb{R}^{W \times H \times C} \rightarrow \mathbb{R}^N$ is the function that maps images from pixel space with dimensionality $W \times H \times C$ to a high-level feature space with dimensionality N given by a stack of pre-trained DAE layers up to a certain layer depth. Note that by replacing the pixel based reconstruction loss in Eq. 1 with high-level feature reconstruction loss in Eq. 2 we are no longer optimising the variational lower bound, and β -VAE_{DAE} with $\beta = 1$ loses its equivalence to the Variational Autoencoder (VAE) framework as proposed by (Kingma & Welling, 2014; Rezende et al., 2014). In this setting, the only way to interpret β is as a mixing coefficient that balances the capacity of the latent channel \mathbf{z} of β -VAE_{DAE} against the pressure to match the high-level features within the pre-trained DAE.

2.4. Reinforcement Learning Algorithms

We used various RL algorithms (DQN, A3C and Episodic Control: Mnih et al., 2015; 2016; Blundell et al., 2016) to learn the source policy π^S during stage two of the pipeline using the latent states s^z acquired by β -VAE based models during stage one of the DARLA pipeline.

Deep Q Network (DQN) (Mnih et al., 2015) is a variant of the Q-learning algorithm (Watkins, 1989) that utilises deep learning. It uses a neural network to parametrise an approximation for the action-value function $Q(s, a; \theta)$ using parameters θ .

Asynchronous Advantage Actor-Critic (A3C) (Mnih et al., 2016) is an asynchronous implementation of the advantage actor-critic paradigm (Sutton & Barto, 1998; Degris & Sutton, 2012), where separate threads run in parallel and perform updates to shared parameters. The different threads each hold their own instance of the environment and have different exploration policies, thereby decorrelating parameter updates without the need for experience replay. Therefore, A3C is an online algorithm, whereas DQN learns its policy offline, resulting in different learning dynamics be-

tween the two algorithms.

Model-Free Episodic Control (EC) (Blundell et al., 2016) was proposed as a complementary learning system to the other RL algorithms described above. The EC algorithm relies on near-determinism of state transitions and rewards in RL environments; in settings where this holds, it can exploit these properties to memorise which action led to high returns in similar situations in the past. Since in its simplest form EC relies on a lookup table, it learns good policies much faster than value-function-approximation based deep RL algorithms like DQN trained via gradient descent - at the cost of generality (i.e. potentially poor performance in non-deterministic environments).

We also compared our approach to that of *UNREAL* (Jaderberg et al., 2017), a recently proposed RL algorithm which also attempts to utilise unsupervised data in the environment. The UNREAL agent takes as a base an LSTM A3C agent (Mnih et al., 2016) and augments it with a number of unsupervised auxiliary tasks that make use of the rich perceptual data available to the agent besides the (sometimes very sparse) extrinsic reward signals. This auxiliary learning tends to improve the representation learnt by the agent. See Sec. A.6 in Supplementary Materials for further details of the algorithms above.

3. Tasks

We evaluate the performance of DARLA on different task and environment setups that probe subtly different aspects of domain adaptation. As a reminder, in Sec. 2.2 we defined $\hat{\mathcal{S}}$ as a state space that contains all possible conjunctions of high-level factors of variation necessary to generate any naturalistic observation in any $D_i \in \mathcal{M}$. During domain adaptation scenarios agent observation states are generated according to $s_S^o \sim \mathbf{Sim}_S(\hat{\mathcal{S}}_S)$ and $s_T^o \sim \mathbf{Sim}_T(\hat{\mathcal{S}}_T)$ for the source and target domains respectively, where $\hat{\mathcal{S}}_S$ and $\hat{\mathcal{S}}_T$ are sampled by some distributions or processes \mathcal{G}_S and \mathcal{G}_T according to $\hat{\mathcal{S}}_S \sim \mathcal{G}_S(\hat{\mathcal{S}})$ and $\hat{\mathcal{S}}_T \sim \mathcal{G}_T(\hat{\mathcal{S}})$.

We use DeepMind Lab (Beattie et al., 2016) to test a version of domain adaptation setup where the source and target domain observation simulators are equal ($\mathbf{Sim}_S = \mathbf{Sim}_T$), but the processes used to sample $\hat{\mathcal{S}}_S$ and $\hat{\mathcal{S}}_T$ are different ($\mathcal{G}_S \neq \mathcal{G}_T$). We use the Jaco arm with a matching MuJoCo simulation environment (Todorov et al., 2012) in two domain adaptation scenarios: simulation to simulation (sim2sim) and simulation to reality (sim2real). The sim2sim domain adaptation setup is relatively similar to DeepMind Lab i.e. the source and target domains differ in terms of processes \mathcal{G}_S and \mathcal{G}_T . However, there is a significant point of difference. In DeepMind Lab, all values of factors in the target domain, $\hat{\mathcal{S}}_T$, are previously seen in the source domain; however, $\hat{\mathcal{S}}_S \neq \hat{\mathcal{S}}_T$ as the *conjunctions* of

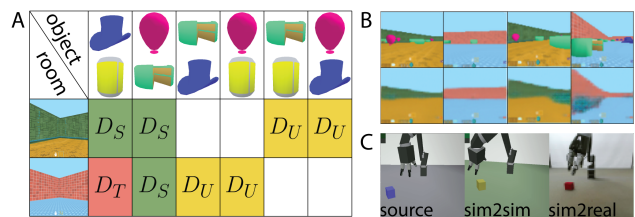


Figure 2. **A:** DeepMind Lab (Beattie et al., 2016) transfer task setup. Different conjunctions of $\{\text{room}, \text{object}_1, \text{object}_2\}$ were used during different parts of the domain adaptation curriculum. During stage one, D_U (shown in yellow), we used a minimal set spanning all objects and all rooms whereby each object is seen in each room. Note there is no extrinsic reward signal or notion of ‘task’ in this phase. During stage two, D_S (shown in green), the RL agents were taught to pick up cans and balloons and avoid hats and cakes. The objects were always presented in pairs hat/can and cake/balloon. The agent never saw the hat/can pair in the pink room. This novel room/object conjunction was presented as the target domain adaptation condition D_T (shown in red) where the ability of the agent to transfer knowledge of the objects’ value to a novel environment was tested. **B:** β -VAE reconstructions (bottom row) using frames from DeepMind Lab (top row). Due to the increased $\beta > 1$ necessary to disentangle the data generative factors of variations the model lost information about objects. See Fig. 3 for a model appropriately capturing objects. **C:** *left* – sample frames from MuJoCo simulation environments used for vision (phase 1, S_U) and source policy training (phase 2, S_S); *middle* – sim2sim domain adaptation test (phase 3, S_T); and *right* – sim2real domain adaptation test (phase 3, S_T).

these factor values are different. In sim2sim, by contrast, novel factor *values* are experienced in the target domain (this accordingly also leads to novel factor conjunctions). Hence, DeepMind Lab may be considered to be assessing domain *interpolation* performance, whereas sim2sim tests domain *extrapolation*.

The sim2real setup, on the other hand, is based on identical processes $\mathcal{G}_S = \mathcal{G}_T$, but different observation simulators $\mathbf{Sim}_S \neq \mathbf{Sim}_T$ corresponding to the MuJoCo simulation and the real world, which results in the so-called ‘perceptual reality gap’ (Rusu et al., 2016). More details of the tasks are given below.

3.1. DeepMind Lab

DeepMind Lab is a first person 3D game environment with rich visuals and realistic physics. We used a standard seek-avoid object gathering setup, where a room is initialised with an equal number of randomly placed objects of two different types. One of the object varieties is ‘good’ (its collection is rewarded +1), while the other is ‘bad’ (its collection is punished -1). The full state space $\hat{\mathcal{S}}$ consisted of all conjunctions of two room types (pink and green based on the colour of the walls) and four object types (hat, can, cake and balloon) (see Fig. 2A). The source domain D_S con-

tained environments with hats/cans presented in the green room, and balloons/cakes presented in either the green or the pink room. The target domain D_T contained hats/cans presented in the pink room. In both domains cans and balloons were the rewarded objects.

1) *Learn to see*: we used β -VAE_{DAE} to learn the disentangled latent state representation s^z that includes both the room and the object generative factors of variation within DeepMind Lab. We had to use the high-level feature space of a pre-trained DAE within the β -VAE_{DAE} framework (see Section 2.3.1), instead of the pixel space of vanilla β -VAE, because we found that objects failed to reconstruct when using the values of β necessary to disentangle the generative factors of variation within DeepMind Lab (see Fig. 2B).

β -VAE_{DAE} was trained on observations s_U^o collected by an RL agent with a simple wall-avoiding policy π_U (otherwise the training data was dominated by close up images of walls). In order to enable the model to learn $\mathcal{F}(s_U^o) \approx \hat{S}$, it is important to expose the agent to at least a minimal set of environments that span the range of values for each factor, and where no extraneous correlations are added between different factors³(see Fig. 2A, yellow). See Section A.3.1 in Supplementary Materials for details of β -VAE_{DAE} training.

2) *Learn to act*: the agent was trained with the algorithms detailed in Section 2.4 on a seek-avoid task using the source domain (D_S) conjunctions of object/room shown in Fig. 2A (green). Pre-trained β -VAE_{DAE} from stage one was used as the ‘vision’ part of various RL algorithms (DQN, A3C and Episodic Control: Mnih et al., 2015; 2016; Blundell et al., 2016) to learn a source policy π_S that picks up balloons and avoids cakes in both the green and the pink rooms, and picks up cans and avoids hats in the green rooms. See Section A.3.1 in Supplementary Materials for more details of the various versions of DARLA we have tried, each based on a different base RL algorithm.

3) *Transfer*: we tested the ability of DARLA to transfer the seek-avoid policy π_S it had learnt on the source domain in stage two using the domain adaptation condition D_T illustrated in Figure 2A (red). The agent had to continue picking up cans and avoid hats in the pink room, even though these objects had only been seen in the green room during source policy training. The optimal policy π_T is one that maintains the reward polarity from the source domain (cans are good and hats are bad). For further details, see Appendix A.2.1.

³In our setup of DeepMind Lab domain adaptation task, the object and environment factors are supposed to be independent. In order to ensure that β -VAE_{DAE} learns a factorised representation that reflects this ground truth independence, we present observations of all possible conjunctions of room and individual object types.

3.2. Jaco Arm and MuJoCo

We used frames from an RGB camera facing a robotic Jaco arm, or a matching rendered camera view from a MuJoCo physics simulation environment (Todorov et al., 2012) to investigate the performance of DARLA in two domain adaptation scenarios: 1) simulation to simulation (sim2sim), and 2) simulation to reality (sim2real). The sim2real setup is of particular importance, since the progress that deep RL has brought to control tasks in simulation (Schulman et al., 2015; Mnih et al., 2016; Levine & Abbeel, 2014; Heess et al., 2015; Lillicrap et al., 2015; Schulman et al., 2016) has not yet translated as well to reality, despite various attempts (Tobin et al., 2017; Tzeng et al., 2016; Daftry et al., 2016; Finn et al., 2015; Rusu et al., 2016). Solving control problems in reality is hard due to sparse reward signals, expensive data acquisition and the attendant danger of breaking the robot (or its human minders) during exploration.

In both sim2sim and sim2real, we trained the agent to perform an object reaching policy where the goal is to place the end effector as close to the object as possible. While conceptually the reaching task is simple, it is a hard control problem since it requires correct inference of the arm and object positions and velocities from raw visual inputs.

1) *Learn to see*: β -VAE was trained on observations collected in MuJoCo simulations with the same factors of variation as in D_S . In order to enable the model to learn $\mathcal{F}(s_U^o) \approx \hat{s}$, a reaching policy was applied to phantom objects placed in random positions - therefore ensuring that the agent learnt the independent nature of the arm position and object position (see Fig. 2C, left);

2) *Learn to act*: a feedforward-A3C based agent with the vision module pre-trained in stage one was taught a source reaching policy π_S towards the real object in simulation (see Fig. 2C (left) for an example frame, and Sec. A.4 in Supplementary Materials for a fuller description of the agent). In the source domain D_S the agent was trained on a distribution of camera angles and positions. The colour of the tabletop on which the arm rests and the object colour were both sampled anew every episode.

3) *Transfer*: **sim2sim**: in the target domain, D_T , the agent was faced with a new distribution of camera angles and positions with little overlap with the source domain distributions, as well as a completely held out set of object colours (see Fig. 2C, middle). **sim2real**: in the target domain D_T the camera position and angle as well as the tabletop colour and object colour were sampled from the same distributions as seen in the source domain D_S , but the target domain D_T was now the real world. Many details present in the real world such as shadows, specularly, multiple light sources and so on are not modelled in the simulation;

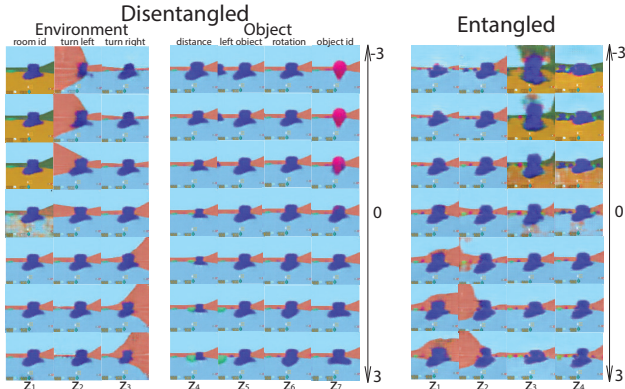


Figure 3. Plot of traversals of various latents of an entangled and a disentangled version of β -VAE_{DAE} using frames from DeepMind Lab (Beattie et al., 2016).

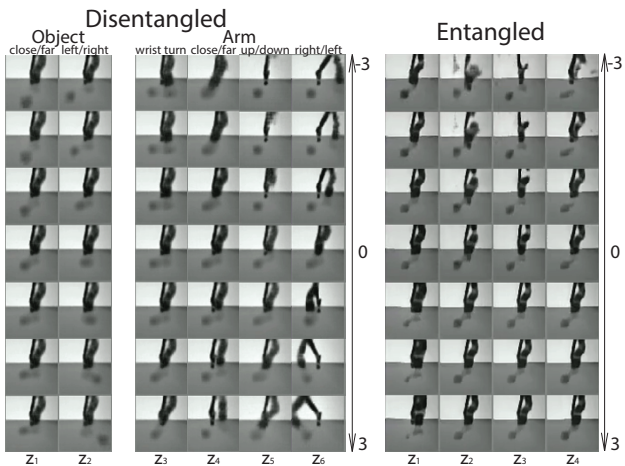


Figure 4. Plot of traversals of β -VAE on MuJoCo. Using a disentangled β -VAE model, single latents directly control for the factors responsible for the object or arm placements.

the physics engine is also not a perfect model of reality. Thus sim2real tests the ability of the agent to cross the perceptual-reality gap and generalise its source policy π_S to the real world (see Fig. 2C, right). For further details, see Appendix A.2.2.

4. Results

We evaluated the robustness of DARLA’s policy π_S learnt on the source domain to various shifts in the input data distribution. In particular, we used domain adaptation scenarios based on the DeepMind Lab seek-avoid task and the Jaco arm reaching task described in Sec. 3. On each task we compared DARLA’s performance to that of various baselines. We evaluated the importance of learning ‘good’ vision during stage one of the pipeline, i.e one that maps the input observations s^o to disentangled representations $s^z \approx \hat{s}$. In order to do this, we ran the DARLA pipeline with different vision models: the encoders of a

disentangled β -VAE⁴ (the original DARLA), an *entangled* β -VAE (DARLA_{ENT}), and a denoising autoencoder (DARLA_{DAE}). Apart from the nature of the learnt representations s^z , DARLA and all versions of its baselines were equivalent throughout the three stages of our proposed pipeline in terms of architecture and the observed data distribution (see Sec. A.3 in Supplementary Materials for more details).

Figs. 3-4 display the degree of disentanglement learnt by the vision modules of DARLA and DARLA_{ENT} on DeepMind Lab and MuJoCo. DARLA’s vision learnt to independently represent environment variables (such as room colour-scheme and geometry) and object-related variables (change of object type, size, rotation) on DeepMind Lab (Fig. 3, left). Disentangling was also evident in MuJoCo. Fig. 4, left, shows that DARLA’s single latent units z_i learnt to represent different aspects of the Jaco arm, the object, and the camera. By contrast, in the representations learnt by DARLA_{ENT}, each latent is responsible for changes to both the environment and objects (Fig. 3, right) in DeepMind Lab, or a mixture of camera, object and/or arm movements (Fig. 4, right) in MuJoCo.

The table in Fig. 5 shows the average performance (across different seeds) in terms of rewards per episode of the various agents on the target domain with no fine-tuning of the source policy π_S . It can be seen that DARLA is able to zero-shot-generalise significantly better than DARLA_{ENT} or DARLA_{DAE}, highlighting the importance of learning a disentangled representation $s^z = s^z_{\hat{s}}$ during the unsupervised stage one of the DARLA pipeline. In particular, this also demonstrates that the improved domain transfer performance is not simply a function of increased exposure to training observations, as both DARLA_{ENT} and DARLA_{DAE} were exposed to the same data. The results are mostly consistent across target domains and in most cases DARLA is significantly better than the second-best-performing agent. This holds in the sim2real task⁵, where being able to perform zero-shot policy transfer is highly valuable due to the particular difficulties of gathering data in the real world.

DARLA’s performance is particularly surprising as it actually preserves *less* information about the raw observations s^o than DARLA_{ENT} and DARLA_{DAE}. This is due to the nature of the β -VAE and how it achieves disentangling; the disentangled model utilised a significantly higher value of the hyperparameter β than the entangled model (see Appendix A.3 for further details), which constrains the ca-

⁴In this section of the paper, we use the term β -VAE to refer to a standard β -VAE for the MuJoCo experiments, and a β -VAE_{DAE} for the DeepMind Lab experiments (as described in stage 1 of Sec. 3.1).

⁵See <https://youtu.be/sZqrWfI0wQ4> for example sim2sim and sim2real zero-shot transfer policies of DARLA and baseline A3C agent.

Table 1. Transfer performance

VISION TYPE	DQN	DEEPMIND LAB		JACO (A3C)	
		A3C	EC	SIM2SIM	SIM2REAL
BASILINE AGENT	1.86 ± 3.91	5.32 ± 3.36	-0.41 ± 4.21	97.64 ± 9.02	94.56 ± 3.55
UNREAL	-	4.13 ± 3.95	-	-	-
DARLA _{FT}	13.36 ± 5.8	1.4 ± 2.16	-	86.59 ± 5.53	99.25 ± 2.3
DARLA _{ENT}	3.45 ± 4.47	15.66 ± 5.19	5.69 ± 3.73	84.77 ± 4.42	59.99 ± 15.05
DARLA _{DAE}	7.83 ± 4.47	6.74 ± 2.81	5.59 ± 3.37	85.15 ± 7.43	100.72 ± 4.7
DARLA	10.25 ± 5.46	19.7 ± 5.43	11.41 ± 3.52	100.85 ± 2.92	108.2 ± 5.97

DARLA’S PERFORMANCE IS SIGNIFICANTLY DIFFERENT FROM ALL BASELINES UNDER WELCH’S UNEQUAL VARIANCES T-TEST WITH $p < 0.01$ ($N \in [60, 150]$).

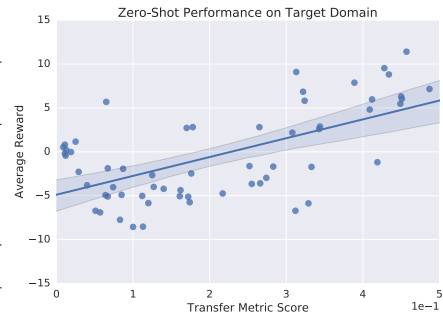


Figure 5. **Table:** Zero-shot performance (avg. reward per episode) of the source policy π_S in target domains within DeepMind Lab and Jaco/MuJoCo environments. Baseline agent refers to vanilla DQN/A3C/EC (DeepMind Lab) or A3C (Jaco) agents. See main text for more detailed model descriptions. **Figure:** Correlation between zero-shot performance transfer performance on the DeepMind Lab task obtained by EC based DARLA and the level of disentanglement as measured by the transfer/disentanglement score ($r = 0.6$, $p < 0.001$)

capacity of the latent channel. Indeed, DARLA’s β -VAE only utilises 8 of its possible 32 Gaussian latents to store observation-specific information for MuJoCo/Jaco (and 20 in DeepMind Lab), whereas DARLA_{ENT} utilises all 32 for both environments (as does DARLA_{DAE}).

Furthermore, we examined what happens if DARLA’s vision (i.e. the encoder of the disentangled β -VAE) is allowed to be fine-tuned via gradient updates while learning the source policy during stage two of the pipeline. This is denoted by DARLA_{FT} in the table in Fig. 5. We see that it exhibits significantly worse performance than that of DARLA in zero-shot domain adaptation using an A3C-based agent in all tasks. This suggests that a favourable initialisation does not make up for subsequent overfitting to the source domain for the on-policy A3C. However, the off-policy DQN-based fine-tuned agent performs very well. We leave further investigation of this curious effect for future work.

Finally, we compared the performance of DARLA to an UNREAL (Jaderberg et al., 2017) agent with the same architecture. Despite also exploiting the unsupervised data available in the source domain, UNREAL performed worse than baseline A3C on the DeepMind Lab domain adaptation task. This further demonstrates that use of unsupervised data in itself is not a panacea for transfer performance; it must be utilised in a careful and structured manner conducive to learning disentangled latent states $s^z = s^z_S$.

In order to quantitatively evaluate our hypothesis that disentangled representations are essential for DARLA’s performance in domain adaptation scenarios, we trained various DARLAs with different degrees of learnt disentanglement in s^z by varying β (of β -VAE) during stage one of the pipeline. We then calculated the correlation between the performance of the EC-based DARLA on the DeepMind Lab domain adaptation task and the transfer metric, which approximately measures the quality of disentanglement of DARLA’s latent representations s^z (see Sec. A.5.2

in Supplementary Materials). This is shown in the chart in Fig. 5; as can be seen, there is a strong positive correlation between the level of disentanglement and DARLA’s zero-shot domain transfer performance ($r = 0.6$, $p < 0.001$).

Having shown the robust utility of disentangled representations in agents for domain adaptation, we note that there is evidence that they can provide an important additional benefit. We found significantly improved speed of learning of π_S on the source domain itself, as a function of how disentangled the model was. The gain in data efficiency from disentangled representations for source policy learning is not the main focus of this paper so we leave it out of the main text; however, we provide results and discussion in Section A.7 in Supplementary Materials.

5. Conclusion

We have demonstrated the benefits of using disentangled representations in a deep RL setting for domain adaptation. In particular, we have proposed DARLA, a multi-stage RL agent. DARLA first learns a visual system that encodes the observations it receives from the environment as disentangled representations, in a completely unsupervised manner. It then uses these representations to learn a robust source policy that is capable of zero-shot domain adaptation.

We have demonstrated the efficacy of this approach in a range of domains and task setups: a 3D naturalistic first-person environment (DeepMind Lab), a simulated graphics and physics engine (MuJoCo), and crossing the simulation to reality gap (MuJoCo to Jaco sim2real). We have also shown that the effect of disentangling is consistent across very different RL algorithms (DQN, A3C, EC), achieving significant improvements over the baseline algorithms (median 2.7 times improvement in zero-shot transfer across tasks and algorithms). To the best of our knowledge, this is the first comprehensive empirical demonstration of the strength of disentangled representations for domain adaptation in a deep RL setting.

References

- Abadi, Martin, Agarwal, Ashish, and et al, Paul Barham. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *Preliminary White Paper*, 2015.
- Barreto, André, Munos, Rémi, Schaul, Tom, and Silver, David. Successor features for transfer in reinforcement learning. *CoRR*, abs/1606.05312, 2016. URL <http://arxiv.org/abs/1606.05312>.
- Beattie, Charles, Leibo, Joel Z., Teplyashin, Denis, Ward, Tom, and et al, Marcus Wainwright. Deepmind lab. *arxiv*, 2016.
- Bengio, Y., Courville, A., and Vincent, P. Representation learning: A review and new perspectives. In *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 2013.
- Blundell, Charles, Uria, Benigno, Pritzel, Alexander, Li, Yazhe, Ruderman, Avraham, Leibo, Joel Z, Rae, Jack, Wierstra, Daan, and Hassabis, Demis. Model-free episodic control. *arXiv*, 2016.
- Candy, T. Rowan, Wang, Jingyun, and Ravikumar, Sowmya. Retinal image quality and postnatal visual experience during infancy. *Optom Vis Sci*, 86(6):556–571, 2009.
- Chen, Xi, Duan, Yan, Houthoofd, Rein, Schulman, John, Sutskever, Ilya, and Abbeel, Pieter. Infogan: Interpretable representation learning by information maximizing generative adversarial nets. *arXiv*, 2016.
- Cheung, Brian, Levezey, Jesse A., Bansal, Arjun K., and Olshausen, Bruno A. Discovering hidden factors of variation in deep networks. In *Proceedings of the International Conference on Learning Representations, Workshop Track*, 2015.
- Cohen, T. and Welling, M. Transformation properties of learned visual representations. In *ICLR*, 2015.
- Cohen, Taco and Welling, Max. Learning the irreducible representations of commutative lie groups. *arXiv*, 2014.
- Daftry, Shreyansh, Bagnell, J. Andrew, and Hebert, Martial. Learning transferable policies for monocular reactive mav control. *International Symposium on Experimental Robotics*, 2016.
- Degrís, Thomas, Pilarski Patrick M and Sutton, Richard S. Model-free reinforcement learning with continuous action in practice. *American Control Conference (ACC)*, 22:21772182, 2012.
- Desjardins, G., Courville, A., and Bengio, Y. Disentangling factors of variation via generative entangling. *arXiv*, 2012.
- Dosovitskiy, Alexey and Brox, Thomas. Generating images with perceptual similarity metrics based on deep networks. *arxiv*, 2016.
- Finn, Chelsea, Tan, Xin Yu, Duan, Yan, Darrell, Trevor, Levine, Sergey, and Abbeel, Pieter. Deep spatial autoencoders for visuomotor learning. *arxiv*, 2015.
- Finn, Chelsea, Yu, Tianhe, Fu, Justin, Abbeel, Pieter, and Levine, Sergey. Generalizing skills with semi-supervised reinforcement learning. *ICLR*, 2017.
- Garnelo, Marta, Arulkumaran, Kai, and Shanahan, Murray. Towards deep symbolic reinforcement learning. *arXiv*, 2016.
- Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., and Bengio, Y. Generative adversarial nets. *NIPS*, pp. 26722680, 2014.
- Goroshin, Ross, Mathieu, Michael, and LeCun, Yann. Learning to linearize under uncertainty. *NIPS*, 2015.
- Guez, Arthur, Silver, David, and Dayan, Peter. Efficient bayesian adaptive reinforcement learning using sample-based search. *NIPS*, 2012.
- Gupta, Abhishek, Devin, Coline, Liu, YuXuan, Abbeel, Pieter, and Levine, Sergey. Learning invariant feature spaces to transfer skills with reinforcement learning. *ICLR*, 2017.
- Heess, Nicolas, Wayne, Gregory, Silver, David, Lillicrap, Timothy P., Erez, Tom, and Tassa, Yuval. Learning continuous control policies by stochastic value gradients. *NIPS*, 2015.
- Higgins, Irina, Matthey, Loic, Pal, Arka, Burgess, Christopher, Glorot, Xavier, Botvinick, Matthew, Mohamed, Shakir, and Lerchner, Alexander. Beta-vae: Learning basic visual concepts with a constrained variational framework. In *ICLR*, 2017.
- Hinton, G., Krizhevsky, A., and Wang, S. D. Transforming autoencoders. *International Conference on Artificial Neural Networks*, 2011.
- Jaderberg, Max, Mnih, Volodymyr, Czarnecki, Wojciech Marian, Schaul, Tom, Leibo, Joel Z, Silver, David, and Kavukcuoglu, Koray. Reinforcement learning with unsupervised auxiliary tasks. *ICLR*, 2017.
- Karaletsos, Theofanis, Belongie, Serge, and Rtsch, Gunnar. Bayesian representation learning with oracle constraints. *ICLR*, 2016.
- Kingma, D. P. and Ba, Jimmy. Adam: A method for stochastic optimization. *arXiv*, 2014.
- Kingma, D. P. and Welling, M. Auto-encoding variational bayes. *ICLR*, 2014.
- Krizhevsky, Alex, Sutskever, Ilya, and Hinton, Geoffrey E. Imagenet classification with deep convolutional neural networks. *NIPS*, 2012.
- Kulkarni, Tejas, Whitney, William, Kohli, Pushmeet, and Tenenbaum, Joshua. Deep convolutional inverse graphics network. *NIPS*, 2015.
- Lake, Brenden M., Ullman, Tomer D., Tenenbaum, Joshua B., and Gershman, Samuel J. Building machines that learn and think like people. *arXiv*, 2016.
- Larsen, Anders Boesen Lindbo, Snderby, Sren Kaae, Larochelle, Hugo, and Winther, Ole. Autoencoding beyond pixels using a learned similarity metric. *ICML*, 2016.
- Leat, Susan J., Yadav, Naveen K., and Irving, Elizabeth L. Development of visual acuity and contrast sensitivity in children. *Journal of Optometry*, 2009.
- Levine, Sergey and Abbeel, Pieter. Learning neural network policies with guided policy search under unknown dynamics. *NIPS*, 2014.

- Lillicrap, Timothy P., Hunt, Jonathan J., Pritzel, Alexander, Heess, Nicolas, Erez, Tom, Tassa, Yuval, Silver, David, and Wierstra, Daan. Continuous control with deep reinforcement learning. *CoRR*, 2015.
- Littman, Michael L., Sutton, Richard S., and Singh, Satinder. Predictive representations of state. *NIPS*, 2001.
- Marr, D. Simple memory: A theory for archicortex. *Philosophical Transactions of the Royal Society of London*, pp. 2381, 1971.
- McClelland, James L., McNaughton, Bruce L., and O'Reilly, Randall C. Why there are complementary learning systems in the hippocampus and neocortex: insights from the successes and failures of connectionist models of learning and memory. *Psychological review*, 102:419, 1995.
- Mnih, Volodymyr, Kavukcuoglu, Koray, Silver, David S, and Rusu, Andrei A. et al. Human-level control through deep reinforcement learning. *Nature*, 2015.
- Mnih, Volodymyr, Badia, Adri Puigdomnech, Mirza, Mehdi, Graves, Alex, Lillicrap, Timothy P., Harley, Tim, Silver, David, and Kavukcuoglu, Koray. Asynchronous methods for deep reinforcement learning. *ICML*, 2016. URL <https://arxiv.org/pdf/1602.01783.pdf>.
- Niekum, Scott, Chitta, Sachin, Barto, Andrew G, Marthi, Bhaskara, and Osentoski, Sarah. Incremental semantically grounded learning from demonstration. *Robotics: Science and Systems*, 2013.
- Norman, Kenneth A and O'Reilly, Randall C. Modeling hippocampal and neocortical contributions to recognition memory: a complementary-learning-systems approach. *Psychological review*, 110:611, 2003.
- Pan, Sinno Jialin and Yang, Quiang. A survey on transfer learning. *IEEE Transactions on Knowledge and Data Engineering*, 2009.
- Parisotto, Emilio, Ba, Jimmy, and Salakhutdinov, Ruslan. Actor-mimic: Deep multitask and transfer reinforcement learning. *CoRR*, 2015.
- Pathak, Deepak, Krähenbühl, Philipp, Donahue, Jeff, Darrell, Trevor, and Efros, Alexei A. Context encoders: Feature learning by inpainting. *CoRR*, abs/1604.07379, 2016. URL <http://arxiv.org/abs/1604.07379>.
- Peng, J. Efficient dynamic programming-based learning for control. *PhD thesis, Northeastern University, Boston.*, 1993.
- Peng, Jing and Williams, Ronald J. Incremental multi-step q-learning. *Machine Learning*, 22:283290, 1996.
- Puterman, Martin L. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. John Wiley & Sons, Inc., New York, NY, USA, 1st edition, 1994. ISBN 0471619779.
- Raffin, Antonin, Hfer, Sebastian, Jonschkowski, Rico, Brock, Oliver, and Stulp, Freek. Unsupervised learning of state representations for multiple tasks. *ICLR*, 2017.
- Rajendran, Janarthanan, Lakshminarayanan, Aravind, Khapra, Mitesh M., P, Prasanna, and Ravindran, Balaraman. Attend, adapt and transfer: Attentive deep architecture for adaptive transfer from multiple sources in the same domain. *ICLR*, 2017.
- Reed, Scott, Sohn, Kihyuk, Zhang, Yuting, and Lee, Honglak. Learning to disentangle factors of variation with manifold interaction. *ICML*, 2014.
- Rezende, Danilo J., Mohamed, Shakir, and Wierstra, Daan. Stochastic backpropagation and approximate inference in deep generative models. *arXiv*, 2014.
- Ridgeway, Karl. A survey of inductive biases for factorial Representation-Learning. *arXiv*, 2016. URL <http://arxiv.org/abs/1612.05299>.
- Rippel, Oren and Adams, Ryan Prescott. High-dimensional probability estimation with deep density models. *arXiv*, 2013.
- Rusu, Andrei A., Vecerik, Matej, Rothl, Thomas, Heess, Nicolas, Pascanu, Razvan, and Hadsell, Raia. Sim-to-real robot learning from pixels with progressive nets. *arXiv*, 2016.
- Schmidhuber, Jürgen. Learning factorial codes by predictability minimization. *Neural Computation*, 4(6):863–869, 1992.
- Schulman, John, Levine, Sergey, Moritz, Philipp, Jordan, Michael I., and Abbeel, Pieter. Trust region policy optimization. *ICML*, 2015.
- Schulman, John, Moritz, Philipp, Levine, Sergey, Jordan, Michael, and Abbeel, Pieter. High-dimensional continuous control using generalized advantage estimation. *ICLR*, 2016.
- Sutherland, Robert J and Rudy, Jerry W. Configural association theory: The role of the hippocampal formation in learning, memory, and amnesia. *Psychobiology*, 17:129144, 1989.
- Sutton, Richard S. and Barto, Andrew G. *Reinforcement Learning: An Introduction*. MIT Press, 1998.
- Talvitie, Erik and Singh, Satinder. An experts algorithm for transfer learning. In *Proceedings of the 20th international joint conference on Artificial intelligence*, 2007.
- Tobin, Josh, Fong, Rachel, Ray, Alex, Schneider, Jonas, Zaremba, Wojciech, and Abbeel, Pieter. Domain randomization for transferring deep neural networks from simulation to the real world. *arxiv*, 2017.
- Todorov, E., Erez, T., and Tassa, Y. Mujoco: A physics engine for model-based control. *IROS*, 2012.
- Tulving, Endel, Hayman, CA, and Macdonald, Carol A. Long-lasting perceptual priming and semantic learning in amnesia: a case experiment. *Journal of Experimental Psychology: Learning, Memory, and Cognition*, 17:595, 1991.
- Tzeng, Eric, Devin, Coline, Hoffman, Judy, Finn, Chelsea, Abbeel, Pieter, Levine, Sergey, Saenko, Kate, and Darrell, Trevor. Adapting deep visuomotor representations with weak pairwise constraints. *WAFR*, 2016.
- Vincent, Pascal, Larochelle, Hugo, Lajoie, Isabelle, Bengio, Yoshua, and Manzagol, Pierre-Antoine. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *NIPS*, 2010.
- Warde-Farley, David and Bengio, Yoshua. Improving generative adversarial networks with denoising feature matching. *ICLR*, 2017.

Watkins, Christopher John Cornish Hellaby. *Learning from delayed rewards*. PhD thesis, University of Cambridge, Cambridge, UK, 1989.

Whitney, William F., Chang, Michael, Kulkarni, Tejas, and Tenenbaum, Joshua B. Understanding visual concepts with continuation learning. *arXiv*, 2016. URL <http://arxiv.org/pdf/1602.06822.pdf>.

Yang, Jimei, Reed, Scott, Yang, Ming-Hsuan, and Lee, Honglak. Weakly-supervised disentangling with recurrent transformations for 3d view synthesis. *NIPS*, 2015.

A. Supplementary Materials

A.1. The Reinforcement Learning Paradigm

The reinforcement learning (RL) paradigm consists of an agent receiving a sequence of observations s_t^o which are some function of environment states $s_t \in \mathcal{S}$ and may be accompanied by rewards $r_{t+1} \in R$ conditional on the actions $a_t \in \mathcal{A}$, chosen at each time step t (Sutton & Barto, 1998). We assume that these interactions can be modelled as a Markov Decision Process (MDP) (Puterman, 1994) defined as a tuple $D \equiv (\mathcal{S}, \mathcal{A}, \mathcal{T}, R, \gamma)$. $\mathcal{T} = p(s|s_t, a_t)$ is a transition function that models the distribution of all possible next states given action a_t is taken in state s_t for all $s_t \in \mathcal{S}$ and $a_t \in \mathcal{A}$. Each transition $s_t \xrightarrow{a_t} s_{t+1}$ may be accompanied by a reward signal $r_{t+1}(s_t, a_t, s_{t+1})$. The goal of the agent is to learn a policy $\pi(a_t|s_t)$, a probability distribution over actions $a_t \in \mathcal{A}$, that maximises the expected return i.e. the discounted sum of future rewards $R_t = \mathbb{E}[\sum_{\tau=1}^{T-t} \gamma^{\tau-1} r_{t+\tau}]$. T is the time step at which each episode ends, and $\gamma \in [0, 1)$ is the discount factor that progressively down-weights future rewards. Given a policy $\pi(a|s)$, one can define the value function $V_\pi(s) = \mathbb{E}[R_t|s_t = s, \pi]$, which is the expected return from state s following policy π . The action-value function $Q_\pi(s, a) = \mathbb{E}[R_t|s_t = s, a_t = a, \pi]$ is the expected return for taking action a in state s at time t , and then following policy π from time $t + 1$ onward.

A.2. Further task details

A.2.1. DEEPMIND LAB

As described in Sec 3.1, in each source episode of DeepMind Lab the agent was presented with one of three possible room/object type conjunctions, chosen at random. These are marked D_S in Fig 2. The setup was a seek-avoid style task, where one of the two object types in the room gave a reward of +1 and the other gave a reward of -1. The agent was allowed to pick up objects for 60 seconds after which the episode would terminate and a new one would begin; if the agent was able to pick up all the ‘good’ objects in less than 60 seconds, a new episode was begun immediately. The agent was spawned in a random location in the room at the start of each new episode.

During transfer, the agent was placed into the held out conjunction of object types and room background; see D_T in Fig 2.

Visual pre-training was performed in other conjunctions of object type and room background denoted D_U in Fig 2.

The observation size of frames in the DeepMind Lab task was $84 \times 84 \times 3$ ($H \times W \times C$).

A.2.2. MUJoCo/JACO ARM EXPERIMENTS

As described in Sec 3.2, the source task consisted of an agent learning to control a simulated arm in order to reach toward an object. A shaping reward was used, with a maximum value of 1 when the centre of the object fell between the pinch and grip sites of the end effector, or within a 10cm distance of the two. Distances on the x and y dimensions counted double compared to distances on the z dimension.

During each episode the object was placed at a random drop point within a 40x40cm area, and the arm was set to a random initial start position high above the work-space, independent of the object’s position. Each episode lasted for 150 steps, or 7.5 seconds, with a control step of 50ms. Observations s_U^o were sampled

randomly across episodes. Overall, 4 million frames of dimensions $64 \times 64 \times 3$ ($H \times W \times C$) were used for this stage of the curriculum. For each episode the camera position and orientation were randomly sampled from an isotropic normal distribution centred around the approximate position and orientation of the real camera, with standard deviation 0.01. No precise measurements were used to match the two. Work-space table colour was sampled uniformly between -5% and $+5\%$ around the midpoint, independently for each RGB channel; object colours were sampled uniformly at random in RGB space, rejecting colours which fell within a ball around 10 held-out intensities (radius 10% of range); the latter were only used for simulated transfer experiments, i.e. in D_T in the sim2sim experiments. Additionally, Gaussian noise with standard deviation 0.01 was added to the observations s_T^o in the sim2sim task.

For the real Jaco arm and its MuJoCo simulation counterpart, each of the nine joints could independently take 11 different actions (a linear discretisation of the continuous velocity action space). In simulation Gaussian noise with standard deviation 0.1 was added to each discrete velocity output; delays in the real setup between observations and action execution were simulated by randomly mixing velocity outputs from two previous steps instead of emitting the last output directly. Speed ranges were between -50% and 50% of the Jaco arm’s top speed on joints 1 through 6 starting at the base, while the fingers could use a full range. For safety reasons the speed ranges have been reduced by a factor of 0.3 while evaluating agents on the Jaco arm, without significant performance degradation.

A.3. Vision model details

A.3.1. DENOISING AUTOENCODER FOR β -VAE

A denoising autoencoder (DAE) was used as a model to provide the feature space for the β -VAE reconstruction loss to be computed over (for motivation, see Sec. 2.3.1). The DAE was trained with occlusion-style masking noise in the vein of (Pathak et al., 2016), with the aim for the DAE to learn a semantic representation of the input frames. Concretely, two values were independently sampled from $U[0, W]$ and two from $U[0, H]$ where W and H were the width and height of the input frames. These four values determined the corners of the rectangular mask applied; all pixels that fell within the mask were set to zero.

The DAE architecture consisted of four convolutional layers, each with kernel size 4 and stride 2 in both the height and width dimensions. The number of filters learnt for each layer was $\{32, 32, 64, 64\}$ respectively. The bottleneck layer consisted of a fully connected layer of size 100 neurons. This was followed by four deconvolutional layers, again with kernel sizes 4, strides 2, and $\{64, 64, 32, 32\}$ filters. The padding algorithm used was ‘SAME’ in TensorFlow (Abadi et al., 2015). ReLU non-linearities were used throughout.

The model was trained with loss given by the L2 distance of the outputs from the original, un-noised inputs. The optimiser used was Adam (Kingma & Ba, 2014) with a learning rate of $1e-3$.

A.3.2. β -VAE WITH PERCEPTUAL SIMILARITY LOSS

After training a DAE, as detailed in the previous section⁶, a β -VAE_{DAE} was trained with perceptual similarity loss given by

⁶In principle, the β -VAE_{DAE} could also have been trained end-to-end in one pass, but we did not experiment with this.

Eq. 2, repeated here:

$$\mathcal{L}(\theta, \phi; \mathbf{x}, \mathbf{z}, \beta) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \|J(\hat{\mathbf{x}}) - J(\mathbf{x})\|_2^2 - \beta D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})\|p(\mathbf{z})) \quad (3)$$

Specifically, the input was passed through the β -VAE and a sampled⁷ reconstruction was passed through the pre-trained DAE up to a designated layer. The L2 distance of this representation from the representation of the original input passed through the same layers of the DAE was then computed, and this formed the training loss for the β -VAE part of the β -VAE_{DAE}⁸. The DAE weights remained frozen throughout.

The β -VAE architecture consisted of an encoder of four convolutional layers, each with kernel size 4, and stride 2 in the height and width dimensions. The number of filters learnt for each layer was $\{32, 32, 64, 64\}$ respectively. This was followed by a fully connected layer of size 256 neurons. The latent layer comprised 64 neurons parametrising 32 (marginally) independent Gaussian distributions. The decoder architecture was simply the reverse of the encoder, utilising deconvolutional layers. The decoder used was Gaussian, so that the number of output channels was $2C$, where C was the number of channels that the input frames had. The padding algorithm used was ‘SAME’ in TensorFlow. ReLU non-linearities were used throughout.

The model was trained with the loss given by Eq. 3. Specifically, the disentangled model used for DARLA was trained with a β hyperparameter value of 1 and the layer of the DAE used to compute the perceptual similarity loss was the last deconvolutional layer. The entangled model used for DARLA_{ENT} was trained with a β hyperparameter value of 0.1 with the last deconvolutional layer of the DAE used to compute the perceptual similarity loss.

The optimiser used was Adam with a learning rate of $1e-4$.

A.3.3. β -VAE

For the MuJoCo/Jaco tasks, a standard β -VAE was used rather than the β -VAE_{DAE} used for DeepMind Lab. The architecture of the VAE encoder, decoder and the latent size were exactly as described in the previous section A.3.2. β for the disentangled β -VAE in DARLA was 175. β for the entangled model DARLA_{ENT} was 1, corresponding to the standard VAE of (Kingma & Welling, 2014).

The optimizer used was Adam with a learning rate of $1e-4$.

A.3.4. DENOISING AUTOENCODER FOR BASELINE

For the baseline model DARLA_{DAE}, we trained a denoising autoencoder with occlusion-style masking noise as described in Appendix Section A.3.1. The architecture used matched that exactly of the β -VAE described in Appendix Section A.3.2 - however, all stochastic nodes were replaced with deterministic neurons.

The optimizer used was Adam with a learning rate of $1e-4$.

⁷It is more typical to use the mean of the reconstruction distribution, but this does not induce any pressure on the Gaussians parametrising the decoder to reduce their variances. Hence full samples were used instead.

⁸The representations were taken after passing through the layer but before passing through the following non-linearity. We also briefly experimented with taking the L2 loss post-activation but did not find a significant difference.

A.4. Reinforcement Learning Algorithm Details

A.4.1. DEEPMIND LAB

The action space in the DeepMind Lab task consisted of 8 discrete actions.

DQN: in DQN, the convolutional (or ‘vision’) part of the Q-net was replaced with the encoder of the β -VAE_{DAE} from stage 1 and frozen. DQN takes four consecutive frames as input in order to capture some aspect of environment dynamics in the agent’s state. In order to match this in our setup with a pre-trained vision stack \mathcal{F}_U , we passed each observation frame $s_{\{1..4\}}^o$ through the pre-trained model $s_{\{1..4\}}^z = \mathcal{F}_U(s_{\{1..4\}}^o)$ and then concatenated the outputs together to form the k -dimensional (where $k = 4|s^z|$) input to the policy network. In this case the size of s^z was 64 for DARLA as well as DARLA_{ENT}, DARLA_{DAE} and DARLA_{FT}.

On top of the frozen convolutional stack, two ‘policy’ layers of 512 neurons each were used, with a final linear layer of 8 neurons corresponding to the size of the action space in the DeepMind Lab task. ReLU non-linearities were used throughout. All other hyperparameters were as reported in (Mnih et al., 2015).

A3C: in A3C, as with DQN, the convolutional part of the network that is shared between the policy net and the value net was replaced with the encoder of the β -VAE_{DAE} in DeepMind Lab tasks. All other hyperparameters were as reported in (Mnih et al., 2016).

Episodic Control: for the Episodic Controller-based DARLA we used mostly the same hyperparameters as in the original paper by (Blundell et al., 2016). We explored the following hyperparameter settings: number of nearest neighbours $\in \{10, 50\}$, return horizon $\in \{100, 400, 800, 1800, 500000\}$, kernel type $\in \{\text{inverse, gaussian}\}$, kernel width $\in \{1e-6, 1e-5, 1e-4, 1e-3, 1e-2, 1e-1, 0.5, 0.99\}$ and we tried training EC with and without Peng’s $Q(\lambda)$ (Peng, 1993). In practice we found that none of the explored hyperparameter choices significantly influenced the results of our experiments. The final hyperparameters used for all experiments reported in the paper were the following: number of nearest neighbours: 10, return horizon: 400, kernel type: inverse, kernel width: 1e-6 and no Peng’s $Q(\lambda)$ (Peng, 1993).

UNREAL: We used a vanilla version of UNREAL, with parameters as reported in (Jaderberg et al., 2017).

A.4.2. MUJoCo/JACO ARM EXPERIMENTS

For the real Jaco arm and its MuJoCo simulation, each of the nine joints could independently take 11 different actions (a linear discretisation of the continuous velocity action space). Therefore the action space size was 99.

DARLA for MuJoCo/Jaco was based on feedforward A3C (Mnih et al., 2016). We closely followed the simulation training setup of (Rusu et al., 2016) for feed-forward networks using raw visual-input only. In place of the usual conv-stack, however, we used the encoder of the β -VAE as described in Appendix A.3.3. This was followed by a linear layer with 512 units, a ReLU non-linearity and a collection of 9 linear and softmax layers for the 9 independent policy outputs, as well as a single value output layer that outputted the value function.

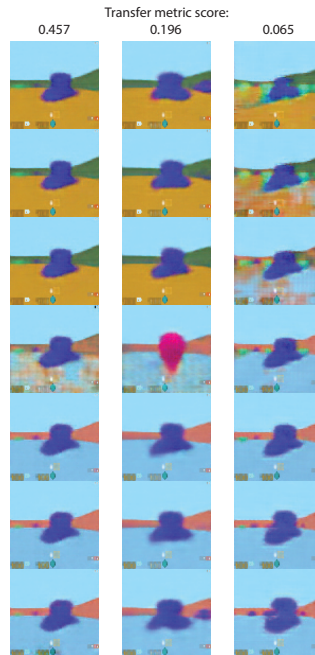


Figure 6. Traversals of the latent corresponding to room background for models with different transfer metric scores (shown top). Note that in the entangled model, many other objects appear and blue hat changes shape in addition to the background changing. For the model with middling transfer score, both the object type and background alter; whereas for the disentangled model, very little apart from the background changes.

A.5. Disentanglement Evaluation

A.5.1. VISUAL HEURISTIC DETAILS

In order to choose the optimal value of β for the β -VAE-DAE models and evaluate the fitness of the representations s_U^z learnt in stage 1 of our pipeline (in terms of disentanglement achieved), we used the visual inspection heuristic described in (Higgins et al., 2017). The heuristic involved clustering trained β -VAE based models based on the number of informative latents (estimated as the number of latents z_i with average inferred standard deviation below 0.75). For each cluster we examined the degree of learnt disentanglement by running inference on a number of seed images, then traversing each latent unit $z_{\{i\}}$ one at a time over three standard deviations away from its average inferred mean while keeping all other latents $z_{\{\setminus i\}}$ fixed to their inferred values. This allowed us to visually examine whether each individual latent unit z_i learnt to control a single interpretable factor of variation in the data. A similar heuristic has been the *de rigueur* method for exhibiting disentanglement in the disentanglement literature (Chen et al., 2016; Kulkarni et al., 2015).

A.5.2. TRANSFER METRIC DETAILS

In the case of DeepMind Lab, we were able to use the ground truth labels corresponding to the two factors of variation of the object type and the background to design a proxy to the disentanglement metric proposed in (Higgins et al., 2017). The procedure used

consisted of the following steps:

1) Train the model under consideration on observations s_V^o to learn \mathcal{F}_U , as described in stage 1 of the DARLA pipeline.

2) Learn a linear model $\mathcal{L} : S_V^z \rightarrow M \times N$ from the representations $s_V^z = \mathcal{F}_V(s_V^o)$, where $M \in \{0, 1\}$ corresponds to the set of possible rooms and $N \in \{0, 1, 2, 3\}$ corresponds to the set of possible objects⁹. Therefore we are learning a low-VC dimension classifier to predict the room and the object class from the latent representation of the model. Crucially, the linear model \mathcal{L} is trained on only a subset of the Cartesian product $M \times N$ e.g. on $\{\{0, 0\}, \{0, 3\}, \{1, 1\}, \{1, 2\}\}$. In practice, we utilised a softmax classifier each for M and N and trained this using backpropagation with a cross-entropy loss, keeping the unsupervised model (and therefore \mathcal{F}_U) fixed.

3) The trained linear model \mathcal{L} 's accuracy is evaluated on the held out subset of the Cartesian product $M \times N$.

Although the above procedure only measures disentangling up to linearity, and only does so for the latents of object type and room background, we nevertheless found that the metric was highly correlated with disentanglement as determined via visual inspection (see Fig. 6).

A.6. Background on RL Algorithms

In this Appendix, we provide background on the different RL algorithms that the DARLA framework was tested on in this paper.

A.6.1. DQN

(DQN) (Mnih et al., 2015) is a variant of the Q-learning algorithm (Watkins, 1989) that utilises deep learning. It uses a neural network to parametrise an approximation for the action-value function $Q(s, a; \theta)$ using parameters θ . These parameters are updated by minimising the mean-squared error of a 1-step look-ahead loss $\mathcal{L}_Q = \mathbb{E}[(r_t + \gamma \max_{a'} Q(s', a'; \theta^-) - Q(s, a; \theta))^2]$, where θ^- are parameters corresponding to a frozen network and optimisation is performed with respect to θ , with θ^- being synced to θ at regular intervals.

A.6.2. A3C

Asynchronous Advantage Actor-Critic (A3C) (Mnih et al., 2016) is an asynchronous implementation of the advantage actor-critic paradigm (Sutton & Barto, 1998; Degris & Sutton, 2012), where separate threads run in parallel and perform updates to shared parameters. The different threads each hold their own instance of the environment and have different exploration policies, thereby decorrelating parameter updates without the need for experience replay.

A3C uses neural networks to approximate both policy $\pi(a|s; \theta)$ and value $V_\pi(s; \theta)$ functions using parameters θ using n-step look-ahead loss (Peng & Williams, 1996). The algorithm is trained using an advantage actor-critic loss function with an entropy regularisation penalty: $\mathcal{L}_{A3C} \approx \mathcal{L}_{VR} + \mathcal{L}_\pi - \mathbb{E}_{s \sim \pi}[\alpha H(\pi(a|s; \theta))]$, where H is entropy. The parameter updates are performed after every t_{max} actions or when a terminal state is reached. $\mathcal{L}_{VR} =$

⁹For the purposes of this metric, we utilised rooms with only single objects, which we denote by the subscript V e.g. the observation set S_V^o .

$\mathbb{E}_{s \sim \pi} [(R_{t:t+n} + \gamma^n V(s_{t+n+1}; \theta) - V(s_t; \theta))^2]$ and $\mathcal{L}_\pi = \mathbb{E}_{s \sim \pi} [\log \pi(a|s; \theta)(Q^\pi(s, a; \theta) - V^\pi(s; \theta))]$. Unlike DQN, A3C uses an LSTM core to encode its history and therefore has a longer term memory permitting it to perform better in partially observed environments. In the version of A3C used in this paper for the DeepMind Lab task, the policy net additionally takes the last action a_{t-1} and last reward r_{t-1} as inputs along with the observation s_t^o , as introduced in (Jaderberg et al., 2017).

A.6.3. UNREAL

The *UNREAL* agent (Jaderberg et al., 2017) takes as a base an LSTM A3C agent (Mnih et al., 2016) and augments it with a number of unsupervised auxiliary tasks that make use of the rich perceptual data available to the agent besides the (sometimes very sparse) extrinsic reward signals. This auxiliary learning tends to improve the representation learnt by the agent. While training the base agent, its observations, rewards, and actions are stored in a replay buffer, which is used by the auxiliary learning tasks. The tasks include: 1) pixel control - the agent learns how to control the environment by training auxiliary policies to maximally change pixel intensities in different parts of the input; 2) reward prediction - given a replay buffer of observations within a short time period of an extrinsic reward, the agent has to predict the reward obtained during the next unobserved timestep using a sequence of three preceding steps; 3) value function replay - extra training of the value function to promote faster value iteration.

A.6.4. EPISODIC CONTROL

In its simplest form EC is a lookup table of states and actions denoted as $Q^{EC}(s, a)$. In each state EC picks the action with the highest Q_{EC} value. At the end of each episode $Q^{EC}(s, a)$ is set to R_t if $(s_t, a_t) \notin Q^{EC}$, where R_t is the discounted return. Otherwise $Q^{EC}(s, a) = \max\{Q^{EC}(s, a), R_t\}$. In order to generalise its policy to novel states that are not in Q^{EC} , EC uses a non-parametric nearest neighbours search $\widehat{Q^{EC}}(s, a) = \frac{1}{k} \sum_{i=1}^k Q^{EC}(s^i, a)$, where $s^i, i = 1, \dots, k$ are k states with the smallest distance to the novel state s . Like DQN, EC takes a concatenation of four frames as input.

The EC algorithm is proposed as a model of fast hippocampal instance-based learning in the brain (Marr, 1971; Sutherland & Rudy, 1989), while the deep RL algorithms described above are more analogous to slow cortical learning that relies on generalised statistical summaries of the input distribution (McClelland et al., 1995; Norman & O'Reilly, 2003; Tulving et al., 1991).

A.7. Source Task Performance Results

The focus of this paper is primarily on zero-shot domain adaptation performance. However, it is also interesting to analyse the effect of the DARLA approach on source domain policy performance. In order to compare the models' behaviour on the source task, we examined the training curves (see Figures 7-10) and noted in particular their:

1. Asymptotic task performance, i.e. the rewards per episode at the point where π_S has converged for the agent under consideration.
2. Data efficiency, i.e. how quickly the training curve was able to achieve convergence.

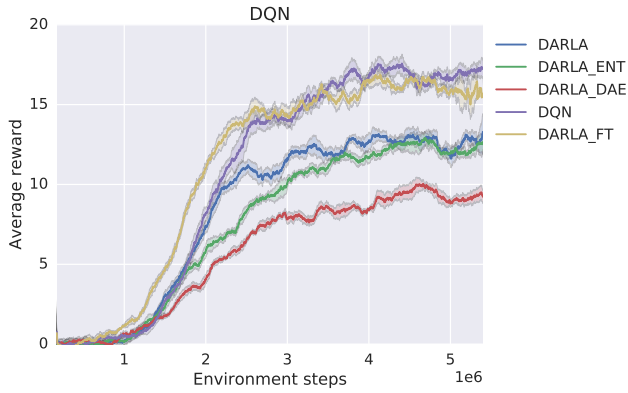


Figure 7. Source task training curves for DQN. Curves show average and standard deviation over 20 random seeds.

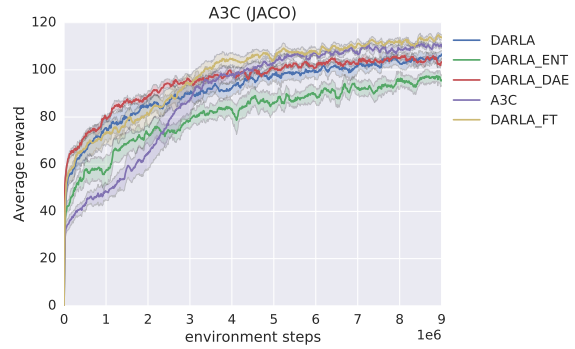


Figure 10. Training curves for various baselines on the source MuJoCo reaching task

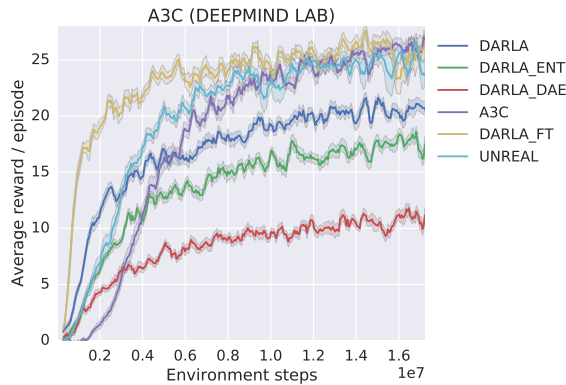


Figure 8. Source task performance training curves for A3C and UNREAL. DARLA shows accelerated learning of the task compared to other architectures. Results show average and standard deviation over 20 random seeds, each using 16 workers.

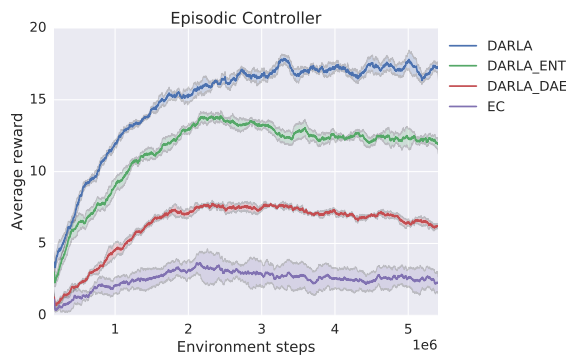


Figure 9. Source task training curves for EC. Results show average and standard deviation over 20 random seeds.

We note the following consistent trends across the results:

1. Using DARLA provided an initial boost in learning performance, which depended on the degree of disentanglement of the representation. This was particularly observable in A3C,

see Fig. 8.

2. Baseline algorithms where \mathcal{F} could be fine-tuned to the source task were able to achieve higher asymptotic performance. This was particularly notable on DQN and A3C (see Figs. 7 and 8) in DeepMind Lab. However, in both those cases, DARLA was able to learn very reasonable policies on the source task which were on the order of 20% lower than the fine-tuned models – arguably a worthwhile sacrifice for a subsequent median 270% improvement in target domain performance noted in the main text.
3. Allowing DARLA to fine-tune its vision module (DARLA_{FT}) boosted its source task learning speed, and allowed the agent to asymptote at the same level as the baseline algorithms. As discussed in the main text, this comes at the cost of significantly reduced domain transfer performance on A3C. For DQN, however, finetuning appears to offer the best of both worlds.
4. Perhaps most relevantly for this paper, even if solely examining source task performance, DARLA outperforms both DARLA_{ENT} and DARLA_{DAE} on both asymptotic performance and data efficiency – suggesting that disentangled representations have wider applicability in RL beyond the zero-shot domain adaptation that is the focus of this paper.