

---

# Learning Options via Compression

---

**Yiding Jiang\***

Carnegie Mellon University  
yidingji@cs.cmu.edu

**Evan Zheran Liu\***

Stanford University  
evanliu@cs.stanford.edu

**Benjamin Eysenbach**

Carnegie Mellon University  
beysenba@cs.cmu.edu

**J. Zico Kolter**

Carnegie Mellon University  
zkolter@cs.cmu.edu

**Chelsea Finn**

Stanford University  
cbfinn@cs.stanford.edu

## Abstract

Identifying statistical regularities in solutions to some tasks in multi-task reinforcement learning can accelerate the learning of new tasks. Skill learning offers one way of identifying these regularities by decomposing pre-collected experiences into a sequence of skills. A popular approach to skill learning is maximizing the likelihood of the pre-collected experience with latent variable models, where the latent variables represent the skills. However, there are often many solutions that maximize the likelihood equally well, including degenerate solutions. To address this underspecification, we propose a new objective that combines the maximum likelihood objective with a penalty on the description length of the skills. This penalty incentivizes the skills to maximally extract common structures from the experiences. Empirically, our objective learns skills that solve downstream tasks in fewer samples compared to skills learned from only maximizing likelihood. Further, while most prior works in the offline multi-task setting focus on tasks with low-dimensional observations, our objective can scale to challenging tasks with high-dimensional image observations.

## 1 Introduction

While learning tasks from scratch with reinforcement learning (RL) is often sample inefficient [30, 8], leveraging datasets of pre-collected experience from various tasks can accelerate the learning of new tasks. This experience can be used in numerous ways, including to learn a dynamics model [18, 38], to learn compact representations of observations [87], to learn behavioral priors [77], or to extract meaningful skills or options [79, 42, 2, 59] for hierarchical reinforcement learning (HRL) [6]. Our work studies this last approach. The central idea is to extract commonly occurring behaviors from the pre-collected experience as skills, which can then be used in place of primitive low-level actions to accelerate learning new tasks via planning [70, 51] or RL [55, 47, 52]. For example, in a navigation domain, learned skills may correspond to navigating to different rooms or objects.

Prior methods learn skills by maximizing the likelihood of the pre-collected experience [43, 42, 2, 90]. However, this maximum likelihood objective (or the lower bounds on it) is underspecified: it often admits many solutions, only some of which would help learning new tasks (see Figure 1). For example, one degenerate solution is to learn a single skill for each entire trajectory; another degenerate solution is to learn skills that operate for a single timestep, equivalent to the original action space. Both of these decompositions can perfectly reconstruct the pre-collected experiences (and, hence, maximize likelihood), but they are of limited use for learning to solve new tasks. Overall, the maximum likelihood objective cannot distinguish between such decompositions and potentially more useful

---

\*Equal contribution

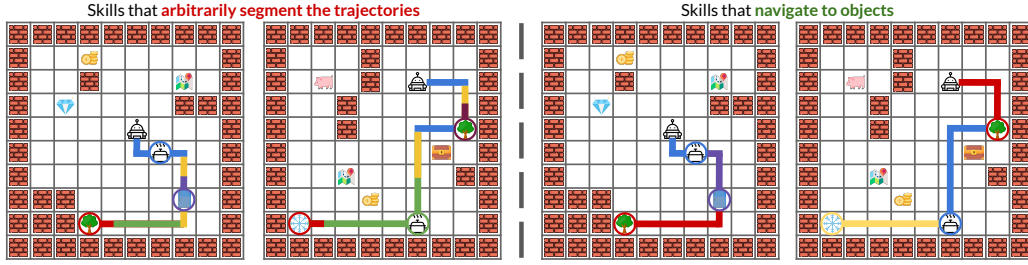


Figure 1: Skill learning via maximizing likelihood is *underspecified*: the skills (represented by the different colored segments) on both the left and the right maximize the likelihood (i.e., encode the data), but the skills on the right are more likely to be useful for a new task. Our approach factors out common temporal structures, yielding the skills on the right. The visualized tasks are from Kipf et al. [42].

**decompositions**, and we find that this underspecification problem can empirically occur even on simple tasks in our experiments.

To address this underspecification problem, we propose to introduce an additional objective which biases skill learning to acquire skills that can accelerate the learning of new tasks. To construct such an objective, we make two key observations:

- Skills must extract reusable temporally-extended structure from the pre-collected experiences to accelerate learning of new tasks.
- *Compressing* data also requires identifying and extracting common structure.

Hence, we hypothesize that compression is a good objective for skill learning. We formalize this notion of compression with the principle of *minimum description length* (MDL) [69]. Concretely, we combine the maximum likelihood objective (used in prior works) with a **new term to minimize the number of bits required to represent the pre-collected experience with skills**, which incentivizes the skills to find common structure. Additionally, while **prior compression-based methods typically involve discrete optimization and hence are not differentiable**, we also provide a method to minimize our compression objective via gradient descent (Section 4), which enables optimizing our objective with neural networks.

Overall, our main contribution is an objective for extracting skills from offline experience, which addresses the underspecification problem with the maximum likelihood objective. We call the resulting approach LOVE (Learning Options Via comprESSION). Using multi-task benchmarks from prior work [42], we find that LOVE can learn skills that enable faster RL and are more semantically meaningful compared to skills learned with prior methods. We also find that LOVE can scale to tasks with high-dimensional image observations, whereas most prior works focus on tasks with low-dimensional observations.

## 2 Related Works

We study the problem of using offline experience from one set of tasks to quickly learn new tasks. While we use the experience to learn skills, prior works also consider other approaches of leveraging the offline experience, including to learn a dynamics model [38], to **learn a compact representation of observations** [87], and to learn behavioral priors for exploring new tasks [77].

We build on a rich literature on extracting skills [79] from *offline* experience [62, 21, 43, 74, 42, 7, 66, 92, 91, 46, 2, 72, 71, 59, 90, 53, 93, 68, 80, 88]. These works predominantly learn skills using a latent variable model, **where the latent variables partition the experience into skills**, and the overall model is learned by maximizing (a lower bound on) the likelihood of the experiences. This approach is structurally similar to a long line of work in hierarchical Bayesian modeling [24, 9, 20, 49, 37, 13, 48, 39, 28, 15]. We also follow this approach and build off of VTA [39] for the latent variable model, but differ by introducing a new compression term to address the underspecification problem with maximizing likelihood. **Several prior approaches [90, 91, 23] also use compression, but yield open-loop skills, whereas we aim to learn closed-loop skills that can react to the state.** Additionally, Zhang et al. [90] show that maximizing likelihood with variational inference can itself be seen as a form of compression, but this work compresses the latent variables at each time step, which does not in general yield optimal compression as our objective does.

Beyond learning skills from offline experience, prior works also consider learning skills with additional supervision [3, 63, 75, 35, 78], from online interaction without reward labels [25, 19, 17, 73, 5], and from online interaction with reward labels [81, 44, 4, 84, 29, 58]. Additionally, a large body of work on meta-reinforcement learning [16, 85, 67, 94, 50] also leverages prior experience to quickly learn new tasks, though not necessarily through skill learning.

### 3 Preliminaries

**Problem setting.** We consider the problem of using an offline dataset of experience to quickly solve new RL tasks, where each task is identified by a reward function. The offline dataset is a set of trajectories  $\mathcal{D} = \{\tau_i\}_{i=1}^N$ , where each trajectory is a sequence of states  $\mathbf{x} \in \mathcal{X}$  and actions  $\mathbf{a} \in \mathcal{A}$ :  $\tau_i = \{(\mathbf{x}_1, \mathbf{a}_1), (\mathbf{x}_2, \mathbf{a}_2), \dots\}$ . Each trajectory is collected by some unknown policy interacting with a Markov decision process (MDP) with dynamics  $\mathcal{P}(\mathbf{x}_{t+1} | \mathbf{x}_t, \mathbf{a}_t)$ . Following prior work [18, 1, 14], we do not assume access to the rewards (i.e., the task) of the offline trajectories.

Using this dataset, our aim is to quickly solve a new task. The new task involves interacting with the same MDP as the data collection policy in the offline dataset, except with a new reward function  $\mathcal{R}(\mathbf{x}, \mathbf{a})$ . The objective is to learn a policy that maximizes the expected returns in as few numbers of environment interactions as possible.

**Variational Temporal Abstraction.** Our method builds upon the graphical model from variational temporal abstraction (VTA) [39], which is a method for decomposing non-control sequential data (i.e., data without actions) into subsequences. We overview VTA below and extend it to handle actions in Section 4.1.

VTA assumes that each observation  $\mathbf{x}_t$  is associated with an abstract representation  $\mathbf{s}_t$  and a subsequence descriptor  $\mathbf{z}_t$ . An additional, binary random variable  $\mathbf{m}_t$  indicates whether a new subsequence begins at the current observation. VTA uses the following generative model:

1. Determine whether to begin a new subsequence:  $\mathbf{m}_t \sim p(\mathbf{m}_t | \mathbf{s}_{t-1})$ . A new subsequence always begins on the first time step, i.e.,  $\mathbf{m}_1 = 1$ .
2. Sample the descriptor:  $\mathbf{z}_t \sim p(\mathbf{z}_t | \mathbf{z}_{<t}, \mathbf{m}_t)$ . If  $\mathbf{m}_t = 0$ , the previous descriptor is copied, i.e.,  $\mathbf{z}_t = \mathbf{z}_{t-1}$ .
3. Sample the next state abstract representation:  $\mathbf{s}_t \sim p(\mathbf{s}_t | \mathbf{s}_{<t}, \mathbf{z}_t, \mathbf{m}_t)$ .
4. Sample the observation:  $\mathbf{x}_t \sim p(\mathbf{x}_t | \mathbf{s}_t)$ .

We visualize this sampling procedure in Figure 2. Formally, we can write this generative model as:

$$p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}, \mathbf{s}_{1:T}, \mathbf{m}_{1:T}) = \prod_{t=1}^T p(\mathbf{x}_t | \mathbf{s}_t) p(\mathbf{m}_t | \mathbf{s}_{t-1}) p(\mathbf{s}_t | \mathbf{s}_{<t}, \mathbf{z}_t, \mathbf{m}_t) p(\mathbf{z}_t | \mathbf{z}_{<t}, \mathbf{m}_t).$$

VTA maximizes the likelihood of the observed data under this latent variable model using the standard evidence lower bound (ELBO):

$$\log p(\mathbf{x}_{1:T}) \geq \sum_{\mathbf{m}_{1:T}} \int_{\mathbf{s}_{1:T}, \mathbf{z}_{1:T}} q_{\theta}(\mathbf{z}_{1:T}, \mathbf{s}_{1:T}, \mathbf{m}_{1:T} | \mathbf{x}_{1:T}) \cdot \log \frac{p_{\theta}(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}, \mathbf{s}_{1:T}, \mathbf{m}_{1:T})}{q_{\theta}(\mathbf{z}_{1:T}, \mathbf{s}_{1:T}, \mathbf{m}_{1:T} | \mathbf{x}_{1:T})}$$

This lower bound holds for any choice of  $q_{\theta}$ . VTA chooses to factor this distribution as:

$$q_{\theta}(\mathbf{z}_{1:T}, \mathbf{s}_{1:T}, \mathbf{m}_{1:T} | \mathbf{x}_{1:T}) = \prod_{t=1}^T q_{\theta}(\mathbf{m}_t | \mathbf{x}_{1:t}) q_{\theta}(\mathbf{z}_t | \mathbf{z}_{t-1}, \mathbf{m}_t, \mathbf{x}_{1:T}) q_{\theta}(\mathbf{s}_t | \mathbf{z}_t, \mathbf{m}_t, \mathbf{x}_{1:t}).$$

As discussed above and as we observe in Section 6, the maximum likelihood objective is under-specified and may yield degenerate or unhelpful solutions, which we address in the next section. Specifically, we use the graphical model from VTA, but introduce a new objective that compresses the latent variables  $\mathbf{z}$ , while also maximizing the ELBO.

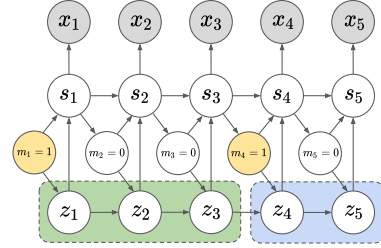


Figure 2: The VTA graphical model (adapted from Kim et al. [39]). The model decomposes the data  $\mathbf{x}_{1:T}$  into subsequences demarcated by when the boundary variables  $\mathbf{m}_t = 1$  (highlighted in yellow). Each subsequence is assigned a descriptor  $\mathbf{z}$ , where all  $\mathbf{z}_t$  within a subsequence are the same (outlined in dashed lines).

## 4 LOVE: Learning Options via Compression

In this section, we describe our method for learning skills from pre-collected experience. We first extend the VTA graphical model to handle experience labeled with actions, and then introduce a compression objective that encourages extracting common structure from the experience.

### 4.1 A Graphical Model for Interaction Data

We now extend the VTA graphical model [39] to handle sequential data labeled with actions, where descriptors  $\mathbf{z}$  now represent skills. From a high level, the model partitions a trajectory into subsequences with the boundary variables  $\mathbf{m}$  and labels each subsequence as a skill  $\mathbf{z}$ .

We wish to learn a state-conditional policy rather than the joint distribution of the whole trajectory. To do this, we write a new generative model for the actions  $\mathbf{a}_{1:T}$  conditional on the state  $\mathbf{x}_{1:T}$ :

$$p(\mathbf{z}_{1:T}, \mathbf{s}_{1:T}, \mathbf{m}_{1:T}, \mathbf{a}_{1:T} \mid \mathbf{x}_{1:T}) = \prod_{t=1}^T p(\mathbf{a}_t \mid \mathbf{s}_t) p(\mathbf{m}_t \mid \mathbf{s}_{t-1}) p(\mathbf{s}_t \mid \mathbf{x}_t, \mathbf{z}_t) p(\mathbf{z}_t \mid \mathbf{x}_t, \mathbf{z}_{t-1}, \mathbf{m}_{t-1}).$$

This differs from the original VTA generative model in two ways: (1) we introduce a  $p(\mathbf{a}_t \mid \mathbf{s}_t)$  term that indicates how actions are sampled; and (2) the distributions over  $\mathbf{z}_t$  and  $\mathbf{s}_t$  do not depend on all previous  $\mathbf{z}_{1:t}$  and  $\mathbf{s}_{1:t}$  to encode the Markov property. We then augment the variational distribution such that the posterior over skills  $\mathbf{z}$  also depends on actions:

$$q_{\theta}(\mathbf{z}_{1:T}, \mathbf{s}_{1:T}, \mathbf{m}_{1:T} \mid \mathbf{x}_{1:T}, \mathbf{a}_{1:T}) = \prod_{t=1}^T \underbrace{q_{\theta}(\mathbf{m}_t \mid \mathbf{x}_{1:t})}_{\text{termination policy}} \underbrace{q_{\theta}(\mathbf{z}_t \mid \mathbf{z}_{t-1}, \mathbf{m}_t, \mathbf{x}_{1:T}, \mathbf{a}_{1:T})}_{\text{skill posterior}} \underbrace{q_{\theta}(\mathbf{s}_t \mid \mathbf{z}_t, \mathbf{x}_t)}_{\text{state abstraction posterior}}.$$

Overall, this yields a model with 3 learned components:

- A *state abstraction posterior*  $q_{\theta}(\mathbf{s}_t \mid \mathbf{z}_t, \mathbf{x}_t)$  and *decoder*  $p_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t)$ , which together form a distribution over actions conditioned on the skill  $\mathbf{z}_t$  and current state  $\mathbf{x}_t$ . We refer to these jointly as the *skill policy*  $\pi_{\theta}(\mathbf{a}_t \mid \mathbf{z}_t, \mathbf{x}_t) = \mathbb{E}_{\mathbf{s}_t \sim q_{\theta}(\mathbf{s}_t \mid \mathbf{z}_t, \mathbf{x}_t)} [p_{\theta}(\mathbf{a}_t \mid \mathbf{s}_t)]$ .
- A *termination policy*  $q_{\theta}(\mathbf{m}_t \mid \mathbf{x}_{1:t})$ , which determines whether the previous skill terminates.
- A *skill posterior*  $q_{\theta}(\mathbf{z}_t \mid \mathbf{z}_{t-1}, \mathbf{m}_t, \mathbf{x}_{1:T}, \mathbf{a}_{1:T})$ , which outputs the current skill  $\mathbf{z}_t$  conditioned on all states  $\mathbf{x}_{1:T}$  and actions  $\mathbf{a}_{1:T}$ . This depends on  $\mathbf{m}_t$  and  $\mathbf{z}_{t-1}$ , since the boundary variables determine whether the previous skill  $\mathbf{z}_{t-1}$  terminates: when  $\mathbf{m}_t = 0$ , the previous skill does not terminate and  $\mathbf{z}_t = \mathbf{z}_{t-1}$ .

Once this model is learned, the skill policy and termination policy represent the skills, without a need for the skill posterior: Given a skill variable  $\mathbf{z}$ , the skill policy encodes what actions the skill takes, and the termination policy determines when the skill stops. In the next section, we describe our objective for learning this model. Then, in Section 5, we describe how we use the skill policy and termination policy to learn new tasks.

### 4.2 Discovering Structure via Compression

As previously discussed, the maximum likelihood objective is underspecified for skill learning, because many skills can maximize likelihood, independent of whether they are useful for learning new tasks. In this section, we address this with a new objective that attempts to measure how useful skills will be for learning new tasks in terms of compression. Our objective is based on the intuition that effectively compressing a sequence of data requires factoring out common structure, and factoring out common structure is critical for learning useful skills.

We measure the complexity of a skill decomposition as the amount of information required to communicate the sequence of skills that encode the pre-collected experience. The latent variable model introduced in the previous section encodes each trajectory as a sequence of skills  $\mathbf{z}_{1:T}$  and boundaries  $\mathbf{m}_{1:T}$ . However, using the skill and termination policies, it is possible to recover each trajectory from only the skill variables at the boundary points: i.e., at the time steps  $\{t \in [T] \mid \mathbf{m}_t = 1\}$ . For a prior on skills  $p_{\mathbf{z}}$ , an optimal code requires  $-\log p_{\mathbf{z}}(\mathbf{z}_t)$  bits to send a skill  $\mathbf{z}_t$  [12, Chapter 5.2]. Hence,

the expected code length of communicating a trajectory  $\tau_{1:T} = \{(x_1, a_1), \dots, (x_T, a_T)\}$  is:

$$\text{INFOCOST}(\theta; p_z) = -\mathbb{E}_{\substack{\tau_{1:T}, \\ \mathbf{m}_{1:T}, \\ \mathbf{z}_{1:T}}} \left[ \sum_{t=1}^T \log p_z(z_t) \mathbf{m}_t \right],$$

where the expectation is under trajectories  $\tau_{1:T}$  from the pre-collected experience and sampling  $\mathbf{z}_{1:T}$  from the skill posterior and  $\mathbf{m}_{1:T}$  from the termination policy. The choice of prior that minimizes the average code length is one that equals the empirical distribution of skills under the pre-collected experience [12, Chapter 5.3]:

$$p_z^* \triangleq \arg \min_{p_z} \text{INFOCOST}(\theta; p_z), \quad \text{where } p_z^*(z) = \frac{1}{n_s} \mathbb{E}_{\substack{\tau_{1:T}, \\ \mathbf{m}_{1:T}, \\ \mathbf{z}_{1:T}}} \left[ \sum_{t=1}^T \delta(z_t = z) \mathbf{m}_t \right],$$

and  $n_s \triangleq \mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}} \left[ \sum_{t=1}^T \mathbf{m}_t \right]$ . In Appendix A.1, we show that this marginal  $p_z^*$  is a proper density for both continuous and discrete  $z$ . Substituting in this optimal choice of prior, we can show that the code length can be expressed as the marginal entropy  $\mathcal{H}_{p_z^*}[z]$  times the number of skills per trajectory  $n_s$  (see Appendix A.2 for proof):

$$\mathcal{L}_{\text{CL}}(\theta) \triangleq \min_{p(z_t)} \text{INFOCOST}(\theta; p(z_t)) = -\mathbb{E}_{\substack{\tau_{1:T}, \\ \mathbf{m}_{1:T}, \\ \mathbf{z}_{1:T}}} \left[ \sum_{t=1}^T \log p_z^*(z_t) \mathbf{m}_t \right] = n_s \mathcal{H}_{p_z^*}[z]. \quad (1)$$

Intuitively, this is equal to the average code length of a skill multiplied by the the average number of skills per trajectory. Note that compression is only beneficial if the model also achieves high likelihood of the data. We capture this by solving the following constrained optimization problem:

$$\min_{\theta} \mathcal{L}_{\text{CL}}(\theta) \quad \text{s.t. } \mathcal{L}_{\text{ELBO}}(\theta) \leq C, \quad (2)$$

where  $\mathcal{L}_{\text{ELBO}}(\theta)$  a negated evidence lower bound on the log-likelihood (detailed in Appendix B).

**Remarks.** We discuss a connection between this objective and variational inference in Appendix A.3. Additionally, while this work focuses on the RL setting, our objective generally applies to sequential modeling problems. We believe that it could be useful for many applications beyond option learning.

### 4.3 Connections to Minimal Description Length

Our approach closely relates to the minimum description length (MDL) principle [69]. This principle equates learning with finding regularity in data, which can be used to compress the data. Informally, the best model is the one that encodes the data with the lowest description length. Given a model  $\theta$  that encodes the data  $\mathcal{D}$  into some message, one way to formalize the description length of the data  $L(\mathcal{D})$  is with a crude two-part code [26]. This decomposes the description length of the data as the length of the message plus the length of the model:

$$L(\mathcal{D}) = \underbrace{L(\mathcal{D} | \theta)}_{\text{message length}} + \underbrace{L(\theta)}_{\text{model length}}. \quad (3)$$

In our case, the message is the latent variables  $z_t$  at the boundary points  $\{t \in [T] \mid \mathbf{m}_t = 1\}$ , which can be decoded into the data  $\mathcal{D}$  with the skill and termination policies (representing the model). Hence, our approach can be seen as an instance of minimizing the description length  $L(\mathcal{D})$ . Optimizing our objective in Equation 1 is equivalent to minimizing the message length  $L(\mathcal{D} | \theta)$ . While we do not directly attempt to minimize the model length term  $L(\theta)$ , many works [60, 82] indicate that deep learning has implicit regularization, which biases optimization toward low complexity solutions without an explicit regularizer. In general, computing the true description length or appropriate notion of complexity for neural networks is a tall order [61, 89, 36]; however, there is a rich space of methods to be explored here and we therefore leave the extension of our approach to directly minimizing the model length term for future work.

### 4.4 A Practical Implementation

**Model.** We instantiate our model by defining discrete skills  $z \in [K]$ , state abstractions  $\mathbf{s} \in \mathbb{R}^d$ , and binary boundary variables  $\mathbf{m} \in \{0, 1\}$ . We parameterize all components of our model as neural networks. See Appendix D for architecture details.



**Optimization.** We apply Gumbel-softmax [34, 54] to optimize over the discrete random variables  $z$  and  $m$ . We rewrite the compression objective  $\mathcal{L}_{\text{CL}}$  as a product of  $n_s$  and  $\mathcal{H}_{p_z^*}[z_t]$  in Equation 1 to improve the stability of optimization. This rewriting allows computing a gradient in terms of a distribution over the skill variables  $p_z^*$ , rather than samples  $z_t$ , which yields a more accurate finite sample gradient. Empirically, this leads to stabler optimization and convergence to better solutions.

We approximate the optimal skill prior  $p_z^*$  using our skill posterior as:

$$p_z^*(z) \approx \frac{1}{\widehat{\mathbb{E}}_{\substack{\tau_{1:T}, \\ m_{1:T}}} \left[ \sum_{t=1}^T m_t \right]} \widehat{\mathbb{E}}_{\substack{\tau_{1:T}, \\ m_{1:T}, \\ z_{1:T}}} \left[ \sum_{t=1}^T q_{\theta}(z_t = z \mid z_{t-1}, m_t, x_{1:T}, a_{1:T}) m_t \right],$$

where  $\widehat{\mathbb{E}}$  denotes the empirical expectation over minibatches of  $x_{1:T}, a_{1:T}$  from  $\mathcal{D}$  and sampling  $m_{1:T}$  from the termination policy. We solve the constrained optimization problem (Equation 2) with dual gradient descent on the Lagrangian. In Appendix G, we summarize the overall training procedure in Algorithm 1 and report details about the Lagrangian.

**Enforcing a minimum skill length.** Though degenerate solutions such as skills that only take a single action score poorly in our compression objective, empirically, such solutions create local optima that are difficult to escape. To avoid this, we mask the boundary variables during training to ensure that each skill  $z_t$  operates for at least  $T_{\min} = 3$  time steps. We find that when the skills are at least minimally temporally extended, optimization of the compression objective appears to be stabler and achieves better values. We remove these masks at test time, when we learn a task with our skills.

## 5 Using the Learned Skills for Hierarchical RL

We now describe how we quickly learn new tasks, given the skills learned from the pre-collected experience. Overall, we simply augment the agent’s action space with the learned skills and learn a new policy that can take either low-level actions or learned skills, similar to Kipf et al. [42]. However, our compression objective may result in unused skills, since this decreases the encoding cost of the used skills. Therefore, we first filter down to only the skills that are used to compress the pre-collected experience by selecting the skills where the marginal is over some threshold  $\alpha$ :

$$\mathcal{Z} = \{k \in [K] \mid p_z^*(k) > \alpha\}.$$

Then, on a new task with action space  $\mathcal{A}$ , we train an agent using the augmented action space  $\mathcal{A}^+ = \mathcal{A} \cup \mathcal{Z}$ . When the agent selects a skill  $z \in \mathcal{Z}$ , we follow the procedure in Algorithm 2 in the Appendix. We take actions following the skill policy  $\pi_{\theta}(a_t \mid z, x_t)$  (lines 2–3), until the termination policy  $q_{\theta}(m_t \mid x_{1:t})$  outputs a termination  $m_t = 1$  (line 5). At that point, the agent observes the next state  $x_{t+1}$  and selects the next action or skill.

## 6 A Didactic Experiment: Frame Prediction

Before studying the RL setting, we first illustrate the effect of LOVE in the simpler setting of sequential data without actions. In this setting, we use the original VTA model, which partitions a sequence  $x_{1:T}$  into contiguous subsequences with the boundary variables  $m$  and labels each subsequence with a descriptor  $z$ , as described in Section 3. Here, VTA’s objective is to maximize the likelihood of the sequence  $x_{1:T}$ , and our compression objective applies to this setting completely unmodified. We compare LOVE and only maximizing likelihood, i.e., VTA [39], by measuring how well the learned subsequences correspond to underlying patterns in the data. To isolate the effect of the compression objective, we do not enforce a minimum skill length.

**Dataset.** We consider two simple datasets *Simple Colors* and *Conditional Colors*, which consist of sequences of  $32 \times 32$  monochromatic frames, with repeating underlying patterns. *Simple Colors* consists of 3 patterns: 3 consecutive yellow frames, 3 consecutive blue frames and 3 consecutive green frames. The dataset is generated by sampling these patterns with probability 0.4, 0.4, and 0.2 respectively and concatenating the results. Each pattern is sampled independent of history. The dataset is then divided into sequences of 6 patterns, equal to  $6 \times 3 = 18$  frames.

In contrast to *Simple Colors*, *Conditional Colors* tests learning patterns that depend on history. It consists of 4 patterns: the 3 patterns from *Simple Colors* with an additional pattern of 3 consecutive



Figure 3: Learned boundaries and descriptors on *Simple Colors*. LOVE recovers the patterns, while VTA learns boundaries that break up patterns or span multiple patterns (outlined in red).

Table 1: Effect of LOVE on the *Simple Colors* and *Conditional Colors* datasets (5 seeds). All approaches achieve similar values for the likelihood objective, but LOVE better recovers the correct boundaries.

	<i>Simple Colors</i>		<i>Conditional Colors</i>	
	VTA	LOVE	VTA	LOVE
$\mathcal{L}_{\text{ELBO}}$	$2868 \pm 43$	$2838 \pm 19$	$2832 \pm 7.7$	$2827 \pm 1.9$
Precision	$0.87 \pm 0.19$	<b><math>0.99 \pm 0.01</math></b>	$0.84 \pm 0.22$	<b><math>0.99 \pm 0.01</math></b>
Recall	$0.79 \pm 0.13$	<b><math>0.85 \pm 0.03</math></b>	<b><math>0.82 \pm 0.16</math></b>	<b><math>0.83 \pm 0.06</math></b>
F1	$0.82 \pm 0.13$	<b><math>0.91 \pm 0.02</math></b>	$0.83 \pm 0.19$	<b><math>0.90 \pm 0.03</math></b>
Code Length	$7.58 \pm 1.3$	<b><math>6.34 \pm 0.75</math></b>	$9.17 \pm 1.76$	<b><math>6.83 \pm 0.51</math></b>

purple frames. The dataset is generated by uniformly sampling the patterns from *Simple Colors*. To make the current timestep pattern dependent on the previous timestep, the yellow frames are re-colored to purple, if the previous pattern was blue or yellow. As in *Simple Colors*, the dataset is divided into sequences of 6 patterns.

In both datasets, the optimal encoding strategy is to learn 3 subsequence descriptors (i.e., skills), one for each pattern in *Simple Colors*. In *Conditional Colors*, the descriptor corresponding to yellow in *Simple Colors* either outputs yellow or purple, depending on the history. This encoding strategy achieves an expected code length of 6.32 nats in *Simple Colors* and 6.59 nats in *Conditional Colors*.

**Results.** We visualize the learned descriptors and boundaries in Figure 3. LOVE successfully segments the sequences into the patterns and assigns a consistent descriptor  $z$  to each pattern. In contrast, despite the simple structure of this problem, VTA learns subsequences that span over multiple patterns or last for fewer timesteps than a single pattern.

Quantitatively, we measure the (1) precision, recall, and F1 scores of the boundary prediction, (2) the ELBO of the maximum likelihood objective and (3) the average code length  $\mathcal{L}_{\text{CL}}$  in Table 1. While both VTA and LOVE achieve approximately the same negated ELBO (lower is better), LOVE recovers the correct boundaries with much higher precision, recall, and F1 scores. This illustrates that the underspecification problem can occur even in simple sequential data. Additionally, we find that LOVE achieves an encoding cost close to the optimal value. One interesting, though rare, failure mode is that LOVE sometimes even achieves a lower encoding cost than the optimal value by over-weighting the compression term and imperfectly reconstructing the data. In Appendix C, we conduct an ablation study on the weights of  $\mathcal{L}_{\text{CL}}$  in optimization.

## 7 Experiments

We aim to answer three questions: (1) Does LOVE learn semantically meaningful skills? (2) Do skills learned with LOVE accelerate learning for new tasks? (3) Can LOVE scale to high-dim pixel observations? To answer these questions, we learn skills from demonstrations and use these skills to learn new tasks, first on a 2D multi-task domain, and later on a 3D variant.

**Multi-task domain.** We consider the multi-task  $10 \times 10$  grid world introduced by Kipf et al. [42], a representative skill learning approach (Figure 4). This domain features a challenging control problem that requires the agent to collect up to 5 objects, and we later add a perception challenge to it with the 3D variant. We use our own custom implementation, since the code from Kipf et al. [42] is not publicly available. In each task, 6 objects are randomly sampled from a set of  $N_{\text{obj}} = 10$  different objects and placed at random positions in the grid. Additionally, impassable walls are randomly added to the grid, and the agent is also placed at a random initial grid cell. Each task also includes an instruction list of  $N_{\text{pick}} = 3$  or 5 objects that the agent must pick up in order.

Within each task, the agent’s actions are to move up, down, left, right, or to pick up an item at the agent’s grid cell. The state consists of two parts: (1) a  $10 \times 10 \times (N_{\text{obj}} + 2)$  grid observation,

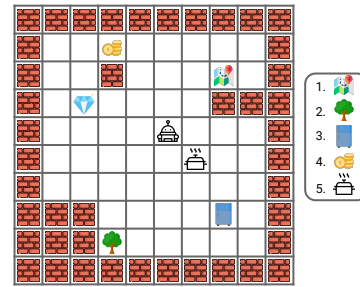


Figure 4: Multi-task grid world instance from Kipf et al. [42]. The agent must pick up a sequence of objects in a specified order.

indicating if the agent, a wall, or any of the  $N_{\text{obj}}$  different object types is present at each of the  $10 \times 10$  grid cells; (2) the instruction corresponding to the next object to pick up, encoded as a one-hot integer. Following Kipf et al. [42], we consider two reward variants: In the *sparse reward* variant, the agent only receives +1 reward after picking up *all*  $N_{\text{pick}}$  objects in the correct order. In the *dense reward* variant, the agent receives +1 reward after picking up each specified object in the correct order. Our dense reward variant is slightly harder than the variant in Kipf et al. [42], which gives the agent +1 reward for picking up objects in *any* order. The agent receives 0 reward in all other time steps. The episode ends when the agent has picked up all the objects in the correct order, or after 50 timesteps.

**Pre-collected experience.** We follow the setting in Kipf et al. [42]. We set the pre-collected experience to be 2000 demonstrations generated via breadth-first search on randomly generated tasks with only  $N_{\text{pick}} = 3$  and test if the agent can generalize to  $N_{\text{pick}} = 5$  when learning a new task. These demonstrations are not labeled with rewards and also *do not* contain the instruction list observations.

**Points of comparison.** To study our compression objective, we compare with two representatives of learning skills via the maximum likelihood objective:

- VTA [39], modified to handle interaction data as in Section 4.1.
- Discovery of deep options (DDO) [21]. We implement DDO’s maximum likelihood objective and graphical model on top of VTA’s variational inference optimization procedure, instead of expectation-gradients used in the original paper.

For fairness, we also compare with variants of these that implement the minimum skill length constraint from Section 4.4. CompILE [42] is another approach that learns skill by maximizing likelihood and was introduced in the same paper as this grid world. Because the implementation is unavailable, we do not compare with it. However, we note that CompILE requires additional supervision that LOVE, VTA, and DDO do not: *namely, it requires knowing how many skills each demonstration is composed of*. This supervision can be challenging to obtain without already knowing what the skills should be.

Since latent variable models are prone to local optima [45], it is common to learn such models with multiple restarts [57]. We therefore run each method with 3 random initializations and pick the best model according to the compression objective  $\mathcal{L}_{\text{CL}}$  for LOVE and according to the ELBO of the maximum likelihood objective  $\mathcal{L}_{\text{ELBO}}$  for the others. Notably, this does not require any additional demonstrations or experience.

We begin by analyzing the skills learned from the pre-collected experience before analyzing the utility of these skills for learning downstream tasks in the next sections.

## 7.1 Analyzing the Learned Skills

We analyze the learned skills by comparing them to a natural decomposition that partitions the demonstration in  $N_{\text{pick}}$  sequences of moving to and picking up an object. Specifically, we measure the precision and recall of each method in outputting the correct boundaries of these  $N_{\text{pick}}$  sequences.

We find that only LOVE recovers the correct boundaries with both high precision and recall (Table 2). Visually examining the

skills learned with LOVE shows that each skill moves to an object and picks it up. We visualize LOVE’s skills in Appendix H. In contrast, maximizing likelihood with either DDO or VTA learns degenerate skills that output only a single action, leading to high recall, but low precision. Adding the minimum skill length constraint to these approaches prevents such degenerate solutions, but still does not lead to the correct boundaries. Skills learned by VTA with the minimum skill length constraint exhibit no apparent structure, but skills learned by DDO with the minimum skill length constraint appear to move toward objects, though they terminate at seemingly random places. These results further illustrate the underspecification problem. All approaches successfully learn to reconstruct the

Table 2: Precision and recall of outputting the boundaries between picking up different objects. Only LOVE learns skills that move to and pick up an object.

	Precision	Recall	F1
DDO [21]	0.19	<b>1</b>	0.32
DDO + min. skill length	0.26	0.53	0.35
VTA [39]	0.19	<b>0.99</b>	0.32
VTA + min. skill length	0.27	0.53	0.36
LOVE (ours)	<b>0.90</b>	0.94	<b>0.92</b>



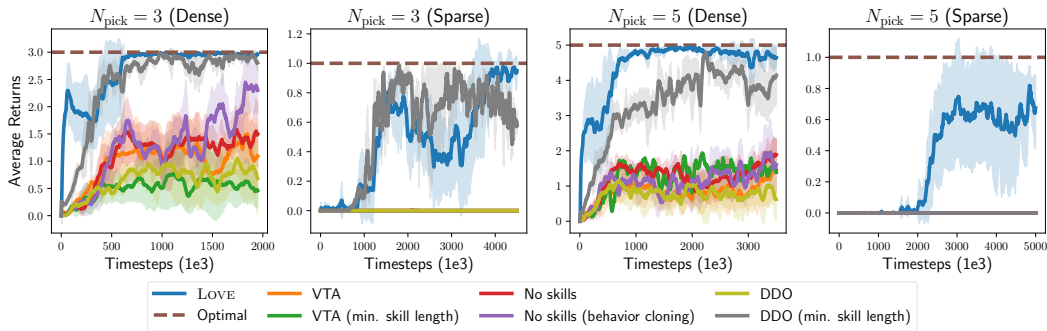


Figure 5: **Sample efficient learning.** We plot returns vs. timesteps of environment interactions for 4 settings in the grid world with 1-stddev error bars (5 seeds). Only LOVE achieves high returns across all 4 settings.

demonstrations and achieve high likelihood, but only LOVE and DDO with the minimum skill length constraint learn skills that appear to be semantically meaningful.

## 7.2 Learning New Tasks

Next, we evaluate the utility of the skills for learning new tasks. To evaluate the skills, we sample a new task in each of 4 settings: with sparse or dense rewards and with  $N_{\text{pick}} = 3$  or  $N_{\text{pick}} = 5$ . Then we learn the new tasks following the procedure in Section 5: we augment the action space with the skills and train an agent over the augmented action space. To understand the impact of skills, we also compare with a low-level policy that learns these tasks using only the original action space, and the same low-level policy that incorporates the demonstrations by pre-training with behavioral cloning. We parametrize the policy for all approaches with dueling double deep Q-networks [56, 86, 83] with  $\epsilon$ -greedy exploration. We report the returns of evaluation episodes with  $\epsilon = 0$ , which are run every 10 episodes. We report returns averaged over 5 seeds. See Appendix F.4 for full details.

**Results.** Overall, LOVE learns new tasks across all 4 settings comparably to or faster than both skill methods based on maximizing likelihood and incorporating demonstrations via behavior cloning (Figure 5). More specifically, when  $N_{\text{pick}} = 3$ , DDO with the min. skill length constraint performs comparably to LOVE, despite not segmenting the demonstrations into sequences of moving to and picking up objects. We observe that while a single DDO skill does not move to and pick up an object, the skills consistently move toward objects, which likely helps with exploration. Imitating the demonstrations with behavior cloning also accelerates learning when  $N_{\text{pick}} = 3$  with dense rewards, though not as much as LOVE or DDO, and yields insignificant benefit in all other settings. Skill learning with VTA yields no benefit over no skill learning at all.

Recall that  $N_{\text{pick}} = 3$  in all of the demonstrations. LOVE learns new tasks in the generalization setting of  $N_{\text{pick}} = 5$  much faster than the other methods. With dense rewards, DDO (min. skill length) also eventually solves the task, but requires over  $8\times$  more timesteps. The other approaches learn to pick up some objects, but never achieve optimal returns of 5. With sparse rewards, LOVE is the only approach that achieves high returns, while all other approaches achieve 0 returns. This likely occurs because this setting creates a challenging exploration problem, so exploring with low-level actions or skills learned with VTA or DDO rarely achieves reward. By contrast, exploring with skills learned with LOVE achieves rewards much more frequently, which enables LOVE to ultimately solve the task.

## 7.3 Scaling to High-dimensional Observations

Prior works in the offline multi-task setting have learned skills from low-dimensional states and then transferred them to high-dim pixel observations [68], or have learned hierarchical models from pixel observations [39, 21]. However, to the best of our knowledge, prior works in this setting have not directly learned skills from high-dim pixel observations, only from low-dimensional states. Hence, in this section, we test if LOVE can scale up to high-dim pixel observations by considering a challenging 3D image-based variant of the multi-task domain above, illustrated in Figure 6. The goal is still to pick up  $N_{\text{pick}}$  objects in the correct order, where objects are blocks of different colors. Each task includes 3 objects randomly sampled from a set of  $N_{\text{obj}} = 6$  different colors. The observations are  $400 \times 60 \times 3$  egocentric panoramic RGB arrays. The actions are to move forward / backward, turn left / right, and pick up.

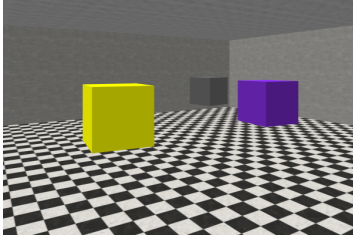


Figure 6: 3D visual multi-task domain.

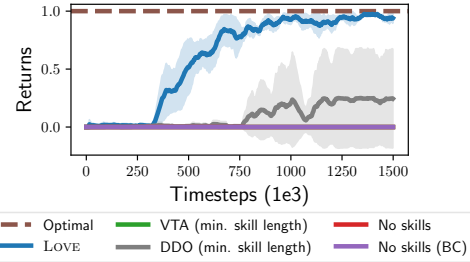


Figure 7: Returns on the 3D visual multi-task domain with 1-stddev error bars (5 seeds). Only LOVE achieves near-optimal returns.

We consider only the harder *sparse reward* and generalization setting, where the agent must pick up  $N_{\text{pick}} = 4$  objects in the correct order, after receiving 2,000 pre-collected demonstrations of the approximate shortest path of picking up  $N_{\text{pick}} = 2$  objects in the correct order. We use the same hyperparameters as in multi-task domain and only change the observation encoder and action decoder (full details in Appendix D.3). Figure 7 shows the results. Again, LOVE learns skills that each navigate to and pick up an object, which quickly achieves optimal returns, indicating its ability to scale to high-dim observations. In contrast, the other approaches perform much worse on pixel observations than the previous low-dim observations. DDO also makes progress on the task, but requires far more samples, as it again learns skills that only operate for a few timesteps, though they move toward the objects. All other approaches fail to learn at all, including the variants of VTA and DDO without the min. skill length, which are not plotted.

## 8 Conclusion

We started by highlighting the underspecification problem: maximizing likelihood does not necessarily yield skills that are useful for learning new tasks. To overcome this problem, we drew a connection between skill learning and compression and proposed a new objective for skill learning via compression and a differentiable model for optimizing our objective. Empirically, we found that the underspecification problem occurs even on simple tasks and learning skills with our objective allows learning new tasks much faster than learning skills by maximizing likelihood.

Still, important future work remains. LOVE applies when there are useful and consistent structures that can be extracted from multiple trajectories. This is often present in multi-task demonstrations, which solve related tasks in similar ways. However, an open challenge for adapting to general offline data like D4RL [22], is to ensure that the learned skills do not overfit to potentially noisy or unhelpful behaviors often present in offline data. In addition, we showed a connection between skill learning and compression by minimizing the description length of a crude two-part code. However, we only proposed a way to optimize the message length term and not the model length. Completing this connection by accounting for model length may therefore be an promising direction for future work.

**Reproducibility.** Our code is publicly available at <https://github.com/yidingjiang/love>.

**Acknowledgements.** We would like to thank Karol Hausman, Abhishek Gupta, Archit Sharma, Sergey Levine, Annie Xie, Zhe Dong, Samuel Sokota for valuable discussions during the course of this work, and Victor Akinwande, Christina Baek, Swaminathan Gurusamy for comments on the draft. We would also like to thank the anonymous reviewers for their valuable feedback. Last, but certainly not least, YJ and EZL thank the beautiful beach of Kona for initiating this project.

YJ is supported by funding from the Bosch Center for Artificial Intelligence. EZL is supported by a National Science Foundation Graduate Research Fellowship under Grant No. DGE-1656518. CF is a Fellow in the CIFAR Learning in Machines and Brains Program. This work was also supported in part by the ONR grant N00014-21-1-2685, Google, and Intel. Icons in this work were made by ThoseIcons, FreePik, VectorsMarket, from FlatIcon.

## References

- [1] Pulkit Agrawal, Ashvin Nair, Pieter Abbeel, Jitendra Malik, and Sergey Levine. Learning to poke by poking: Experiential learning of intuitive physics. *arXiv preprint arXiv:1606.07419*, 2016.
- [2] Anurag Ajay, Aviral Kumar, Pulkit Agrawal, Sergey Levine, and Ofir Nachum. Opal: Offline primitive discovery for accelerating offline reinforcement learning. *arXiv preprint arXiv:2010.13611*, 2020.
- [3] Jacob Andreas, Dan Klein, and Sergey Levine. Modular multitask reinforcement learning with policy sketches. In *International Conference on Machine Learning*, pages 166–175. PMLR, 2017.
- [4] Pierre-Luc Bacon, Jean Harb, and Doina Precup. The option-critic architecture. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 31, 2017.
- [5] Akhil Bagaria, Jason K Senthil, and George Konidaris. Skill discovery for exploration and planning using deep skill graphs. In *International Conference on Machine Learning*, pages 521–531. PMLR, 2021.
- [6] Andrew G Barto and Sridhar Mahadevan. Recent advances in hierarchical reinforcement learning. *Discrete event dynamic systems*, 13(1):41–77, 2003.
- [7] Ritwik Bera, Vinicius G Goecks, Gregory M Gremillion, John Valasek, and Nicholas R Waytowich. Podnet: A neural network for discovery of plannable options. *arXiv preprint arXiv:1911.00171*, 2019.
- [8] Christopher Berner, Greg Brockman, Brooke Chan, Vicki Cheung, Przemysław Debiak, Christy Dennison, David Farhi, Quirin Fischer, Shariq Hashme, Chris Hesse, et al. Dota 2 with large scale deep reinforcement learning. *arXiv preprint arXiv:1912.06680*, 2019.
- [9] David M Blei and Pedro J Moreno. Topic segmentation with an aspect hidden markov model. In *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 343–348, 2001.
- [10] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*, 2014.
- [11] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (elus). *arXiv preprint arXiv:1511.07289*, 2015.
- [12] Thomas M Cover. *Elements of information theory*. John Wiley & Sons, 1999.
- [13] Hanjun Dai, Bo Dai, Yan-Ming Zhang, Shuang Li, and Le Song. Recurrent hidden semi-markov model. 2016.
- [14] Sudeep Dasari, Frederik Ebert, Stephen Tian, Suraj Nair, Bernadette Bucher, Karl Schmeckpeper, Siddharth Singh, Sergey Levine, and Chelsea Finn. Robonet: Large-scale multi-robot learning. *arXiv preprint arXiv:1910.11215*, 2019.
- [15] Zhe Dong, Bryan Seybold, Kevin Murphy, and Hung Bui. Collapsed amortized variational inference for switching nonlinear dynamical systems. In *International Conference on Machine Learning*, pages 2638–2647. PMLR, 2020.
- [16] Yan Duan, John Schulman, Xi Chen, Peter L Bartlett, Ilya Sutskever, and Pieter Abbeel.  $RI^2$ : Fast reinforcement learning via slow reinforcement learning. *arXiv preprint arXiv:1611.02779*, 2016.
- [17] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations*, 2019. URL <https://openreview.net/forum?id=SJx63jRqFm>.

- [18] Chelsea Finn, Ian Goodfellow, and Sergey Levine. Unsupervised learning for physical interaction through video prediction. *Advances in neural information processing systems*, 29:64–72, 2016.
- [19] Carlos Florensa, Yan Duan, and Pieter Abbeel. Stochastic neural networks for hierarchical reinforcement learning. *arXiv preprint arXiv:1704.03012*, 2017.
- [20] Emily B Fox, Erik B Sudderth, Michael I Jordan, and Alan S Willsky. A sticky hdp-hmm with application to speaker diarization. *The Annals of Applied Statistics*, pages 1020–1056, 2011.
- [21] Roy Fox, Sanjay Krishnan, Ion Stoica, and Ken Goldberg. Multi-level discovery of deep options. *arXiv preprint arXiv:1703.08294*, 2017.
- [22] Justin Fu, Aviral Kumar, Ofir Nachum, George Tucker, and Sergey Levine. D4rl: Datasets for deep data-driven reinforcement learning. *arXiv preprint arXiv:2004.07219*, 2020.
- [23] Francisco M Garcia, Bruno C da Silva, and Philip S Thomas. A compression-inspired framework for macro discovery. *arXiv preprint arXiv:1711.09048*, 2017.
- [24] Zoubin Ghahramani and Geoffrey E Hinton. Variational learning for switching state-space models. *Neural computation*, 12(4):831–864, 2000.
- [25] Karol Gregor, Danilo Jimenez Rezende, and Daan Wierstra. Variational intrinsic control. *arXiv preprint arXiv:1611.07507*, 2016.
- [26] Peter Grunwald. A tutorial introduction to the minimum description length principle. *arXiv preprint math/0406077*, 2004.
- [27] Jean Harb, Pierre-Luc Bacon, Martin Klissarov, and Doina Precup. When waiting is not an option: Learning options with a deliberation cost. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32, 2018.
- [28] James Harrison, Apoorva Sharma, Chelsea Finn, and Marco Pavone. Continuous meta-learning without tasks. *arXiv preprint arXiv:1912.08866*, 2019.
- [29] 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*, 2018. URL <https://openreview.net/forum?id=rk07ZXZRb>.
- [30] Nicolas Heess, Dhruva TB, Srinivasan Sriram, Jay Lemmon, Josh Merel, Greg Wayne, Yuval Tassa, Tom Erez, Ziyu Wang, SM Eslami, et al. Emergence of locomotion behaviours in rich environments. *arXiv preprint arXiv:1707.02286*, 2017.
- [31] Geoffrey E Hinton and Richard S Zemel. Autoencoders, minimum description length, and helmholtz free energy. *Advances in neural information processing systems*, 6:3–10, 1994.
- [32] Antti Honkela and Harri Valpola. Variational learning and bits-back coding: an information-theoretic view to bayesian learning. *IEEE transactions on Neural Networks*, 15(4):800–810, 2004.
- [33] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *International conference on machine learning*, pages 448–456. PMLR, 2015.
- [34] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- [35] Yiding Jiang, Shixiang Gu, Kevin Murphy, and Chelsea Finn. Language as an abstraction for hierarchical deep reinforcement learning. *arXiv preprint arXiv:1906.07343*, 2019.
- [36] Yiding Jiang\*, Behnam Neyshabur\*, Hossein Mobahi, Dilip Krishnan, and Samy Bengio. Fantastic generalization measures and where to find them. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=SJgIPJBFvH>.

- [37] Matthew J Johnson, David K Duvenaud, Alex Wiltschko, Ryan P Adams, and Sandeep R Datta. Composing graphical models with neural networks for structured representations and fast inference. *Advances in neural information processing systems*, 29:2946–2954, 2016.
- [38] Rahul Kidambi, Aravind Rajeswaran, Praneeth Netrapalli, and Thorsten Joachims. Morel: Model-based offline reinforcement learning. *arXiv preprint arXiv:2005.05951*, 2020.
- [39] Taesup Kim, Sungjin Ahn, and Yoshua Bengio. Variational temporal abstraction. *Advances in Neural Information Processing Systems*, 32:11570–11579, 2019.
- [40] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [41] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*, 2013.
- [42] Thomas Kipf, Yujia Li, Hanjun Dai, Vinicius Zambaldi, Alvaro Sanchez-Gonzalez, Edward Grefenstette, Pushmeet Kohli, and Peter Battaglia. Compile: Compositional imitation learning and execution. In *International Conference on Machine Learning*, pages 3418–3428. PMLR, 2019.
- [43] Sanjay Krishnan, Roy Fox, Ion Stoica, and Ken Goldberg. Ddco: Discovery of deep continuous options for robot learning from demonstrations. In *Conference on robot learning*, pages 418–437. PMLR, 2017.
- [44] Tejas D Kulkarni, Karthik Narasimhan, Ardavan Saeedi, and Josh Tenenbaum. Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation. *Advances in neural information processing systems*, 29:3675–3683, 2016.
- [45] M Pawan Kumar, Benjamin Packer, and Daphne Koller. Self-paced learning for latent variable models. In *NIPS*, volume 1, page 2, 2010.
- [46] Sang-Hyun Lee and Seung-Woo Seo. Learning compound tasks without task-specific knowledge via imitation and self-supervised learning. In *International Conference on Machine Learning*, pages 5747–5756. PMLR, 2020.
- [47] Youngwoon Lee, Shao-Hua Sun, Sriram Somasundaram, Edward S Hu, and Joseph J Lim. Composing complex skills by learning transition policies. In *International Conference on Learning Representations*, 2018.
- [48] Scott Linderman, Matthew Johnson, Andrew Miller, Ryan Adams, David Blei, and Liam Paninski. Bayesian learning and inference in recurrent switching linear dynamical systems. In *Artificial Intelligence and Statistics*, pages 914–922. PMLR, 2017.
- [49] Scott W Linderman, Andrew C Miller, Ryan P Adams, David M Blei, Liam Paninski, and Matthew J Johnson. Recurrent switching linear dynamical systems. *arXiv preprint arXiv:1610.08466*, 2016.
- [50] Evan Z Liu, Aditi Raghunathan, Percy Liang, and Chelsea Finn. Decoupling exploration and exploitation for meta-reinforcement learning without sacrifices. In *International conference on machine learning*, pages 6925–6935. PMLR, 2021.
- [51] Evan Zheran Liu, Ramtin Keramati, Sudarshan Seshadri, Kelvin Guu, Panupong Pasupat, Emma Brunskill, and Percy Liang. Learning abstract models for strategic exploration and fast reward transfer. *arXiv preprint arXiv:2007.05896*, 2020.
- [52] Siqi Liu, Guy Lever, Zhe Wang, Josh Merel, SM Eslami, Daniel Hennes, Wojciech M Czarnecki, Yuval Tassa, Shayegan Omidshafiei, Abbas Abdolmaleki, et al. From motor control to team play in simulated humanoid football. *arXiv preprint arXiv:2105.12196*, 2021.
- [53] Yuchen Lu, Yikang Shen, Siyuan Zhou, Aaron Courville, Joshua B Tenenbaum, and Chuang Gan. Learning task decomposition with ordered memory policy network. *arXiv preprint arXiv:2103.10972*, 2021.



- [54] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- [55] Amy McGovern and Richard S Sutton. Macro-actions in reinforcement learning: An empirical analysis. *Computer Science Department Faculty Publication Series*, page 15, 1998.
- [56] Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A Rusu, Joel Veness, Marc G Bellemare, Alex Graves, Martin Riedmiller, Andreas K Fidjeland, Georg Ostrovski, et al. Human-level control through deep reinforcement learning. *nature*, 518(7540):529–533, 2015.
- [57] Kevin P Murphy. *Machine learning: a probabilistic perspective*. MIT press, 2012.
- [58] Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement learning. *arXiv preprint arXiv:1805.08296*, 2018.
- [59] Taewook Nam, Shao-Hua Sun, Karl Pertsch, Sung Ju Hwang, and Joseph J Lim. Skill-based meta-reinforcement learning. In *Fifth Workshop on Meta-Learning at the Conference on Neural Information Processing Systems*, 2021. URL <https://openreview.net/forum?id=jsV1-AQVEFY>.
- [60] Behnam Neyshabur, Ryota Tomioka, and Nathan Srebro. In search of the real inductive bias: On the role of implicit regularization in deep learning. *arXiv preprint arXiv:1412.6614*, 2014.
- [61] Behnam Neyshabur, Srinadh Bhojanapalli, David McAllester, and Nathan Srebro. Exploring generalization in deep learning. *arXiv preprint arXiv:1706.08947*, 2017.
- [62] Scott Niekum, Sachin Chitta, Andrew G Barto, Bhaskara Marthi, and Sarah Osentoski. Incremental semantically grounded learning from demonstration. In *Robotics: Science and Systems*, volume 9, pages 10–15607. Berlin, Germany, 2013.
- [63] Junhyuk Oh, Satinder Singh, Honglak Lee, and Pushmeet Kohli. Zero-shot task generalization with multi-task deep reinforcement learning. In *International Conference on Machine Learning*, pages 2661–2670. PMLR, 2017.
- [64] Aaron van den Oord, Sander Dieleman, Heiga Zen, Karen Simonyan, Oriol Vinyals, Alex Graves, Nal Kalchbrenner, Andrew Senior, and Koray Kavukcuoglu. Wavenet: A generative model for raw audio. *arXiv preprint arXiv:1609.03499*, 2016.
- [65] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning. *arXiv preprint arXiv:1711.00937*, 2017.
- [66] Karl Pertsch, Youngwoon Lee, and Joseph J Lim. Accelerating reinforcement learning with learned skill priors. *arXiv preprint arXiv:2010.11944*, 2020.
- [67] Kate Rakelly, Aurick Zhou, Chelsea Finn, Sergey Levine, and Deirdre Quillen. Efficient off-policy meta-reinforcement learning via probabilistic context variables. In *International conference on machine learning*, pages 5331–5340. PMLR, 2019.
- [68] Dushyant Rao, Fereshteh Sadeghi, Leonard Hasenclever, Markus Wulfmeier, Martina Zambelli, Giulia Vezzani, Dhruva Tirumala, Yusuf Aytar, Josh Merel, Nicolas Heess, et al. Learning transferable motor skills with hierarchical latent mixture policies. *arXiv preprint arXiv:2112.05062*, 2021.
- [69] Jorma Rissanen. Modeling by shortest data description. *Automatica*, 14(5):465–471, 1978.
- [70] Melrose Roderick, Christopher Grimm, and Stefanie Tellex. Deep abstract q-networks. *arXiv preprint arXiv:1710.00459*, 2017.
- [71] Tanmay Shankar and Abhinav Gupta. Learning robot skills with temporal variational inference. In *International Conference on Machine Learning*, pages 8624–8633. PMLR, 2020.
- [72] Tanmay Shankar, Shubham Tulsiani, Lerrel Pinto, and Abhinav Gupta. Discovering motor programs by recomposing demonstrations. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=rkgHYONYwr>.

- [73] Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware unsupervised discovery of skills. In *International Conference on Learning Representations*, 2020. URL <https://openreview.net/forum?id=HJgLZR4KvH>.
- [74] Arjun Sharma, Mohit Sharma, Nicholas Rhinehart, and Kris M Kitani. Directed-info gail: Learning hierarchical policies from unsegmented demonstrations using directed information. *arXiv preprint arXiv:1810.01266*, 2018.
- [75] Kyriacos Shiarlis, Markus Wulfmeier, Sasha Salter, Shimon Whiteson, and Ingmar Posner. Taco: Learning task decomposition via temporal alignment for control. In *International Conference on Machine Learning*, pages 4654–4663. PMLR, 2018.
- [76] Yuri M Shtarkov. Univresal sequential coding of single messages. *Prob. Pered. Inf.*, 23:175–186, 1987.
- [77] Avi Singh, Huihan Liu, Gaoyue Zhou, Albert Yu, Nicholas Rhinehart, and Sergey Levine. Parrot: Data-driven behavioral priors for reinforcement learning. In *International Conference on Learning Representations*, 2021. URL <https://openreview.net/forum?id=Ysuv-WOFeKR>.
- [78] Sumedh A Sontakke, Sumegh Roychowdhury, Mausoom Sarkar, Nikaash Puri, Balaji Krishnamurthy, and Laurent Itti. Video2skill: Adapting events in demonstration videos to skills in an environment using cyclic mdp homomorphisms. *arXiv preprint arXiv:2109.03813*, 2021.
- [79] Richard S Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for temporal abstraction in reinforcement learning. *Artificial intelligence*, 112(1-2): 181–211, 1999.
- [80] Daniel Tanneberg, Kai Ploeger, Elmar Rueckert, and Jan Peters. Skid raw: Skill discovery from raw trajectories. *IEEE Robotics and Automation Letters*, 6(3):4696–4703, 2021.
- [81] Sebastian Thrun and Anton Schwartz. Finding structure in reinforcement learning. *Advances in neural information processing systems*, 7, 1994.
- [82] Guillermo Valle-Perez, Chico Q Camargo, and Ard A Louis. Deep learning generalizes because the parameter-function map is biased towards simple functions. *arXiv preprint arXiv:1805.08522*, 2018.
- [83] Hado Van Hasselt, Arthur Guez, and David Silver. Deep reinforcement learning with double q-learning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 30, 2016.
- [84] Alexander Sasha Vezhnevets, Simon Osindero, Tom Schaul, Nicolas Heess, Max Jaderberg, David Silver, and Koray Kavukcuoglu. Feudal networks for hierarchical reinforcement learning. In *International Conference on Machine Learning*, pages 3540–3549. PMLR, 2017.
- [85] Jane X Wang, Zeb Kurth-Nelson, Dhruva Tirumala, Hubert Soyer, Joel Z Leibo, Remi Munos, Charles Blundell, Dhharshan Kumaran, and Matt Botvinick. Learning to reinforcement learn. *arXiv preprint arXiv:1611.05763*, 2016.
- [86] Ziyu Wang, Tom Schaul, Matteo Hessel, Hado Hasselt, Marc Lanctot, and Nando Freitas. Dueling network architectures for deep reinforcement learning. In *International conference on machine learning*, pages 1995–2003. PMLR, 2016.
- [87] Mengjiao Yang and Ofir Nachum. Representation matters: Offline pretraining for sequential decision making. *arXiv preprint arXiv:2102.05815*, 2021.
- [88] Mengjiao Yang, Sergey Levine, and Ofir Nachum. Trail: Near-optimal imitation learning with suboptimal data. *arXiv preprint arXiv:2110.14770*, 2021.
- [89] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3): 107–115, 2021.
- [90] Jesse Zhang, Karl Pertsch, Jiefan Yang, and Joseph J Lim. Minimum description length skills for accelerated reinforcement learning. In *Self-Supervision for Reinforcement Learning Workshop - ICLR 2021*, 2021. URL <https://openreview.net/forum?id=r4XxtrIo1m9>.

- [91] Zelin Zhao, Chuang Gan, Jiajun Wu, Xiaoxiao Guo, and Joshua B Tenenbaum. Augmenting policy learning with routines discovered from a demonstration. *arXiv preprint arXiv:2012.12469*, 2020.
- [92] Wenxuan Zhou, Sujay Bajracharya, and David Held. Plas: Latent action space for offline reinforcement learning. *arXiv preprint arXiv:2011.07213*, 2020.
- [93] Yifeng Zhu, Peter Stone, and Yuke Zhu. Bottom-up skill discovery from unsegmented demonstrations for long-horizon robot manipulation. *arXiv preprint arXiv:2109.13841*, 2021.
- [94] Luisa Zintgraf, Kyriacos Shiarlis, Maximilian Igl, Sebastian Schulze, Yarin Gal, Katja Hofmann, and Shimon Whiteson. Varibad: A very good method for bayes-adaptive deep rl via meta-learning. *arXiv preprint arXiv:1910.08348*, 2019.

## A Derivations

### A.1 $p_z^*$ is a valid density

**Proposition A.1.**  $p_z^*(\xi)$  is a proper probability density function.

*Proof.* Recall that

$$p_z^*(\xi) = \frac{\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}, \mathbf{z}_{1:T}} \left[ \sum_{t=1}^T \delta(z_t = \xi) \mathbf{m}_t \right]}{\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}} \left[ \sum_{t=1}^T \mathbf{m}_t \right]}.$$

A function is a probability density function if it is non-negative, continuous, and integrates to 1. By construction,  $p_z^*(\xi)$  is non-negative since  $\mathbf{m}_t \geq 0$  and  $\delta(z_t = \xi) \geq 0$ . Further,  $\delta(z_t = \xi)$  is continuous in  $\xi$  so the numerator  $\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}, \mathbf{z}_{1:T}} \left[ \sum_{t=1}^T \delta(z_t = \xi) \mathbf{m}_t \right]$ , an expectation containing  $\delta(z_t = \xi)$ , is also continuous in  $\xi$ . Finally,

$$\begin{aligned} \int_{\xi} p_z^*(\xi) d\xi &= \frac{1}{\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}} \left[ \sum_{t=1}^T \mathbf{m}_t \right]} \int_{\xi} \mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}, \mathbf{z}_{1:T}} \left[ \sum_{t=1}^T \delta(z_t = \xi) \mathbf{m}_t \right] d\xi \\ &= \frac{1}{\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}} \left[ \sum_{t=1}^T \mathbf{m}_t \right]} \mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}, \mathbf{z}_{1:T}} \left[ \sum_{t=1}^T \left( \int_{\xi} \delta(z_t = \xi) d\xi \right) \mathbf{m}_t \right] \quad (\text{Fubini's theorem}) \\ &= \frac{1}{\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}} \left[ \sum_{t=1}^T \mathbf{m}_t \right]} \mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}, \mathbf{z}_{1:T}} \left[ \sum_{t=1}^T 1 \cdot \mathbf{m}_t \right] \\ &= 1. \end{aligned}$$

Therefore,  $p_z^*(\xi)$  satisfies all 3 properties of a probability density function.  $\square$

**Remarks.** This result generalizes to the case for discrete  $z$  by replacing the delta with an indicator function and integration with summation.

### A.2 Two expressions of $\mathcal{L}_{CL}$

**Proposition A.2.** The expected coding length of a trajectory is equal to the expected number skills multiplied by the marginal entropy:

$$\mathcal{L}_{CL}(\theta) = -\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}, \mathbf{z}_{1:T}} \left[ \sum_{t=1}^T \log p_z^*(z_t) \mathbf{m}_t \right] = n_s \mathcal{H}_{p_z^*}[\mathbf{z}].$$

*Proof.*

$$\begin{aligned} & -\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}, \mathbf{z}_{1:T}} \left[ \sum_{t=1}^T \log p_z^*(z_t) \mathbf{m}_t \right] \\ &= -\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}, \mathbf{z}_{1:T}} \left[ \sum_{t=1}^T \left( \int_{\xi} \delta(z_t = \xi) \log p_z^*(\xi) d\xi \right) \mathbf{m}_t \right] \\ &= -\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}, \mathbf{z}_{1:T}} \left[ \int_{\xi} \sum_{t=1}^T \mathbf{m}_t \delta(z_t = \xi) \log p_z^*(\xi) d\xi \right] \quad (\text{Fubini's theorem}) \\ &= -\int_{\xi} \mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}, \mathbf{z}_{1:T}} \left[ \sum_{t=1}^T \mathbf{m}_t \delta(z_t = \xi) \right] \log p_z^*(\xi) d\xi \quad (\text{Linearity of expectation}) \\ &= -\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}} \left[ \sum_{t=1}^T \mathbf{m}_t \right] \int_{\xi} \frac{\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}, \mathbf{z}_{1:T}} \left[ \sum_{t=1}^T \mathbf{m}_t \delta(z_t = \xi) \right]}{\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}} \left[ \sum_{t=1}^T \mathbf{m}_t \right]} \log p_z^*(\xi) d\xi \end{aligned}$$

$$\begin{aligned}
&= \underbrace{\mathbb{E}_{\tau_{1:T}, \mathbf{m}_{1:T}, \mathbf{z}_{1:T}} \left[ \sum_{t=1}^T \mathbf{m}_t \right]}_{n_s} \underbrace{\int_{\xi} -p_{\mathbf{z}}^*(\xi) \log p_{\mathbf{z}}^*(\xi) d\xi}_{\mathcal{H}_{p_{\mathbf{z}}^*}[\mathbf{z}]} \\
&= n_s \mathcal{H}_{p_{\mathbf{z}}^*}[\mathbf{z}]
\end{aligned}$$

□

**Remarks.** This result generalizes to the case for discrete  $\mathbf{z}$  by replacing the delta with an indicator function and integration with summation. Further, this rewriting not only offers an insightful interpretation of the objective but also has *different* finite-sample gradient compared to the original form. We discuss its implication further in Section 4.4.

### A.3 Connection between LOVE and Variational Inference

Is the compression objective a prior in the framework of variational inference (VI)? The answer to this question is surprisingly nuanced. A bridge that connects VI and source coding is the bits-back coding argument [31, 32]. Under this scheme, the additional cost of communicating a message is equal to the KL divergence between the approximate posterior and the chosen prior distribution. This is because the posterior often does not match the prior and the sender must use more bits to send the message compared to using the prior exactly. In modern deep latent variable models such as VAE [41], the prior is in general chosen to be an isotropic Gaussian distribution or a sequence of isotropic Gaussian distributions that decomposes over the time steps. The bit-back coding argument shows that VI as a coding scheme is asymptotically optimal (in the limit of sending large numbers of messages), but it does not immediately guarantee good *representation learning*. As a thought experiment, imagine we choose the dimension of the latent variable in a VAE to be equal to the dimension of the data itself. It is clear that the model is not incentivized to perform good representation learning, but the bit-back coding scheme still guarantees asymptotic optimality. Instead, good representation can only emerge *if* the prior is chosen properly.

Contrary to this prevailing design choice, our latent variables  $\mathbf{z}_{1:T}$  and  $\mathbf{m}_{1:T}$  are not independent from each other and the timesteps are not independent from each other. In fact, the interaction between the latent variables is crucial for compressing the trajectory optimally in our framework. Therefore, our objective takes into account the *global* structure of the latent variables. Note that  $\mathcal{L}_{\text{CL}}$  also has the natural interpretation of measuring the average number of bits that is required to send a trajectory under the coding scheme of this work. As both quantities measure the cost of the communication, the philosophical parallel between KL divergence and  $\mathcal{L}_{\text{CL}}$  is likely not a coincidence. In practice, we observe that  $\mathcal{L}_{\text{CL}}$  also has comparable regularization effect on the latent code, and in some cases it is possible to learn a good model *without* any KL divergence if  $\mathcal{L}_{\text{CL}}$  is present. From this perspective, our compression objective is indeed a “prior” insofar as it specifies the desired properties of the posterior. Another interpretation of LOVE is that it is optimizing the prior *directly*.

Nonetheless, it is an open question whether  $\mathcal{L}_{\text{CL}}$  has a precise counterpart in Bayesian inference, that is, it is not immediately obvious if it can be written as the KL divergence between the posterior  $q(\mathbf{z}_{1:T}, \mathbf{m}_{1:T} \mid \mathbf{x}_{1:T})$  and some prior distribution  $p(\mathbf{z}_{1:T}, \mathbf{m}_{1:T})$ . One complication, for example, is the fact that the log density of  $\mathbf{m}_{1:T}$  does not participate in the computation of  $\mathcal{L}_{\text{CL}}$ . It may be possible to construct some form of hierarchical priors that enables a fully Bayesian interpretation of  $\mathcal{L}_{\text{CL}}$ , but it is also well-known that the MDL framework can accommodate codes that are not Bayesian (e.g., the Shtarkov normalized maximum likelihood code [76]). We do not offer a definitive answer to this question in this work, but the connection between Bayesian inference and  $\mathcal{L}_{\text{CL}}$  is certainly an interesting direction for future works.

## B Evidence Lower Bound

### B.1 VTA

We reproduce the data generating process of Kim et al. [39] here:

$$p(\mathbf{x}_{1:T}, \mathbf{z}_{1:T}, \mathbf{s}_{1:T}, \mathbf{m}_{1:T}) = \prod_{t=1}^T p(\mathbf{x}_t \mid \mathbf{s}_t) p(m_t \mid \mathbf{s}_t) p(\mathbf{s}_t \mid \mathbf{s}_{1:t}, \mathbf{z}_t, m_{t-1}) p(\mathbf{z}_t \mid \mathbf{z}_{1:t}, m_{t-1}).$$



VTA [39] shows that the ELBO can be written as:

$$\begin{aligned}\mathcal{L}_{\text{ELBO}}(\theta, \phi) \approx & \sum_{t=1}^T \log p_{\phi}(\mathbf{x}_t | \mathbf{s}_t) + D_{\text{KL}}(q_{\theta}(m_t | \mathbf{x}_{1:T}) \| p_{\phi}(m_t | \mathbf{s}_t)) \\ & + D_{\text{KL}}(q_{\theta}(\mathbf{z}_t | \mathbf{s}_t, \mathbf{m}_{1:T}, \mathbf{x}_{1:T}) \| p_{\phi}(\mathbf{z}_t | \mathbf{s}_{t-1}, m_{t-1})) \\ & + D_{\text{KL}}(q_{\theta}(\mathbf{s}_t | \mathbf{s}_{t-1}, \mathbf{m}_{1:T}, \mathbf{x}_{1:T}) \| p_{\phi}(\mathbf{s}_t | \mathbf{s}_{t-1}, \mathbf{z}_t, m_{t-1})).\end{aligned}$$

Like most ELBO, this ELBO can be deconstructed into a reconstruction term:

$$\mathcal{L}_{\text{rec}}(\theta, \phi) = \sum_{t=1}^T \log p_{\phi}(\mathbf{x}_t | \mathbf{s}_t),$$

and the remaining KL term  $\mathcal{L}_{\text{KL}}(\theta, \phi)$ .

## B.2 LOVE

We reproduce the data generating process of LOVE here. Instead of using  $\phi$  for the parameters of priors and decoder, we will merge all parameters into  $\theta$ :

$$p(\mathbf{a}_{1:T}, \mathbf{z}_{1:T}, \mathbf{s}_{1:T}, \mathbf{m}_{1:T} | \mathbf{x}_{1:T}) = \prod_{t=1}^T p(\mathbf{a}_t | \mathbf{s}_t) p(m_t | \mathbf{s}_t) p(\mathbf{s}_t | \mathbf{x}_t, \mathbf{z}_t, m_{t-1}) p(\mathbf{z}_t | \mathbf{x}_t, \mathbf{z}_{1:t}, m_{t-1}).$$

For LOVE, the ELBO is re-written as:

$$\begin{aligned}\mathcal{L}_{\text{ELBO}}(\theta) \approx & \sum_{t=1}^T \log p_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \\ & + D_{\text{KL}}(q_{\theta}(m_t | \mathbf{x}_{1:t}) \| p_{\theta}(m_t | \mathbf{s}_t)) + D_{\text{KL}}(q_{\theta}(\mathbf{z}_t | m_t, \mathbf{x}_{1:T}, \mathbf{a}_{1:T}) \| p(\mathbf{z}_t)) \\ & + D_{\text{KL}}(q_{\theta}(\mathbf{s}_t | \mathbf{z}_t, \mathbf{x}_t) \| p_{\theta}(\mathbf{s}_t | \mathbf{s}_{t-1}, m_{t-1})).\end{aligned}\tag{4}$$

Since  $z_t$  is a categorical distribution, we chose the uniform prior  $p(z_t)$  similar to Oord et al. [65]. Similar to VTA, the loss is also decomposed into a reconstruction term and a KL term.

## C Sensitivity to Compression Objective Weighting

Here we study how sensitive LOVE is to the choice of  $\lambda$ , the coefficient of  $\mathcal{L}_{\text{CL}}$ . We tested 5 values of  $\lambda \in [0.01, 1]$  on *Simple* and *Cond. Colors* and applied dual gradient descent to automatically tune  $\lambda$ . Stopping point is selected based on the shortest code length achieved.

Table 3: The effect of different values of  $\lambda$  and dual GD on F1 scores of *Simple Colors* and *Conditional Colors* over 5 seeds.

	$\lambda=0.01$	$\lambda=0.03$	$\lambda=0.1$	$\lambda=0.3$	$\lambda=1.0$	$\lambda$ (dual GD)	VTA
Simple (F1)	.67 $\pm$ .27	.80 $\pm$ .21	.91 $\pm$ .02	.95 $\pm$ .06	.74 $\pm$ .08	.99 $\pm$ .00	.82 $\pm$ .13
Cond. (F1)	.84 $\pm$ .06	.92 $\pm$ .02	.90 $\pm$ .03	.82 $\pm$ .08	.87 $\pm$ .08	.99 $\pm$ .00	.83 $\pm$ .19

We see that while  $\lambda$  can change the quality of solutions found, the performance is stable within a  $\pm 3\times$  regions around  $\lambda = 0.1$ . Further, using dual GD *consistently* outperforms handpicked  $\lambda$  with an ELBO threshold  $C$  slightly (e.g., 5%) higher than the  $\mathcal{L}_{\text{ELBO}}$  achieved without compression, which can be obtained with minimal prior knowledge about the task. Hence  $\lambda$  likely can be tuned without strong prior knowledge. This contrasts with the other methods such as Kipf et al. [42] which require knowing the number of skills in each demonstration to perform well while achieving the similar or better performance.

## D Architecture

### D.1 Frame Prediction

Our base architecture is similar to Kim et al. [39]<sup>2</sup> and we highlight the difference [blue](#). The parameters of the posteriors are implemented with neural networks. First, the observation  $x_{1:T}$  are encoded into a lower dimension embedding  $\hat{x}_{1:T}$  with a `observation_encoder`. The embedding is then passed through a `boundary_posterior_decoder` that outputs the parameters of  $m_{1:T}$ . The embedding is further passed through a GRU [10],  $f_{s-rnn-fwd}$  which runs on  $\hat{x}_{1:T}$  and give the cell state  $h_{1:T}$ . For  $z_{1:T}$ , we have two GRU's with cell,  $f_{z-rnn-bwd}$  and  $f_{z-rnn-fwd}$ , which run on  $\hat{x}_{1:T}$  in different temporal direction (i.e.,  $f_{z-rnn-bwd}$  sees the future) to generate cell states  $c_{1:T}^{fwd}$  and  $c_{1:T}^{bwd}$ . For each time step  $t$ , a candidate  $z'_t$  is sampled from a  $n_z$ -dimensional categorical distribution<sup>3</sup> whose log probability is a linear projection of  $c_t^{fwd} \| c_t^{bwd}$ . The sampled  $z'_t$  are embedded linearly into a vector of size 128. If  $m_1 = 1$ , the  $z_t = z'_t$ ; otherwise,  $z_t = z_{t-1}$ <sup>4</sup>. Then we concatenate  $h_{t-1} \| z_t$  and linearly project the vector into the mean and standard deviation of a isotropic Gaussian in  $\mathbb{R}^8$ . A  $\hat{s}_t$  is sampled from this Gaussian.

Th model also keeps two “belief” vectors that keep track of the history of the computation. These vectors are modulated by two RNN,  $f_{h-rnn}$  and  $f_{c-rnn}$ . The abstraction belief  $c_t = m_t \cdot f_{c-rnn}(z_t, c_{t-1}) + (1 - m_t) \cdot c_{t-1}$ . The observation belief  $h_t = m_t c_{t-1} + (1 - m_t) \cdot f_{h-rnn}(\hat{s}_t \| c_t \| z_t, h_{t-1})$ . Finally the state abstraction is computed as the linear projection of  $s_t = \text{proj}(h_t \| \hat{s}_t)$ . The projection layer is  $256 \times 128$ . Finally,  $s_t$  is passed through a `observation_decoder` and decode into the reconstruction of  $x_t$ . Refer to Kim et al. [39] for more details and an illustration of the computation graph.

For each convolutional layer or transposed convolution layer, the tuple's values correspond to input channel, output channel, kernel size, stride, padding. With that notation in mind, the specific hyperparameters for each components are:

- `observation_encoder`: 4-layer convolutional neural network with  $\{(3, 128, 4, 2, 1), (128, 128, 4, 2, 1), (128, 128, 4, 2, 1), (128, 128, 4, 1, 0)\}$  with ELU activation [11] and batch normalization [33]. The output is flattened linearly projected to vector of size 128.
- `boundary_posterior_decoder` is a 5-layer 1-D causal temporal convolution [64] with  $\{(128, 128, 5, 1, 2) \times 5, (128, 2, 5, 1, 2)\}$  with ELU activation and batch normalization. The output is the logits for the  $m_{1:T}$ .
- $f_{z-rnn-bwd}, f_{z-rnn-fwd}, f_{s-rnn-fwd}, f_{c-rnn}, f_{h-rnn}$ : GRU with cell size 128
- `observation_decoder`: 4-layer transposed convolutional layers  $\{(128, 128, 4, 1, 0), (128, 128, 4, 2, 1), (128, 128, 4, 2, 1), (128, 3, 4, 2, 1)\}$  with ELU activation and batch normalization on all but the last layer. The last layer has no normalization and has  $\tanh$  activation.

**Removing Boundary Regularization.** An important difference from Kim et al. [39] is that we do not enforce maximum number of subsequence  $N_{\max}$  and the maximum length of subsequence  $l_{\max}$  through the boundary prior (This effectively amounts to setting both to a value larger than the sequence length). These two hyperparameters' effect are similar to that of picking the number of segments in ComplE [42] and provide strong training signal. This assumption is in general problematic since we cannot expect to know a priori what the good values are for these hyperparameters. Kipf et al. [42] demonstrates that the performance can suffer if this kind of supervision is wrong. On the other hand, we do not assume anything about the the duration of the subsequence or how many subsequences there are in each demonstration.

We do, however, enforce a minimum length on the skill. While it is largely an optimization choice, we believe this prior is significantly weaker since useful options are almost always temporally extended. Indeed, with only the minimum length constraint, VTA is unable to learn good skills conducive for downstream tasks.

<sup>2</sup><https://github.com/taesupkim/vta>

<sup>3</sup>Kim et al. [39] uses a isotropic Gaussian.

<sup>4</sup>This is the COPY action in Kim et al. [39].

## D.2 Multi-task Grid World Environment

For trajectories, we make some larger changes to the architecture to encode properties we want in a model for learning options, e.g., Markov property. Though similar in spirit, a large portion of the architecture is different from the base VTA model, so the difference will not be highlighted.

For the posterior, we have an `observation` encoder and an `action` encoder that embeds both the state observation  $\mathbf{x}_{1:T}$  and actions  $\mathbf{a}_{1:T}$  to get embedding  $\mathbf{a}_{1:T}$  and  $\mathbf{x}_{1:T}$  each of size 128. The concatenation  $\mathbf{a}_t \parallel \mathbf{x}_t$  linear projected down to a vector of size 128. The embedding is then passed through a `boundary posterior decoder` that outputs the parameters of  $\mathbf{m}_{1:T}$ . The embedding is further passed through a GRU [10],  $f_{s-rnn-fwd}$  which runs on  $\hat{\mathbf{x}}_{1:T}$  and give the cell state  $\mathbf{h}_{1:T}$ . For  $\mathbf{z}_{1:T}$ , we have two GRU's with cell,  $f_{z-rnn-bwd}$  and  $f_{z-rnn-fwd}$ , which run on  $\hat{\mathbf{x}}_{1:T}$  in different temporal direction (i.e.,  $f_{z-rnn-bwd}$  sees the future) to generate cell states  $\mathbf{c}_{1:T}^{fwd}$  and  $\mathbf{c}_{1:T}^{bwd}$ .

For sampling  $\mathbf{z}$  is done differently for interaction data for better gradient and representation learning. Instead of doing straight-through estimator for the categorical random variable, we opt to use vector-quantization [65] with straight-through estimators. First,  $\mathbf{c}_t^{fwd} \parallel \mathbf{c}_t^{bwd}$  is linearly projected into dimension 128. Then a ReLU is applied and another 128 by 128 linear layer is applied. The VQ codebook is of size  $n_z \times 128$ . Instead of taking argmax of the codebook, we approximate the distribution to be proportional to the softmax over the distance over a temperature  $t_{VQ}$ . We can sample from this distribution and apply the straight-through gradient on the sampled code candidate  $\mathbf{z}'_{1:T}$ . If  $\mathbf{m}_1 = 1$ , the  $\mathbf{z}_t = \mathbf{z}'_t$ ; otherwise,  $\mathbf{z}_t = \mathbf{z}_{t-1}$ .

We construct directly  $\mathbf{s}_t = \mathbf{z}_t \parallel \mathbf{x}_t$ , and linearly project the vector into the mean and standard deviation of a isotropic Gaussian  $\mathbb{R}^8$ .  $\mathbf{s}_t$  is sampled from this Gaussian and then decoded with the `action` decoder into the reconstructed action in  $\mathcal{A}$ .

For each convolutional layer or transposed convolution layer, the tuple's values correspond to input channel, output channel, kernel size, stride, padding. With that notation in mind, the specific hyperparameters for each components are:

- `observation` encoder: This takes as input the  $10 \times 10 \times N_{obj} + 2$  and is implemented as a two 2D convolutional layers with ReLU activations, followed by a linear layer with a ReLU activation and a linear layer with no activation. The convolutional layers have hyperparameters (12, 64, 3, 1, 0) and (64, 64, 3, 1, 0) respectively. The first linear has input dimension  $6 \times 6 \times 64$  and output dimension 128. The second linear layer has input and output dimensions 128.
- `action` encoder: The actions are embedded with an embedding matrix with embedding dimension 128.
- `boundary posterior decoder` is a 5-layer 1-D causal temporal convolution [64] with  $\{(128, 128, 5, 1, 2) \times 5, (128, 2, 5, 1, 2)\}$  with ELU activation and batch normalization. The output is the logits for the  $\mathbf{m}_{1:T}$ .
- $f_{z-rnn-bwd}, f_{z-rnn-fwd}, f_{s-rnn-fwd}$ : GRU with cell size 128
- `action` decoder: The decoder is implemented as three linear layers with ReLU activations, followed by a linear layer with no activation. The (input, output) dimensions of these layers are as follows: (128, 128), (128, 128), (128, 128), (128, 5), where the last layer outputs logits over the actions.

## D.3 3D Navigation Environment

The architecture used in the 3D navigation environment is identical to that of the multi-task environment detailed above, with the exception that the `observation` encoder is changed to handle visual inputs. Specifically, the `observation` encoder is a 3-layer convolutional neural network parameters (3, 32, 5, 2, 0), (32, 32, 5, 2, 0), (32, 32, 4, 2, 0), followed two linear layers (7520, 128), (128, 128) with a ReLU activation in between.

## D.4 Hierarchical Reinforcement Learning

For all approaches, we parametrize the policy as a double dueling deep Q-network [56, 86, 83]. The parametrization of the Q-function consists of a state embedding followed by two linear layer heads for

the state-value  $V_\theta(s)$  function and the advantage  $A_\theta(s, a)$  function. Then the Q-function is computed as  $Q_\theta(s, a) = V_\theta(s) + A_\theta(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a' \in \mathcal{A}} A_\theta(s, a')$ .

The value function  $V_\theta(s)$  and advantage function  $A_\theta(s, a)$  are computed as linear layers with output dimension 1 and  $|\mathcal{A}|$  respectively on top of the embedding  $e(s)$  of the state  $s$ . Recall that the state  $s$  consists of two portions: the observation  $o$  of the grid or 3D environment, and the instruction  $i$  corresponding to the next object to pick up, where the instruction is not present in the demonstrations. The embedding  $e(s)$  of the state  $s$  is computed by embedding the observation  $e(o)$  with the `observation_encoder` defined in the previous sections, and embedding the instruction  $e(i)$  with a 16-dimensional embedding matrix. Then, the embedding  $e(s)$  is a final linear layer with output dimension 128 applied to a ReLU of the concatenation of  $e(o)$  and  $e(i)$ .

## E Hardware

All our experiments are conducted on a GeForce RTX 2080 Ti or a GeForce RTX A6000.

## F Hyperparameters

### F.1 Frame Prediction

For these set of experiments, we use the same architecture and hyperparameters as Kim et al. [39]. The hyperparameters used for both *Simple Colors* and *Conditional Colors* are:

- KL divergence weight  $\beta$  : 1.0
- MDL objective weight  $\lambda$  (fixed) : 0.1
- Mini-batch size: 512
- Training iteration: 30000
- Number of skills  $n_z = 10$

### F.2 Multi-task Grid World Segmentation

We use the following hyperparameters for skill learning from demonstrations on the multi-task grid world environment experiments.

- KL divergence weight  $\beta = 0$ :
  - We found that the compression objective readily provides strong and sufficient regularization for the latent code. Adding additional KL divergence often results in the model not being able to reconstruct the actions properly.
- MDL objective weight  $\lambda$ :
  - We use a adaptive scheduling that approximate dual gradient descent.  $\lambda$  is initialized to 0. After every gradient step,  $\lambda$  is increased by  $2.0 \times 10^{-5}$  if  $\mathcal{L}_{\text{ELBO}} \leq 0.05$  (Since  $\beta = 0$ , this is effectively  $\mathcal{L}_{\text{rec}}$ ); otherwise,  $\lambda$  is decreased by  $2.0 \times 10^{-5}$ . After each update, the value of  $\lambda$  is clipped to  $[0, 0.05]$ . We find this setting provides the most stable training.
- Mini-batch size: 64
- Training iteration: 20000
- Number of skills  $n_z = 10$
- $t_{\text{VQ}} = 0.1$
- Learning rate: 0.0005 with Adam Optimizer

### F.3 3D Visual Navigation Segmentation

We use the same hyperparameter as Multi-task grid world experiments except that the maximum  $\mathcal{L}_{\text{ELBO}}$  is capped at 0.001, i.e.,  $\mathcal{L}_{\text{ELBO}} \leq 0.001$ .

## F.4 Hierarchical Reinforcement Learning

We use the same hyperparameters for training the policy for all approaches. Specifically, these hyperparameter values are as follows:

- $\epsilon$ -greedy schedule: We linearly decay  $\epsilon$  from 1 to 0.01 over  $500K$  timesteps for dense reward settings and over  $5M$  timesteps for sparse reward settings in the multi-task grid world environment. In the 3D visual navigation environment, we decay over  $250K$  timesteps.
- Discount factor  $\gamma$ : We use  $\gamma = 0.99$  in all of the DQN updates.
- Maximum replay buffer size: 50K
- Minimum replay buffer size before updating: 500
- Learning rate: 0.0001 with the Adam optimizer [40]
- Batch size: 32
- Update frequency: every 4 timesteps
- Target syncing frequency: every 50K updates for multi-task grid world environment, every 30K updates for 3D navigation
- Gradient  $\ell_2$ -norm clipping: 10
- Marginal threshold  $\alpha = 0.001$

Since the demonstrations do not contain the instruction list observations, we set those to 0 in the demonstrations for the behavior cloning approach. Though we use  $\alpha = 0.001$  in all of the experiments for consistency, we also found that slightly higher values of  $\alpha$  yielded greater sample efficiency for LOVE.

## G Algorithm

The Lagrangian for the unconstrained optimization problem is:

$$\min_{\theta} \max_{\lambda \geq 0} \lambda \mathcal{L}_{\text{CL}}(\theta) + (\mathcal{L}_{\text{ELBO}}(\theta) - C).$$

In standard KKT condition, the dual variable is introduced on the ELBO, but the two problems are equivalent by inverting the dual variable. Since the only constraint on  $\lambda$  is that it is non-negative, this transformation does not change the optimal solution. Further note that we find that in many cases, choosing a fixed constant value for  $\lambda$  is sufficient for solving the problem.

---

**Algorithm 1 Learning Options via Compression (LOVE):** We highlight differences between our method and prior work (VTA) in blue.

---

```

1: Initialize model parameters  $\theta$ 
2: while not converged do
3:   Sample a trajectory  $(x_0, a_0, \dots, x_T)$  from the pre-collected experience  $\mathcal{D}$ 
4:   for  $t = 1, 2, \dots, T$  do
5:     Sample boundary conditioned on entire trajectory  $m_t \sim q_{\theta}(m_t \mid x_{1:t})$ 
6:     Sample skill from skill posterior  $z_t \sim q_{\theta}(z_t \mid m_t, x_{1:t}, a_{1:t})$ 
7:     Sample abstraction from abstraction posterior  $s_t \sim q_{\theta}(s_t \mid z_t, x_t)$ 
8:     Compute probability of the correct action from the decoder  $p_{\theta}(a_t \mid s_t)$ 
9:   end for
10:  Compute lower bound on maximum likelihood objective  $\mathcal{L}_{\text{ELBO}}$  according to Equation 4
11:  Compute compression objective  $\mathcal{L}_{\text{CL}}$  according to Equation 1
12:  Update  $\theta \leftarrow \theta - \eta \nabla(\mathcal{L}_{\text{ELBO}} + \lambda \mathcal{L}_{\text{CL}})$  and  $\lambda$ 
13: end while

```

---



---

**Algorithm 2** Executing skill  $z$ 

---

```
1: while skill has not terminated do  
2:   Compute state abstraction  $\mu_{s_t} = \mathbb{E}_{s \sim q_\theta(s_t | z, x_t)} [s]$   
3:   Take action  $a_t = \arg \max_a p_\theta(a | \mu_{s_t})$   
4:   Observe next state  $x_{t+1}$   
5:   Terminate if  $\mathbb{E}_{m_{t+1} \sim q_\theta(m_{t+1} | x_{1:t+1})} [m_{t+1}] > 0.5$   
6: end while
```

---

## H Visualization of Learned Skills

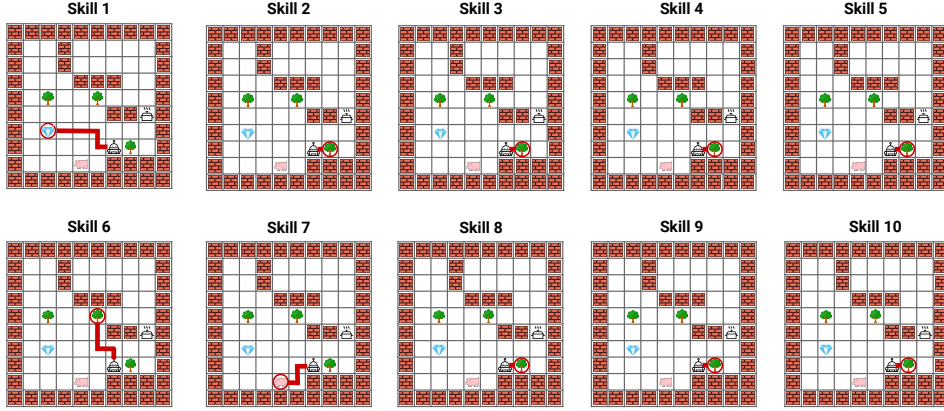


Figure 8: Skill visualization for all 10 skills of LOVE.

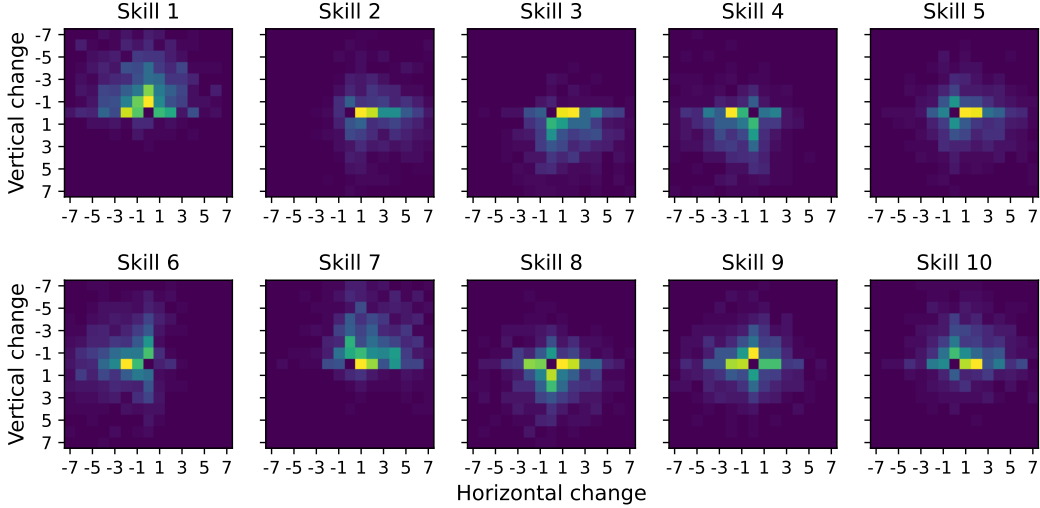


Figure 9: Heat map the agent's end position minus the agent's start position after applying each skill. In other words, this shows what direction the agent moves in after following each skill.

In Figure 8, we visualize the behavior of following each of LOVE's 10 skills on an example task. Each skill moves to and picks up an object, and the total set of skills covers 3 of the 4 object types in the task. These skills pick up both nearby objects, such as the tree in skill 2, as well as more distant objects, such as the diamond in skill 1. This contrasts DDO and VTA, which do not learn skills that move to and pick up objects.

To further understand what each of LOVE's skills does, we analyze whether there is a correlation between each skill and either the type of object it picks up, or the location of the object it picks up. We find that there appears to be little correlation between each skill and the type of object it picks up. Plotting the frequency of picking up each object type by skill shows that each skill picks up each object type with roughly the same frequency. Instead, there appears to be a correlation between each skill and the location of the object it picks up, illustrated in Figure 9. For example, skill 1 appears to specialize in moving to and picking up objects that are above the agent, while skill 7 tends to pick up objects that are up and to the right of the agent. It is also interesting to note that some skills are much more biased to move in certain direction (e.g., skill 1, 6, 7) while some appear to be more general (e.g., 4, 5, 8, 9) and move in any direction. Note in Figure 8, it may seem LOVE's skills seem short / redundant. This is due to the fact that LOVE's skills each pick up an object, so they

naturally appear short when the agent is close to the objects, as in Figure 8. However, when the agent is far away from the objects, the skills still pick up the objects and are much longer. From Figure 9, we can see from the heatmap that even the skills that appear to do same thing in Figure 8 have very distinct behaviors depending on the state they are in. Also, recall LOVE imposes a sparse distribution over skills (Section 5). After filtering, skill 3, 7 and 10 in Figure 8 are dropped since they have low marginal probability and are not used to describe a significant part of the trajectories. Figure 8 includes redundant skills LOVE does not use; the used skills in the sparse distribution have little redundancy.

## I Option Critic Results

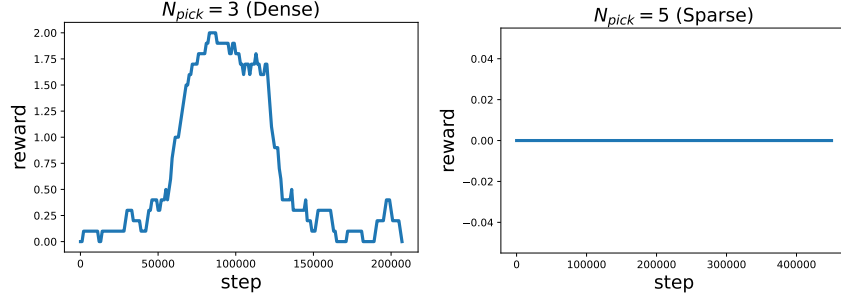


Figure 10: Results of running option critic [4] on 2 of the settings we considered in Figure 5.

Option critic [4] is a classical online HRL algorithm. We show the results of running option critic<sup>5</sup> on two of the four RL tasks we considered in Figure 10. The best return achieved over the training is shown in Table 4. We set the number of option to 8 following Bacon et al. [4] and leave other hyperparameters untouched.  $N_{\text{pick}} = 3$  with dense reward is the easily setting and  $N_{\text{pick}} = 5$  with sparse reward is the hardest setting. We see that in  $N_{\text{pick}} = 3$ , option critic is able to reach reward of 2 (maximum possible is 3) at around 1M environment steps but the performance deteriorates afterwards. In  $N_{\text{pick}} = 5$ , the algorithm fails to make any meaningful progress. These observations may suggest online HRL algorithms like option critic may be insufficient for solving these tasks.

	$N_{\text{pick}} = 3$ (Dense)	$N_{\text{pick}} = 5$ (Sparse)
Option-Critic	2.0	0.0
LOVE (Ours)	<b>3.0</b>	<b>0.7</b>

Table 4: Comparison between LOVE and Bacon et al. [4] on the two RL tasks we consider in this work. Each entry is the maximum average return achieved during the course of training for both algorithms. For  $N_{\text{pick}} = 3$  (Dense), 3.0 is the maximum possible return. For  $N_{\text{pick}} = 5$  (Sparse), 1.0 is the maximum possible return.

## J Comparison to Zhang et al. [90]

	<i>Simple Colors</i>		<i>Conditional Colors</i>		<i>Navigation</i>	
	MDL	LOVE	MDL	LOVE	MDL	LOVE
Precision	0.87	<b>0.99</b>	0.84	<b>0.99</b>	0.79	<b>0.90</b>
Recall	0.78	<b>0.85</b>	0.82	<b>0.83</b>	0.34	<b>0.94</b>
F1	0.82	<b>0.91</b>	0.83	<b>0.90</b>	0.48	<b>0.92</b>

Table 5: Comparison between LOVE and Zhang et al. [90] on the three segmentation tasks we consider in this work. We refer to Zhang et al. [90] as **MDL** in the table.

<sup>5</sup><https://github.com/lweitekamp/option-critic-pytorch>

Zhang et al. [90] learn open-loop skills, which do not condition on the state and therefore cannot adapt to different states. Further, the MDL used in Zhang et al. [90] is equivalent to the variational inference used by VTA (on a different graphical model), which can be seen as greedily compressing each skill independently, rather than compressing a whole trajectory as LOVE does. Table 5 reports segmentation results on the Color domain and grid world, akin to Tables 2 and 3. Zhang et al. [90] performs the same as VTA on the Color domains – as they both use variational inference and there is no state – which achieves lower precision / recall than LOVE. On the grid world navigation, Zhang et al. [90] fails to recover the boundaries, unlike LOVE, because skills that navigate to objects require observing the state. This problem is shared by all open-loop approaches.

## K Number of Initial Skills Ablation

	$K = 2$	$K = 5$	$K = 10$	$K = 15$	$K = 20$	$K = 30$	$K = 50$
Precision	0.27	0.80	0.90	0.96	0.96	0.95	0.95
Recall	0.53	0.91	0.94	0.96	0.95	0.96	0.93
F1	0.35	0.86	0.92	0.96	0.95	0.95	0.94

Table 6: Performance of LOVE on the grid navigation task with varying number of initial skills ( $K$ ).

The number of skills to extract from demonstrations is often not known a priori, so choosing an appropriate value of the number of initial skills before filtering  $K$  is important. Intuitively, we hypothesize that  $K$  can be set conservatively: LOVE requires at least a minimum number of skills in order to fit the behaviors in the demonstrations, but may be able to gracefully prune out skills if  $K$  is set too high via the filtering described in Section 5. We test this intuition by varying  $K \in \{2, 5, 10, 15, 20, 30, 50\}$  and measuring the segmentation performance on the grid world multi-task domain, and find that it appears to hold true in Table 6. For smaller values of  $K$ , such as  $K = 2$  and  $K = 5$ , LOVE’s F1 scores significantly degrade. However, performance remains high for all values where  $K$  is sufficiently large. Hence, we conservatively sizing  $K$  to be a large number.

## L Comparison to Different Regularizers

	entropy	num switch	VTA(3,5)	VTA(3,10)	VTA(5,5)	VTA(5,10)	LOVE
Precision	0.26	0.80	0.34	0.40	0.34	<b>0.95</b>	0.90
Recall	0.53	0.93	0.46	0.60	0.50	0.43	<b>0.94</b>
F1	0.35	0.86	0.39	0.48	0.40	0.37	<b>0.92</b>

Table 7: Comparison of different kinds of regularizers that have been used in the literature for skill segmentation. LOVE outperforms all significantly.

Several prior methods encourage skills to act for more timesteps or prevent skills from switching too frequently, which can similarly help avoid degenerate solutions like LOVE. We note that such prior methods do not necessarily yield skills that help learn new tasks, while the maximum of LOVE’s objective achieves the information-theoretic limit of compressing the trajectories (under a fixed function class), which is not done in other works. We empirically compare LOVE with several such methods on the segmentation of the multi-task grid world domain, including:

1. *Entropy*: This approach regularizes the entropy of skill marginal akin to Shankar et al. [72], which they refer to this as *parsimony*. Specifically, this approach adds a regularization term  $\mathcal{H}_{p_z^*}[z]$  to the objective  $\mathcal{L}_{\text{ELBO}}$ .
2. *Num switches*: This approach penalizes switching between skills, similar to Harb et al. [27], a variant of the Option-Critic framework [4]. Specifically, this approach adds a penalty  $n_s$  to the objective  $\mathcal{L}_{\text{ELBO}}$  equal to the number of switches (i.e., the number of times  $m_t = 1$ ) in a demonstration.

3.  $\text{VTA}(N_{\max}, l_{\max})$ : VTA includes a prior on its boundary variables, which effectively sets the maximum number of skills per episode to be  $N_{\max}$  and sets the maximum number of actions a skill can take to  $l_{\max}$ . We experiment with several settings of  $N_{\max}$  and  $l_{\max}$ , including those that use prior knowledge about the domain not used by LOVE.

The results are summarized in Table 7. We find that regularizing the number of switches performs fairly well, but LOVE achieves higher F1 than other approaches that regularize the skills.