# Adversarial Skill Networks: Unsupervised Robot Skill Learning from Video

Oier Mees*, Markus Merklinger*, Gabriel Kalweit, Wolfram Burgard

*Abstract*— Key challenges for the deployment of reinforcement learning (RL) agents in the real world are the discovery, representation and reuse of skills in the absence of a reward function. To this end, we propose a novel approach to learn a task-agnostic skill embedding space from unlabeled multi-view videos. Our method learns a general skill embedding independently from the task context by using an adversarial loss. We combine a metric learning loss, which utilizes temporal video coherence to learn a state representation, with an entropy-regularized adversarial skill-transfer loss. The metric learning loss learns a disentangled representation by attracting simultaneous viewpoints of the same observations and repelling visually similar frames from temporal neighbors. The adversarial skill-transfer loss enhances re-usability of learned skill embeddings over multiple task domains. We show that the learned embedding enables training of continuous control policies to solve novel tasks that require the interpolation of previously seen skills. Our extensive evaluation with both simulation and real world data demonstrates the effectiveness of our method in learning transferable skills from unlabeled interaction videos and composing them for new tasks. Code, pretrained models and dataset are available at **http://robotskills.cs.uni-freiburg.de**

## I. INTRODUCTION

Intelligent beings have the ability to discover, learn and transfer skills without supervision. Moreover, they can combine previously learned skills to solve new tasks. This stands in contrast to most current "deep reinforcement learning" (RL) methods, which, despite recent progress [1]–[3], typically learn solutions from scratch for every task and often rely on manual, per-task engineering of reward functions. Furthermore, the obtained policies and representations tend to be task-specific and generally do not transfer to new tasks.

The design of reward functions that elicit the desired agent behavior is especially challenging for real-world tasks, particularly when the state of the environment might not be accessible. Additionally, designing a reward often requires the installation of specific sensors to measure as to whether the task has been executed successfully [4], [5]. In many scenarios, the need for task-specific engineering of reward functions prevents us from end-to-end learning from pixels, if the reward function itself requires a dedicated perception pipeline. To address these problems, we propose an unsupervised skill learning method that aims to discover and learn transferable skills by watching videos. The learned embedding is then used to guide an RL-agent in order to
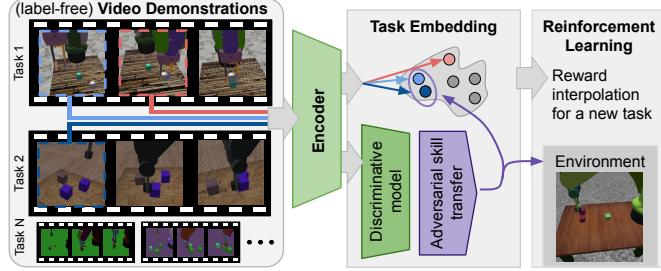


Fig. 1: Given the demonstration of a new task as input, Adversarial Skill Networks yield a distance measure in skill-embedding space which can be used as the reward signal for a reinforcement learning agent for multiple tasks.

solve a wide range of tasks by composing previously seen skills.

Prior work in visual representation learning for deriving reward functions relied on self-supervised objectives [6]–[10] and focused on single tasks. Not only is this inefficient, but also limits the versatility and adaptivity of the systems that can be built. Thus, we consider the problem of learning a multi-skill embedding without human supervision.

In this paper, we present a novel approach called Adversarial Skill Networks (ASN). In order to learn a task-agnostic skill embedding space, our method solely relies on unlabeled multi-view observations. Hence, it does not require correspondences between frames and task IDs nor any additional form of supervision or instrumentation. We combine a metric learning loss, which utilizes temporal video coherence, with an entropy-regularized adversarial skill-transfer loss. Our results indicate that the learned embedding can be used not only to train RL agents for tasks seen during the training of the embedding, but also for novel tasks that require a composition of previously seen skills.

In extensive experiments, we demonstrate both qualitatively and quantitatively that our method learns transferable skill embeddings for simulated and real demonstrations without the requirement of labels. We represent the skill embedding as a latent variable and apply an adversarial entropy regularization technique to ensure that the learned skills are task independent and versatile and that the embedding space is well formed. We show that the learned embedding enables training of continuous control policies with PPO [11] to solve novel tasks that require the interpolation of previously seen skills. Training an RL-agent to re-use skills in an unseen task, by using the learned embedding space as a reward function, solely requires a single video demonstrating the novel task. This makes our method readily applicable in a variety of robotics scenarios.

## II. RELATED WORK

Our work is primarily concerned with learning representations that enable a robot to solve multiple tasks by reusing skills without human supervision, thus falling under the category of self-supervised robot learning [6], [8], [10], [12]. There exists a large body of work for learning representations through autoencoders [13], [14], pre-trained supervised features [10], spatial structure [13], [15] and state estimation from vision [16]. Compared to these approaches, we take multiple tasks into account to learn a skill embedding before training a reinforcement learning agent with a self-supervised vision-based training signal.

Further approaches attempt to derive data-driven reward functions [8]–[10], [17], [18] by providing a label-free training signal from video or images to minimize human supervision. Related to our work, Sermanet *et al.* [10] provide reward functions by identifying key intermediate steps for one task from multiple video examples. Other methods [7], [9], [13], [18] use images of goal examples to construct a task objective for goal reaching tasks such as pushing. Atari video games are solved in [17], [19] by constructing an objective from human demonstration videos. However, it is unclear whether the reward signal reflects a good performance when transferring it to a real-world robotic task. Our model is able to find task specific features and generalizes to unseen objects, viewpoints and backgrounds.

Existing methods for learning reusable skill embeddings make use of entropy-maximization of the policy [20]–[22] and therefore allow for policy interpolation. Haarnoja *et al.* [21] use a composition of soft Q-functions to create a policy that reaches a new goal. Hausmann *et al.* [20] propose a hierarchical reinforcement learning approach that utilizes two embedding networks and an entropy regularization on the policy to cover a latent space with different skill clusters. Orthogonal to our work, these methods rely on previously designed reward functions.

Most related to our approach is the work by Sermanet *et al.* [6] that introduces *Time-Constrative Networks* (TCN) and a triplet loss combined with multi-view metric learning to increase the distance of embeddings for transitions far apart in time. The learned metric can then be used as a reward signal within a RL-setup by minimizing the distance to a visual demonstration. However, TCN focuses only on the single-task setting and does not leverage information from previously learned skills. Dwibedi *et al.* [23] extend TCN using multiple frames (mfTCN). In contrast to our approach, the embedding is not used as a (label-free) reward signal.

In addition to the metric loss, we use an adversarial loss term [24]–[28] as a regularization technique. The adversarial loss was introduced for Generative Adversarial Networks (GAN) [24] and domain adaptation [26]–[28]. Similar to the problem of domain adaptation we have multiple videos for different task domains. For domain adaptation, multiple approaches [26], [27] use a gradient reversal layer and Tzeng *et al.* [28] exploit a GAN-based loss. Springenberg *et al.* [25] introduce an objective function for label free classification by extending GANs to categorical distributions. Our approach uses a similar adversarial loss to learn a reusable skill embedding for different task domains.

In contrast to these previously described approaches, we propose a method to learn skills from video by a composition of metric learning and an entropy-regularized adversarial skill-transfer loss. Our method not only allows for the representation of multiple task-specific reward functions, but also builds upon this information in order to interpolate between learned skills, see Figure 1.

## III. LEARNING A TRANSFERABLE SKILL EMBEDDING

The main incentive of our method is a more general representation of skills that can be re-used and applied to novel tasks. In this work, we define tasks to be composed of a collection of skills. Since we approach this problem in an unsupervised fashion, we do not need any labels describing relations between different task videos or even for different examples of the same task. In our approach, we are interested in the following properties for a learned embedding space:

**i) versatility:** Multiple tasks can be represented in the same embedding space.

**ii) skill representation:** A skill can be described by the embedding of two sequential frames, with a time delay in-between (stride).

**iii) generality:** The learned embedding space should generalize to unseen objects, backgrounds and viewpoints.

**iv) task independent skills:** It should not be possible to distinguish similar skills from different task domains, i.e., the same skill executed in different environments should ideally have an identical embedding.

### A. Adversarial Skill Networks

We propose Adversarial Skill Networks (ASN) to achieve a novel skill representation, which takes these properties into account. We combine a metric learning loss, which utilizes temporal video coherence to learn a state representation, with an entropy-regularized adversarial skill-transfer loss. An overview can be seen in Figure 2.

To transfer similar skills without any label information to unseen tasks, one needs to learn generalized skills that should neither be task- nor domain-specific, following property **iv)**. Since the true class distribution over skills is not known, this problem can naturally be considered as a "soft" probabilistic cluster assignment task.

To solve this, our method introduces a novel entropy regularization by jointly training two networks in an adversarial manner: an encoder network $E$ and a discriminator $D$. Given two sequential frames $(v, w)$, which are separated by a temporal stride $\Delta t$, we define an unlabeled skill embedding $\mathbf{x} = (E(v), E(w))$ and collection of skills $\mathcal{X} = \{\mathbf{x}^1, ..., \mathbf{x}^N\}$ representing the different tasks. The encoder network embeds single frames of the dimension $d_1 \times d_2$ into a lower-dimensional representation of size $n$, i.e. $E: \mathbb{R}^{d_1 \times d_2} \to \mathbb{R}^n$. We compute Euclidean distances in the embedding space to compare the similarity of frames. The discriminator network takes two concatenated embedded frames that define an
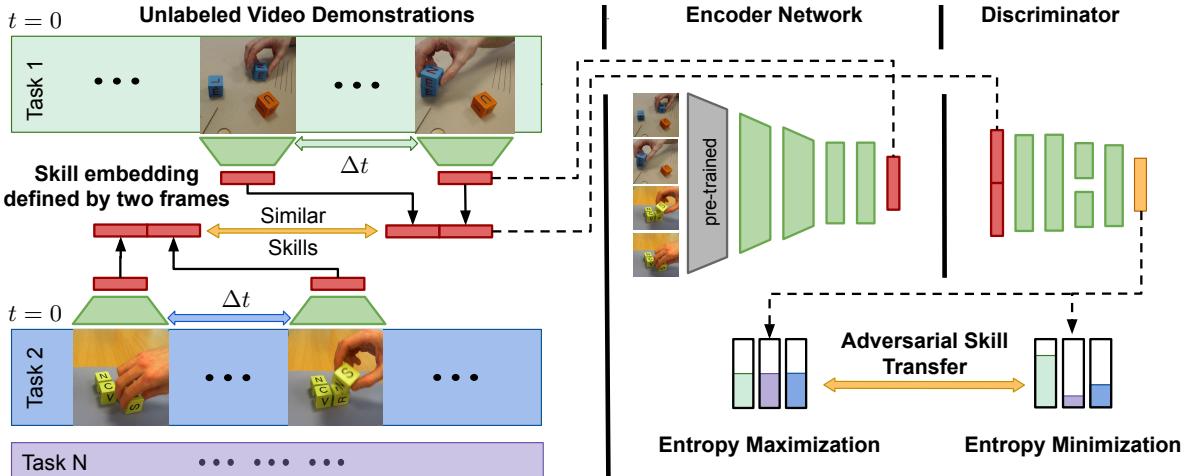
Fig. 2: Structure of Adversarial Skill Networks: We learn a skill metric space in an adversarial framework. The encoding part of the network tries to maximize the entropy to enforce generality. The discriminator, which is not used at test time, tries to minimize the entropy of its prediction to improve recognition of the skills. Finally, maximizing the marginal class entropy over all skills leads to uniform usage of all task classes. Please note that no information about the relation between frames and the tasks they originated from is needed.

unlabeled skill $\mathbf{x}$ as input and outputs $y_c$, the probability of the skill being originated from task $c$. Formally, we require $D(\mathbf{x}) \in \mathbb{R}^C$ to give rise to a conditional distribution over tasks $\sum_{c=1}^{C} p(y_c = c \mid \mathbf{x}, D) = 1$. Although we define this hyper-parameter a priori as the number of tasks contained in a training set, we observed minor performance drops setting it to a value with small deviation from the true number of tasks. Most importantly, ASN does not need a task label for the demonstration videos.

The encoder parameters are updated using a metric learning loss and maximization of the entropy of the discriminator output. In order to capture the temporal task information, we use a modified version of the lifted structure loss [29]. Given two view-pairs $(v_1, v_2)$, synchronized videos from different perspectives, we attract frames that represent the same temporal task state and repulse temporal neighbors, given a constant margin $\lambda$, i.e.

$$
\mathcal{L}_{\text{lifted asn}} = \sum_{i=1}^{M} \left( \log \sum_{y_k = y_i} \left( e^{\lambda - S_{ik}} + \mathbb{1}_{S_{ik} > \xi} \cdot S_{ik} \right) \right.
$$
$$
\left. + \log \sum_{y_k \neq y_i} e^{S_{ik}} \right), \quad (1)
$$

for $M$ frames $(x_1, x_2 \ldots x_M)$ and $S_{ij} = E(x_i) \cdot E(x_j)$, as a dense squared pairwise similarity distance matrix of the batch and $\xi$ a similarity threshold. Additionally, we introduce a constraint that bounds the distance between two positive view-pairs. This constraint is tailored to account for high variance in the learned distance metric. By penalizing large distances of positive view-pairs, we aim at smoother transitions between similar states in a RL setting.

The discriminator network minimizes the entropy given an unlabeled skill embedding $\mathbf{x}$ to be certain about which task $C$ the skill originated from. Note that the discriminator

$D$ is utilized only during training. Without any additional label information about the $C$ classes, we cannot directly specify which class probability $p(y_c = c \mid \mathbf{x}^i, D)$ should be maximized for any given skill $\mathbf{x}$. We make use of information theoretic measures on the predicted class distribution to group the unlabeled skills into well separated categories in the skill embedding space without explicitly modeling $p(\mathbf{x})$. Specifically, if we want the discriminator to be certain for the class distribution $p(y_c = c \mid \mathbf{x}^i, D)$, this corresponds to minimizing the Shannon information entropy $H[p(y_c \mid \mathbf{x}, D)]$, as any draw from this distribution should most of the times result in the same class. On the other hand, if we want the encoder to learn generalized skill representations to meet requirement **iv)**, it should be uncertain of how to classify the unlabeled skills. Thus, the encoder tries to maximize the entropy $H[p(y_c \mid \mathbf{x}, D)]$, which at the optimum will result in a uniform conditional distribution over task classes. Concretely, we define the empirical estimate of the conditional entropy over embedded skill examples $\mathcal{X}$ as:

$$
\mathbb{E}_{\mathbf{x} \sim \mathcal{X}} \left[ H[p(y_c \mid \mathbf{x}, D)] \right] = \frac{1}{N} \sum_{i=1}^{N} H \left[ p(y_c \mid \mathbf{x}^i, D) \right]
$$
$$
= \frac{1}{N} \sum_{i=1}^{N} - \sum_{c=1}^{C} p(y_c = c \mid \mathbf{x}^i, D) \log p(y_c = c \mid \mathbf{x}^i, D).
$$
$$(2)$$

With an additional regularizer we enhance the equal usage of all task classes, corresponding to maximizing a uniform marginal distribution:

$$
H_{\mathbf{x}}[p(y_c \mid D)] = H \left[ \frac{1}{M} \sum_{i=1}^{M} p(y_c \mid \mathbf{x}^i, D) \right], \quad (3)
$$

where $M$ is set to the number of independently drawn samples [25].

In order to disentangle the learned metric and the mapping

to task IDs, we add a sampled latent variable $z = \mu + \sigma \odot \mathcal{E}$ and $\mathcal{E} \sim N(0,1)$, where $D$ estimates $\mu, \sigma$ of a Gaussian distribution. We use the re-parameterization trick to back-propagate through the random node [30]. With the Kullback-Leibler divergence regularization for $z$ we force $D$ to find similar properties describing the skills. Without this objective, similar skills could end up represented far away from each other in the skill embedding space.

This leads to the following objectives for the encoder $E$ and discriminator $D$:

$$\mathcal{L}_{KL} = D_{KL}[p(z \mid \mathbf{x})||p(z)],$$
$$\mathcal{L}_D = -H_{\mathbf{x}}\Big[p(y_c \mid D)\Big] + \mathbb{E}_{\mathbf{x} \sim \mathcal{X}}\Big[H[p(y_c \mid \mathbf{x}, D)]\Big]$$
$$+ \beta\mathcal{L}_{KL} \text{ and} \qquad (4)$$
$$\mathcal{L}_E = H_{\mathbf{x}}\Big[p(y_c \mid D)\Big] + \mathbb{E}_{\mathbf{x} \sim \mathcal{X}}\Big[H[p(y_c \mid \mathbf{x}, D)]\Big]$$
$$- \alpha\mathcal{L}_{\text{lifted asn}}.$$

We therefore optimize the discriminator and the encoder according to:

$$\min_D \mathcal{L}_D \qquad (5)$$

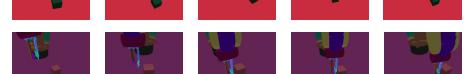and

$$\max_E \mathcal{L}_E. \qquad (6)$$

### B. Implementation Details

The encoder network is inspired by Time-Contrastive Networks (TCN) [6]. We use an Inception network as a feature extractor [31], which is initialized with ImageNet pre-trained weights. The feature extractor is followed by two convolutional layers and a spatial softmax layer for dimension reduction. Finally, after a Fully Connected (FC) layer, the model outputs the embedding vector for a frame. For all experiments, we use $\alpha = 0.1$, $\beta = 1.0$, $\lambda = 1.0$ and an embedding size of 32. The discriminator consists of two FC layers to estimate $\mu$ and $\sigma$ of a Gaussian distribution, followed by two layers to output the task ID. We use dropout for regularization.

We train the encoder and discriminator networks with the Adam optimizer and a learning rate of 0.001. A training batch contains 32 frames from $n = 4$ different view pairs. We load real-world data from video files and sample the simulated data from uncompressed image files. Training directly on images, ensures that our model is not learning any bias introduced by video compression techniques. For frames from the training set, we randomly change brightness, contrast and saturation and randomly mirror frames horizontally. For real-world data, additional training frames are cropped randomly. We train on images of the size $299 \times 299 \times 3$ pixels. For simulated data the discriminator network is only updated with successful task demonstrations, since only they contain the skills we want to transfer. After data augmentation, the frames of a batch are normalized on each RGB channel using the $\mu$ and $\sigma$ of the ImageNet dataset.



(a) Real block tasks



(b) Simulated block tasks

Fig. 3: Visualization of the multi-task datasets used in this work.

## IV. EXPERIMENTAL RESULTS

We evaluate the performance of our ASN model on two data sets, see Figure 3.

The first data set consists of three simulated robot tasks: stacking (A), color pushing (B) and color stacking (C). The data set contains 300 multi-view demonstration videos per task. The tasks are simulated with PyBullet. Of these 300 demonstrations, 150 represent unsuccessful executions of the different tasks. We found it helpful to add unsuccessful demonstrations in the training of the embedding to enable training RL agents on it. Without fake examples, the distances in the embedding space for states not seen during training might be noisy. In the initial phase of training, however, the policy to be learned mostly visits areas of the state-action space which are not covered by the (successful) demonstration. Hence, it is important to have unsuccessful examples in the training set. The test set contains the manipulation of blocks. Within the validation set, the blocks are replaced by cylinders of different colors.

The second data set includes real-world human executions of the simulated robot tasks (A, B and C), as well as demonstrations for a task where one has to first separate blocks in order to stack them (D). Each task contains 60 multi-view demonstration videos, corresponding to 24 minutes of interaction. The test set contains blocks of unseen sizes and textures, as well as unknown backgrounds, in order to evaluate the generality of our approach.
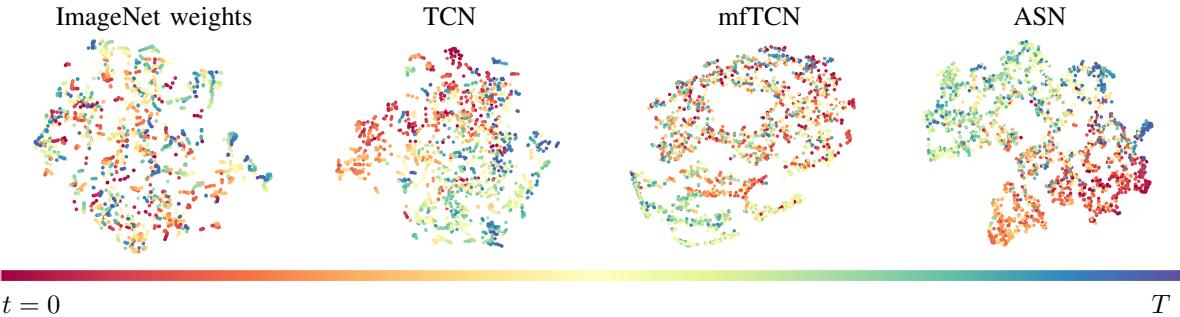
Fig. 4: t-SNE of an **unseen** color stacking video for models trained on block stacking and color pushing. Our ASN model maintains the temporal coherence of the task better than the baselines. The colorbar indicates the temporal task progress.

## A. Quantitative Evaluation

We measure the performance of our models based on the alignment loss, following Sermanet *et al.* [6], to determine how well two views of a video are aligned in time. We take advantage of the fact that frames in both videos are synchronized with each other to get alignment labels for free. We try to sequentially align two view pairs by finding the nearest neighbor in the embedding space normalizing the distances by the demonstration length. After embedding the two videos in our skill embedding space, we search for each frame $t_j^i$ in the first video $i$ the nearest neighbour in the second view and retrieve its time index $t_j^{nn}$. Thus, for video $i$ the aligment loss is defined by:

$$\text{align}_i = \frac{\sum_{j=1}^{F} |t_j^i - t_j^{nn}|}{F}, \quad (7)$$

for a video of length $F$. As Sermanet *et al.* [6] have shown, the alignment loss reflects the quality of the reward signal within a RL setup. Instead of evaluating the alignment loss on the same task as the embedding was trained on, we measure how well view pairs of novel, unseen tasks are aligned. This form of zero-shot evaluation is very challenging, as it requires the combination of previously seen skills. Adversarial Skill Networks yield the best performance for transfer in the simulated robot setting, which can be seen in Table I. Please note the lower bound (0.081) of alignment loss for single-task TCN. In contrast to TCN trained on multiple tasks, our model gets very close (0.099) despite not being trained on the task. Furthermore, our approach outperforms TCN in both the real robot multi-task setup, as well as in the transfer task, which is depicted in Table II. We also compare against the different metric learning losses and show that the lifted loss in combination with the bound for positive view-pairs outperforms other methods.

| Model | Task Combination, Train → Test | | |
| --- | --- | --- | --- |
| | C → C | A,B,C → A,B,C | A,B → C |
| Inception-ImageNet [31] | 0.29 | 0.31 | 0.29 |
| TCN - lifted [6] | 0.081 | 0.058 | 0.112 |
| ASN | - | 0.056 | **0.099** |

TABLE I: Test alignment loss for the simulated robot multi-task dataset, which includes fake examples, Tasks: A: 2 block stack, B: 3 block color push sort, C: 3 block color stack sort.

A visualization of the learned embedding space is depicted in Figure 4. ASN can represent the temporal relations of an unseen task better than TCN or multi-frame TCN, leading to more meaningful distance measures in embedding space. Our proposed multi-task setup is able to reflect skills needed to solve the unseen task and thus can generalize better, while maintaining the temporal coherence of the unseen task.

| Model | Task Combination, Train → Test | | |
| --- | --- | --- | --- |
| | A,B,C → C | A,B → C | A,B,D → C |
| TCN - triplet [6] | 0.186 | 0.21 | 0.218 |
| TCN - lifted [6] | 0.171 | 0.20 | 0.187 |
| TCN - npair [6] | 0.221 | 0.209 | 0.221 |
| mfTCN - lifted [23] | 0.174 | 0.23 | 0.22 |
| ASN - normal lifted | 0.168 | 0.183 | 0.181 |
| ASN | **0.150** | **0.180** | **0.165** |

TABLE II: Test alignment loss real-world block tasks, Tasks: A: 2 block stack, B: 3 block color push sort, C: 3 block color stack sort, D: 4 block separate to stack.

## B. Ablation studies

To analyze the influence of our different building blocks on the learned embedding, we conducted several experiments, see Table III. Our results indicate that it is of benefit to describe a skill with a growing stride, so as to cover macro-actions describing events longer in time. In order to keep the embeddings of these skill frames of higher stride aligned, the KL-divergence is shown to be an effective regularization technique, yielding the lowest alignment loss. Furthermore, it seems to be enough to describe a skill by only the start and end frames. A single frame seems to provide too little information whereas using four frames proves to make the state space too high dimensional.

| Regularization | #Domain frames | Stride | Real block tasks A,B,D → C |
| --- | --- | --- | --- |
| KL | 1 | - | 0.187 |
| KL | 4 | 5 | 0.185 |
| KL | 2 | 5 | 0.168 |
| KL | 2 | 15 | **0.165** |
| FC | 2 | 15 | 0.1987 |
| KL w/o encoder entropy | 2 | 15 | 0.186 |
| KL w/o entropy | 2 | 15 | 0.177 |

TABLE III: Ablation studies: transfer loss for different regularization techniques and skill definitions.
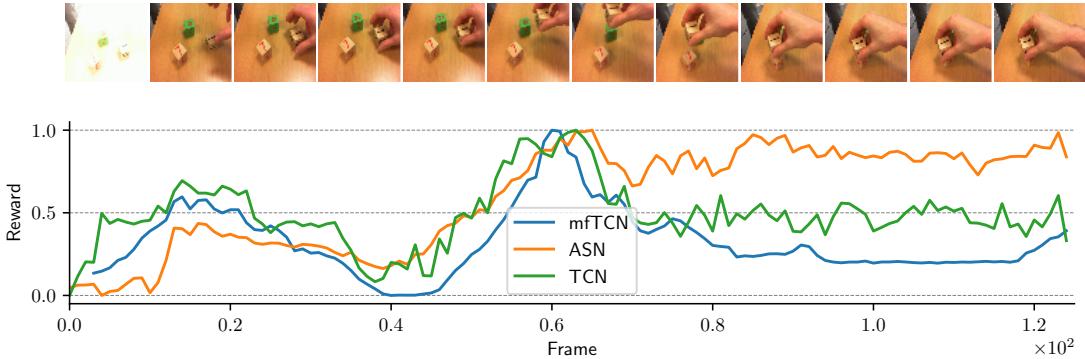
Fig. 5: Reward plot for a novel color stacking task C for models trained on tasks A, B and D. The reward is based on the distance to a single goal frame from a different perspective.

## C. Learning Control Policies

Lastly, we integrate the learned metric within a RL-agent to imitate an unseen task given a single video demonstration. Concretely, for learning a continuous control policy on the color stacking (C) task we train the embedding on the tasks of two block stacking (A) and color pushing (B). Thus, successfully imitating the previously never seen color stacking task requires the interpolation of previously seen skills. Additionally, we also learn a continuous control policy for an unseen color pushing task, given an embedding trained on stacking and color stacking. To train the agents, we use the distance measure in embedding space of the agent view $v_a^t$ and the demonstration frame $v_d^t$ for timestep $t$ as the reward signal for the on-policy optimization algorithm PPO [11]:

$$r^{(t)} = \begin{cases} 10 - d\Big(E(v_a^t), E(v_d^t)\Big) & \text{if } d\Big(E(v_a^t), E(v_d^t)\Big) < \xi \\ 0 & \text{otherwise,} \end{cases}$$
(8)

where $d$ is the euclidean distance and $\xi$ a constant threshold. The agent state consists of the embedding $E(v_a^t)$ and the joint angle of the robot. We train the policy with Adam and a learning rate of $10^{-5}$ and a batch size of 32. To alleviate the problem of exploration and to focus on the quality of reward signal, we take random samples along the given demonstration as initial states and reset the environment if the end effector vastly differs from the demonstration, following Peng *et al.* [32]. The results are depicted in Figure 6.
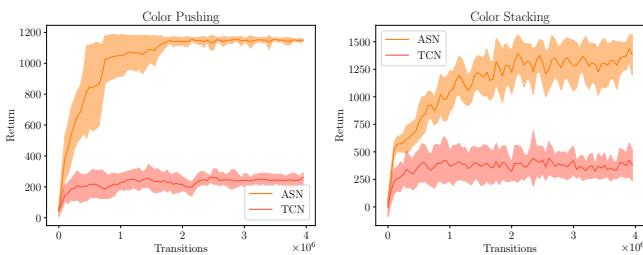


Fig. 6: Results for training a continuous control policy with PPO on the unseen Color Pushing and Color Stacking tasks with the learned reward function. The plot shows mean and standard deviation over five training runs.

Our method succeeds in solving the tasks, whereas the baseline TCN approach converges to a local minima. This demonstrates the effectiveness of our approach in reusing skills for a novel task given a single video demonstration. Please note that training RL agents on tasks which have never been shown during the training of the embedding is very challenging, as it requires the discovery and reuse of task-independent skills.

Additionally, we evaluate the reward signal on an unseen color stacking task (C) for the real-world dataset. We plot a reward signal, which is based on the distance measurement of the task state for each timestep and a single goal frame from a different perspective, see Figure 5. The embedding of all models are trained on tasks A, B and D. To compare the different models we normalize the negative distance outputs for timestep between zero and one. The baseline model already give a similar reward for many initials states and goal states, despite the states being visually different. Our model shows a continuous and incremental reward as the task progresses and saturates as it is completed.

## V. CONCLUSION

We proposed Adversarial Skill Networks, a model to leverage information from multiple label-free demonstrations in order to yield a meaningful embedding for unseen tasks. We showed that our approach is able to reuse learned skills for compositions of tasks and achieves state-of-the-art performance. We demonstrate that the learned embedding enables training of continuous control policies to solve novel tasks that require the interpolation of previously seen skills. Our results show that our model can find a good embedding for vastly different task domains. This is a first step towards discovery, representation and reuse of skills in the absence of a reward function.

Going forward, a natural extension of this work is the application of the learned distance metric in a real-world re-inforcement learning setting and in environments that require a higher degree of interpolation for successful completion. Another promising direction for future work is the evaluation of the proposed approach in a sim-to-real setup [33].

## REFERENCES

[1] Matteo Hessel, Joseph Modayil, Hado Van Hasselt, Tom Schaul, Georg Ostrovski, Will Dabney, Dan Horgan, Bilal Piot, Mohammad Azar, and David Silver. Rainbow: Combining improvements in deep reinforcement learning. In *AAAI Conference on Artificial Intelligence*, 2018.

[2] Gabriel Barth-Maron, Matthew W Hoffman, David Budden, Will Dabney, Dan Horgan, Alistair Muldal, Nicolas Heess, and Timothy Lillicrap. Distributed distributional deterministic policy gradients. In *International Conference on Learning Representations, ICLR*, 2018.

[3] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dharshan Kumaran, Thore Graepel, et al. A general reinforcement learning algorithm that masters chess, shogi, and go through self-play. *Science*, 362, 2018.

[4] Connor Schenck and Dieter Fox. Visual closed-loop control for pouring liquids. In *IEEE International Conference on Robotics and Automation, ICRA*, 2017.

[5] Andrei A Rusu, Matej Večerík, Thomas Rothörl, Nicolas Heess, Razvan Pascanu, and Raia Hadsell. Sim-to-real robot learning from pixels with progressive nets. In *Conference on Robot Learning, CoRL*, 2017.

[6] Pierre Sermanet, Corey Lynch, Yevgen Chebotar, Jasmine Hsu, Eric Jang, Stefan Schaal, and Sergey Levine. Time-contrastive networks: Self-supervised learning from video. *Proceedings of International Conference in Robotics and Automation, ICRA*, 2018.

[7] Manuel Watter, Jost Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. In *Advances in neural information processing systems, NeurIPS*, 2015.

[8] Chelsea Finn and Sergey Levine. Deep visual foresight for planning robot motion. In *IEEE International Conference on Robotics and Automation, ICRA*, 2017.

[9] Tianhe Yu, Gleb Shevchuk, Dorsa Sadigh, and Chelsea Finn. Unsupervised visuomotor control through distributional planning networks. *Proceedings of Robotics: Science and Systems, RSS*, 2019.

[10] Pierre Sermanet, Kelvin Xu, and Sergey Levine. Unsupervised perceptual rewards for imitation learning. *Proceedings of Robotics: Science and Systems, RSS*, 2017.

[11] John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.

[12] Oier Mees, Maxim Tatarchenko, Thomas Brox, and Wolfram Burgard. Self-supervised 3d shape and viewpoint estimation from single images for robotics. In *Proceedings of the International Conference on Intelligent Robots and Systems (IROS)*, Macao, China, 2019.

[13] Chelsea Finn, Xin Yu Tan, Yan Duan, Trevor Darrell, Sergey Levine, and Pieter Abbeel. Deep spatial autoencoders for visuomotor learning. In *IEEE International Conference on Robotics and Automation, ICRA*, 2016.

[14] Christopher P Burgess, Loic Matthey, Nicholas Watters, Rishabh Kabra, Irina Higgins, Matt Botvinick, and Alexander Lerchner. Monet: Unsupervised scene decomposition and representation. *arXiv preprint arXiv:1901.11390*, 2019.

[15] Rico Jonschkowski, Roland Hafner, Jonathan Scholz, and Martin Riedmiller. Pves: Position-velocity encoders for unsupervised learning of structured state representations. *arXiv preprint arXiv:1705.09805*, 2017.

[16] Marcin Andrychowicz, Bowen Baker, Maciek Chociej, Rafal Jozefowicz, Bob McGrew, Jakub Pachocki, Arthur Petron, Matthias Plappert, Glenn Powell, Alex Ray, et al. Learning dexterous in-hand manipulation. *arXiv preprint arXiv:1808.00177*, 2018.

[17] Yusuf Aytar, Tobias Pfaff, David Budden, Thomas Paine, Ziyu Wang, and Nando de Freitas. Playing hard exploration games by watching youtube. In *Advances in Neural Information Processing Systems, NeurIPS*, 2018.

[18] Avi Singh, Larry Yang, Kristian Hartikainen, Chelsea Finn, and Sergey Levine. End-to-end robotic reinforcement learning without reward engineering. *Proceedings of Robotics: Science and Systems, RSS*, 2019.

[19] Borja Ibarz, Jan Leike, Tobias Pohlen, Geoffrey Irving, Shane Legg, and Dario Amodei. Reward learning from human preferences and demonstrations in atari. In *Advances in Neural Information Processing Systems, NeurIPS*, 2018.

[20] Karol Hausman, Jost Tobias Springenberg, Ziyu Wang, Nicolas Heess, and Martin Riedmiller. Learning an embedding space for transferable robot skills. In *International Conference on Learning Representations, ICLR*, 2018.

[21] Tuomas Haarnoja, Vitchyr Pong, Aurick Zhou, Murtaza Dalal, Pieter Abbeel, and Sergey Levine. Composable deep reinforcement learning for robotic manipulation. In *IEEE International Conference on Robotics and Automation, ICRA*, 2018.

[22] Domingo Esteban, Leonel Rozo, and Darwin G Caldwell. Hierarchical reinforcement learning for concurrent discovery of compound and composable policies. *arXiv preprint arXiv:1905.09668*, 2019.

[23] Debidatta Dwibedi, Jonathan Tompson, Corey Lynch, and Pierre Sermanet. Learning actionable representations from visual observations. In *IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS*, 2018.

[24] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, 2014.

[25] Jost Tobias Springenberg. Unsupervised and semi-supervised learning with categorical generative adversarial networks. In *International Conference on Learning Representations, ICLR*, 2016.

[26] Yaroslav Ganin and Victor Lempitsky. Unsupervised domain adaptation by backpropagation. In *International Conference on Machine Learning, ICML*, 2015.

[27] Mingsheng Long, Zhangjie Cao, Jianmin Wang, and Michael I Jordan. Conditional adversarial domain adaptation. In *Advances in Neural Information Processing Systems, NeurIPS*, 2018.

[28] Eric Tzeng, Judy Hoffman, Kate Saenko, and Trevor Darrell. Adversarial discriminative domain adaptation. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2017.

[29] Hyun Oh Song, Yu Xiang, Stefanie Jegelka, and Silvio Savarese. Deep metric learning via lifted structured feature embedding. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR*, 2016.

[30] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations, ICLR*, 2014.

[31] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jonathon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *CVPR*, 2016.

[32] Xue Bin Peng, Pieter Abbeel, Sergey Levine, and Michiel van de Panne. Deepmimic: Example-guided deep reinforcement learning of physics-based character skills. *ACM Transactions on Graphics (TOG)*, 37(4):143, 2018.

[33] Lukas Hermann, Max Argus, Andreas Eitel, Artemij Amiranashvili, Wolfram Burgard, and Thomas Brox. Adaptive curriculum generation from demonstrations for sim-to-real visuomotor control. *arXiv preprint arXiv:1910.07972*, 2019.