# LISA: Learning Interpretable Skill Abstractions from Language

Anonymous Author(s) Affiliation Address email

# Abstract

Learning policies that effectually utilize language instructions in complex, multi-1 task environments is an important problem in imitation learning. While it is 2 possible to condition on the entire language instruction directly, such an approach 3 4 could suffer from generalization issues. To encode complex instructions into skills that can generalize to unseen instructions, we propose Learning Interpretable Skill 5 Abstractions (LISA), a hierarchical imitation learning framework that can learn 6 diverse, interpretable skills from language-conditioned demonstrations. LISA uses 7 vector quantization to learn discrete skill codes that are highly correlated with 8 language instructions and the behavior of the learned policy. In navigation and 9 robotic manipulation environments, LISA outperforms a strong non-hierarchical 10 baseline in the low data regime and is able to compose learned skills to solve 11 tasks containing unseen long-range instructions. Our method demonstrates a more 12 13 natural way to condition on language in sequential decision-making problems and achieve interpretable and controllable behavior with the learned skills. 14

# 15 1 Introduction

Intelligent machines should be able to solve a variety 16 of complex, long-horizon tasks in an environment 17 and generalize to novel scenarios. In the sequential 18 decision-making paradigm, provided expert demon-19 strations, an agent can learn to perform these tasks 20 via multi-task imitation learning (IL). As humans, it 21 is desirable to specify tasks to an agent using a con-22 venient, yet expressive modality and the agent should 23 24 solve the task by taking actions in the environment. 25 There are several ways for humans to specify tasks to an agent, such as task IDs, goal images, and goal 26 demonstrations. However, these specifications tend 27 to be ambiguous, require significant human effort, 28 and can be cumbersome to curate and provide at test 29 time. One of the most natural and versatile ways for 30 31 humans to specify tasks is via natural language.

- 32 The goal of language-conditioned IL is to solve tasks
- <sup>33</sup> in an environment given language-conditioned trajec-
- 34 tories at training time and a natural language instruc-
- <sup>35</sup> tion at test time. This becomes challenging when the
- 36 task involves completing several sub-tasks sequen-
- <sup>37</sup> tially, like the example shown in Figure 1. A crucial



Figure 1: **Overview of LISA.** Given a language instruction, our method learns discrete skill abstractions z, picked from a codebook C. The policy conditioned on the skill code learns to execute distinct behaviors and solve different sub-goals. See GIF.

<sup>38</sup> step towards solving this problem is exploiting the inherent hierarchical structure of natural language.

For example, given the task specification "pull the handle and move black mug right", we can split it into two independent *skills*, *i.e.* "pull the handle" and "move black mug right". If we are able to decompose the problem of solving these complex tasks into learning skills, we can compose these skills to generalize to unseen tasks in the future. This is especially useful in the low-data regime, since *we may not see all possible tasks given the limited dataset, but may see all the constituent sub-tasks.* One of the main goals of language conditioned IL is to utilise language effectively so that we can learn skills as the building blocks of complex behaviours.

Utilising language effectively to learn skills is a non-trivial problem and raises several challenges. (i) 46 The process of learning skills from language-conditioned trajectories is unsupervised as we may not 47 have knowledge about which parts of the trajectory corresponds to each skill. (ii) We need to ensure 48 that the learned skills are useful, *i.e.* encode behavior that can be composed to solve new tasks. (iii) 49 We would like the learned skills to be interpretable by humans, both in terms of the language and 50 51 the behaviours they encode. There are several benefits of interpretability. For example, it allows us to understand which skills our model is good at and which skills it struggles with. In safety critical 52 settings such as robotic surgery or autonomous driving, knowing what each skill does allows us 53 to pick and choose which skills we want to run at test time. It also provides a visual window into 54 a neural network policy which is extremely desirable [54]. There have been prior works such as 55 [38, 47, 13] that have failed to address these challenges and condition on language in a monolithic 56 fashion without learning skills. As a result, they tend to perform poorly on long-horizon composition 57 tasks such as the one in Figure 1. 58

To this end, we propose Learning Interpretable Skill Abstractions from language (LISA), a hierarchi-59 cal imitation learning framework that can learn interpretable skills from language-conditioned 60 61 offline demonstrations. LISA uses a two-level architecture – a skill predictor that predicts quantized skill codes and a policy that uses these skill codes to predict actions. The discrete skill codes learned 62 from language are interpretable (see Figure 3) and can be composed to solve long-range tasks. Using 63 skill quantization maximizes code reuse and enforces a bottleneck to pass information from the 64 language to the policy, enabling unsupervised learning of interpretable skills. We perform experiments 65 on grid world navigation and robotic manipulation tasks and show that our hierarchical method can 66 outperform a strong non-hierarchical baseline in the low-data regime. We analyse these skills qualita-67 tively and quantitatively and find them to be highly correlated to language and behaviour. Finally, 68 using these skills to perform long-range composition tasks on a robotic manipulation environment 69 results in performance that is nearly 2x better than the non-hierarchical version. 70

- 71 Concretely, our contributions are as follows:
- We introduce LISA, a novel hierarchical imitation framework conditioned on language to decompose complex tasks into skills.
- We demonstrate the effectiveness of our approach in the low-data regime where its crucial to break down complex tasks to generalize well.
- We show our method performs well in long-range composition tasks where we may need to
   perform multiple skills sequentially.
- We also show that the learned skills are highly correlated to language and behaviour and can easily be interpreted by humans.

# 80 2 Related Work

### 81 2.1 Imitation Learning

Imitation learning (IL) has a long history, with early works using behavioral cloning [41-43] to 82 learn policies via supervised learning on expert demonstration data. Recent methods have shown 83 significant improvements via learning reward functions [21] or Q-functions [18] from expert data 84 to mimic expert behavior. Nevertheless, these works typically consider a single task. An important 85 problem here is multi-task IL, where the imitator is trained to mimic behavior on a variety of training 86 tasks with the goal of generalizing the learned behaviors to test tasks. A crucial variable in the 87 multi-task IL set-up is how the task is specified, e.g vectorized representations of goal states [37], task 88 IDs [24], and single demonstrations [56, 14, 16, 57]. In contrast, we focus on a multi-task IL setup 89 with task-specification through language, one of the most natural and versatile ways for humans to 90 communicate desired goals and intents. 91

#### 92 2.2 Language Grounding

Several prior works have attempted to ground language with tasks or use language as a source of instructions for learning tasks with varying degrees of success ([32, 55, 4, 39, 5]). [27] is a good reference for works combining language with sequential-decision making.

But apart from a few exceptions, most algorithms in this area use the language instruction in a 96 monolithic fashion and are designed to work for simple goals that requires the agent to demonstrate a 97 single skill. ([40, 9, 20, 6]) or tasks where each constituent sub-goal has to be explicitly specified ([10, 9, 20, 6])98 46, 34, 3, 52, 50, 17, 35, 31]). Some recent works have shown success on using play data [28] 99 or pseudo-expert data such as LOReL [38] and CLIPORT [47]. LOReL and CLIPORT are not 100 hierarchical techniques. [28] can be interpreted as a hierarchical technique that generates latent 101 102 sub-goals as a function of goal images, language instructions and task IDs but the skills learned by LISA are purely a function of language and states alone and do not require goal images or task 103 IDs. [23, 22] and [49] are some examples of works that use a two-level architecture for language 104 conditioned tasks but neither of these methods learn skills that are interpretable. 105

# 106 2.3 Latent-models and Hierarchical Learning

Past works have attempted to learn policies conditioned on latent variables and some of them can 107 be interpreted as hierarchical techniques. For example, [15] learns skills using latent variables that 108 visit different parts of the environment's state space. [45] improved on this by learning skills that 109 were more easily predictable using a dynamics model. But these fall more under the category of skill 110 discovery than hierarchical techniques since the skill code is fixed for the entire trajectory, as is the 111 case with [15]. [29] and [26] are other works that use a latent-variable approach to IL. But these 112 approaches don't necessarily learn a latent variable with the intention of breaking down complex 113 114 tasks into skills. With LISA, we sample several skills per trajectory with the clear intention of each skill corresponding to completing a sub-task for the whole trajectory. Also, none of the methods 115 mentioned here condition on language. 116

There has been some work on hierarchical frameworks for RL to learn high-level action abstractions, called options [51], such as [25, 58, 36] but these works are not goal-conditioned. Unlike LISA, these works don't use language and the options might lack diversity and not correspond to any concrete or interpretable skills. Furthermore, none have used the VQ technique to learn options and often suffer from training instabilities.

# 122 **3** Approach

The key idea of LISA is to learn quantized skill representations that are informative of both language and behaviors, which allows us to break down high-level instructions, specified via language, into discrete, interpretable and composable codes (see Fig. 5 and Fig. 7 for visualizations). These codes enable learning explainable and controllable behaviour, as shown in Fig. 1 and Fig. 3.

Section 3.1 describes the problem formulation, an overview of our framework, and presents our
 language-conditioned model. Section 3.2 provides details on the training approach.

#### 129 3.1 Language-conditioned Skill Learning

#### 130 3.1.1 Problem Setup

We consider general multi-task environments, represented as a task-augmented Markov decision process (MDP) with a family of different tasks  $\mathcal{T}$ . A task  $\mathcal{T}_i$  can be a union of other tasks in  $\mathcal{T}$ . For example, in a navigation environment, a task could be composed of two or more sub-tasks - "pick up ball", "open door" - in any hierarchical order. S,  $\mathcal{A}$  represent state and action spaces. We assume that each full task has a *single* natural language description  $l \in L$ , where L represents the space of language instructions. Any sub-goals for the task are encoded within this single language instruction.

<sup>137</sup> We assume access to an offline dataset  $\mathcal{D}$  of trajectories obtained from an optimal policy for a <sup>138</sup> variety of tasks in an environment with only their language description available. Each trajectory

 $\tau^i = (l^i, \{(s_1^i, a_1^i), (s_2^i, a_2^i), ..., (s_T^i, a_T^i)\})$  consists of the language description and the observations

140  $s_t^i \in S$ , actions  $a_t^i \in A$  taken over T timesteps. The trajectories are not labeled with any rewards.

Our aim is to predict the expert actions  $a_t$ , given a language instruction and past observations.

Note that each trajectory in the training dataset can comprise of any number of sub-tasks. For example, we could have a trajectory to "open a door" and another to "pick up a ball and close the door" in the

# Algorithm 1 Training LISA

**Input:** Dataset  $\mathcal{D}$  of language-paired trajectories **Input:** Num skills K and horizon H 1: Initialize skill predictor  $f_{\phi}$ , policy  $\pi_{\theta}$ 2: Vector Quantization op  $\mathbf{q}(\cdot)$ while not converged do 3: Sample  $\tau = (l, \{s_0, s_1, s_2 \dots s_T\}, \{a_0, a_1, a_2 \dots a_T\})$ 4: Sample  $\tau = (\iota, \{s_0, s_1, s_2\})$ Initialize  $S = \{s_0\}$ for  $k = 0 .. \lfloor \frac{T}{H} \rfloor$  do  $z \leftarrow \mathbf{q}(f_{\phi}(l, S))$ for step t = 1 .. H do 5:  $\triangleright$  List of seen states 6: ▷ Sample a skill every H steps 7: 8:  $a_{kH+t} \leftarrow \pi_{\theta}(z, S[:-H])$  $S \leftarrow S \cup \{s_{kH+t}\}$ end for 9:  $\triangleright$  Use recent H steps ▷ Append last state 10: 11: 12: Train  $f_{\phi}, \pi_{\theta}$  using objective  $\mathcal{L}_{\text{LISA}}$ 13: end for 14: end while

training data. With LISA we aim to solve the task "open a door and pick up the ball" at test time even though we haven't seen this task at training time. In a trajectory with multiple sub-tasks, the training dataset **does not** give us information about where one sub-task ends and where another one begins.

LISA must learn how to identify and stitch together these sub-tasks learned during training, in order to solve a new language instruction such as the one shown in Fig. 1 at test time.

## 149 3.1.2 Hierarchical Skill Abstractions

<sup>150</sup> We visualize the working of LISA in Figure 2.

Our framework consists of two modules: a skill predictor  $f : L \times S \to C$  and a policy  $\pi :$  $S \times C \to A$ . Here,  $C = \{z^1, \dots, z^K\}$  is a learnable codebook of K quantized skill codes.

Our key idea is to break learning behavior from language in two stages: 1) Learn discrete codes z, representing skills, from the full-language instruction to decompose the task into smaller sub-goals 2) Learn a policy  $\pi$  conditioned only on these discrete codes. In LISA, both stages are trained end-to-end.

Given an input  $\tau = (l, \{s_t, a_t\}_{t=1}^T)$ , the skill predictor f predicts a skill code at a timestep 162 163 t as  $\tilde{z} = f(l, (s_t, s_{t-1}, ...))$ . These codes are 164 discretized using a vector quantization operation 165  $\mathbf{q}(\cdot)$  that maps a code  $\tilde{z}$  to its closest codebook 166 entry  $z = \mathbf{q}(\tilde{z})$ . The quantization operation  $\mathbf{q}(\cdot)$ 167 helps in learning discrete codes and acts as a 168 bottleneck on passing language information. We 169 detail its operation in Sec. 3.2. 170



Figure 2: LISA Architecture: The skill predictor f gets the language instruction and a sequence of observations as the input, processed through individual encoders. It predicts quantized skill codes z using a learnable cookbook C, that encodes different sub-goals, and passes them to the policy  $\pi$ . LISA is trained end-to-end.

The chosen skill code z, is persisted for H timesteps where H is called the horizon. After H timesteps, the skill predictor is invoked again to predict a new skill. This enforces the skill to act as a temporal abstraction, i.e. options [51]. The policy  $\pi$  predicts the action  $a_t$  at each timestep t conditioned on a single skill code z that is active at that timestep. For  $\pi$  to correctly predict the original actions, it needs to use the language information encoded in the skill codes.

LISA learns quantized skill codes in a codebook instead of continuous embeddings as this encourages
 reusing and composing these codes together to pass information from the language input to the
 actual behavior. Our learnt discrete skill codes adds interpretability and controllability to the policy's
 behavior.



Figure 3: Behavior with fixed LISA options. We show the word clouds and the behavior of the policy obtained by using a fixed skill code z = 14 for an entire episode. We find that this code encodes the skill "closing the drawer", as indicated by the word cloud. The policy executes this skill with a high degree of success when conditioned on this code for the entire trajectory, across different environment initializations and seeds.

# 180 3.2 Training LISA

**Learning Discrete Skills.** LISA uses Vector Quantization (VQ), inspired from [53]. It is a natural and widely-used method to map an input signal to a low-dimensional discrete learnt representation. VQ learns a codebook  $C \in \{z^1, \ldots, z^K\}$  of K embedding vectors. Given an embedding  $\tilde{z}$  from the skill predictor f, it maps the embedding to the closest vector in the codebook:

$$z = \mathbf{q}(\tilde{z}) =: \underset{z^k \in \mathcal{C}}{\operatorname{arg\,min}} \|\tilde{z} - z^k\|_F$$

This can be classically seen as learning K cluster centers via k-means [19].

Backpropagation through the non-differentiable quantization operation is achieved by a straightthrough gradient estimator, which simply copies the gradients from the decoder to the encoder, such that the model and codebook can be trained end-to-end.

<sup>185</sup> VQ enforces each learnt skill z to lie in C, which can be thought as learning K prototypes or cluster <sup>186</sup> centers for the language embeddings using the seen states. This acts as a bottleneck that efficiently <sup>187</sup> decomposes a language instruction into sub-parts encoded as discrete skills.

LISA Objective. LISA is trained end-to-end using an objective  $\mathcal{L}_{\text{LISA}} = \mathcal{L}_{\text{BC}} + \lambda \mathcal{L}_{\text{VQ}}$ , where  $\mathcal{L}_{\text{BC}}$  is the behavior-cloning loss on the policy  $\pi_{\theta}$ ,  $\lambda$  is the VQ loss weight and  $\mathcal{L}_{\text{VQ}}$  is the vector quantization loss on the skill predictor  $f_{\phi}$  given as:

$$\mathcal{L}_{\mathrm{VQ}}(f) = \mathbb{E}_{\tau}[\|\mathrm{sg}\left[\mathbf{q}(\tilde{z})\right] - \tilde{z}\|_{2}^{2}]$$
(1)

191 with  $\tilde{z} = f_{\phi}(l, (s_t, s_{t-1}, ..)).$ 

sg  $[\cdot]$  denotes the stop-gradient operation.  $\mathcal{L}_{VQ}$  is also called *commitment loss*. It minimizes the conditional entropy of the skill predictor embeddings given the codebook vectors, making the embeddings stick to a single codebook vector.

<sup>195</sup> The codebook vectors are learnt using an exponential moving average update, same as [53].

Avoiding language reconstruction. LISA avoids auxiliary losses for language reconstruction and it's not obvious why the skill codes are properly encoding language. It's known that given a signal Xand a code Z. Reconstructing the signal  $\tilde{X} = f(Z)$  using cross-entropy loss amounts to maximizing a lower bound to the Mutual Information (MI) I(X, Z) between X and Z [1, 7]. In our case, we can write the MI between the skill codes and language using entropies as:  $I(z, l) = H(z) - H(z \mid l)$ , whereas methods that attempt to reconstruct language apply the following decomposition: I(z, l) = $H(l) - H(l \mid z)$  (where H(l), the entropy of language, is independent from LISA's skill encoder).

Thus we can avoid language reconstruction via cross-entropy loss by maximizing I(z, l) directly. In LISA,  $\mathcal{L}_{vq} = -H(z \mid l)$ , and we don't observe a need to place a constraint on H(z) as the codes are diverse, needing to encode enough information to correctly predict the masked actions.<sup>1</sup>

As a result, LISA can maximize the MI between the learnt skills and languages without auxiliary losses and enforcing only  $\mathcal{L}_{vq}$  on the skill codes. We empirically estimate the MI and find that our experiments confirm this in Sec 4.5.

<sup>&</sup>lt;sup>1</sup>In experiments, we tried enforcing a constraint on H(z) by using extra InfoNCE loss term without success.

### 209 3.2.1 LISA Implementation

210 LISA can be be implemented using different network architectures, such as Transformers or MLPs.

In our experiments, we use Transformer architectures with LISA, but we find that out method is effective even with simple architectures choices such as MLPs, as shown in the appendix section F.4. Even when using Transformers for both the skill predictor and the policy network, our compute requirement is comparable to the non-hierarchical Flat Transformer policy as we can get away with using fewer layers in each module.

Language Encoder. We use a pre-trained DistilBERT [44] encoder to generate language embeddings from the text instruction. We further fine-tune the language encoder to the vocabulary of the environment. We use the full language embedding for each word token, and not a pooled representation of the whole text.

Observation Encoder. For image observations, we use convolution layers to generate embeddings.
 For simple state representations, we use MLPs.

**Skill Predictor.** The skill predictor network f is implemented as a small Causal Transformer network that takes in the language embeddings and the observation embeddings at each time step. The language embeddings are concatenated at the beginning of the observation embeddings before being fed into the skill predictor. The network applies a causal mask hiding the future observations.

**Policy Network.** Our policy network  $\pi$ , also implemented as a small Causal Transformer inspired by Decison Transformer (DT) [11]. However, unlike DT, our policy is not conditioned on any reward signal, but on the skill code. The sequence length of  $\pi$  is the horizon H of the skills which is much smaller compared to the length of the full trajectory.

**Flat Baseline.** Our flat baseline is implemented similar to LISA, but without a skill predictor network. The policy here is a Causal Transformer that directly takes the language instruction and past observations as inputs to predict the policy. We found this baseline to be in-efficient at handling long-range language instructions, needing sequence lengths of 1000 on complex environments such as BabyAI-BossLevel in our experiments.

 Table 1: Imitation Results: We show our success rates (in %) compared to the original method and a flat non-hierarchical baseline on each dataset. LISA outperforms all other methods in the low-data regime, and reaches similar performance as the number of demonstrations increases. Best method shown in bold.

 Task
 | Num Demos | Original | Elat Baseline | LISA

Task	Num Demos	Original	Flat Baseline	LISA
BabyAI GoToSeq BabyAI GoToSeq BabyAI GoToSeq	$\begin{array}{c c}1k\\10k\\100k\end{array}$	$\begin{array}{c} 33.3 \pm 1.3 \\ 40.4 \pm 1.2 \\ 47.1 \pm 1.1 \end{array}$	$\begin{array}{c} 49.3 \pm 0.7 \\ 62.1 \pm 1.2 \\ 74.1 \pm 2.3 \end{array}$	$59.4 \pm 0.9 \\ 65.4 \pm 1.6 \\ 77.2 \pm 1.7$
BabyAI SynthSeq BabyAI SynthSeq BabyAI SynthSeq	$\begin{vmatrix} 1k \\ 10k \\ 100k \end{vmatrix}$	$\begin{array}{c} 12.9 \pm 1.2 \\ 32.6 \pm 2.5 \\ 40.4 \pm 3.3 \end{array}$	$\begin{array}{c} 42.3 \pm 1.3 \\ 52.1 \pm 0.5 \\ 64.2 \pm 1.3 \end{array}$	$\begin{array}{c} {\bf 46.3 \pm 1.2} \\ {\bf 53.3 \pm 0.7} \\ {\bf 61.2 \pm 0.6} \end{array}$
BabyAI BossLevel BabyAI BossLevel BabyAI BossLevel	$\begin{array}{c c}1k\\10k\\100k\end{array}$	$\begin{array}{c} 20.7 \pm 4.6 \\ 28.9 \pm 1.3 \\ 45.3 \pm 0.9 \end{array}$	$\begin{array}{c} 44.5\pm 3.3\\ \textbf{60.1}\pm \textbf{5.5}\\ \textbf{72.0}\pm \textbf{4.2} \end{array}$	$\begin{array}{c} {\bf 49.1 \pm 2.4} \\ {\bf 58 \pm 4.1} \\ {\bf 69.8 \pm 3.1} \end{array}$
LOReL - States (fully obs.) LOReL - Images (partial obs.)	$50k \\ 50k$	$6 \pm 1.2^{*}$ 29.5 ± 0.07	$33.3 \pm 5.6 \\ 15 \pm 3.4$	$\begin{array}{c} 66.7 \pm 5.2 \\ 40 \pm 2.0 \end{array}$

# **235 4 Experiments**

In this section, we evaluate LISA on grid-world navigation and robotic manipulation tasks. We compare the performance of LISA with a strong non-hierarchical baseline in the low-data regime. We then analyse our learnt skill abstractions in detail – what they represent, how we can interpret them and how they improve performance on downstream composition tasks. Finally, we show ablation studies on important hyperparameters and architecture choices.

#### 241 4.1 Datasets

Several language-conditioned datasets have been curated as off late. [46, 48, 13, 38, 33, 2, 10, 12] are some examples. Nevertheless, a lot of these datasets focus on complex-state representations and navigation in 3D environments, making them challenging to train on and qualitatively analyze our skills as shown in Fig. 3. We found BabyAI, a grid-world navigation environment and LOReL, a
robotic manipulation environment as two diverse test beds that were very different from each other
and conducive for hierarchical skill learning as well as detailed qualitative and quantitative analysis
of our learned skills and we use them for our experiments.

**BabyAI Dataset.** The BabyAI dataset [13] contains 19 levels of increasing difficulty where each 249 level is set in a grid world and an agent sees a partially observed ego-centric view in a square of size 250 7x7. The agent must learn to perform various tasks of arbitrary difficulty such as moving objects 251 between rooms, opening or closing doors, etc. all with a partially observed state and a language 252 instruction. The language instructions for easy levels are quite simple but get exponentially more 253 challenging for harder levels and contain several skills that the agent must complete in sequence 254 (examples in appendix section C.1). The dataset provides 1 million expert trajectories for each of the 255 19 levels, but we use 0.1 - 10% of these trajectories to train our models. We evaluate our policy on 256 257 100 different instructions from the gym environment for each level, which very likely contains unseen environments and language instructions given the limited data we use for training. More details about 258 this dataset can be found in the appendix and in the BabyAI paper. 259

LOReL Sawyer Dataset. This dataset [38] consists of *pseudo-expert* trajectories (play data) collected 260 from a replay buffer of a random RL policy and has been labeled with post-hoc crowd-sourced 261 language instructions. Hence, the trajectories complete the language instruction provided but may not 262 263 necessarily be optimal. Play data is inexpensive to collect [30] in the real world and it is important for algorithms to be robust to such datasets as well. However, due to the randomness in the trajectories, 264 this makes the dataset extremely difficult to use in a behavior cloning (BC) setting. Despite this, we 265 are able to achieve good performance on this benchmark and are able to learn some very useful skills. 266 The LOReL Sawyer dataset contains 50k trajectories of length 20 on a simulated environment with a 267 Sawyer robot. We evaluate on the same set of 6 tasks that the original paper does for our results in 268 Table 10: close drawer, open drawer, turn faucet right, turn faucet left, move black mug right, move 269 white mug down. More details can be found in the appendix section C.2 and in the LOReL paper. 270

#### 271 4.2 Baselines

**Original.** These refer to the baselines from the original paper for each dataset. For BabyAI, we 272 trained their non-hierarchical RNN based method on different number of trajectories. Similarly, on 273 LOReL we compare with the performance of language-conditioned BC. The original LOReL method 274 uses a planning algorithm on a learned reward function to get around the sub-optimal nature of the 275 trajectories. We found the BC baseline as a more fair comparison, as LISA is trained using BC as 276 well. Nonetheless, we compare with the original LOReL planner in Section 4.7 for composition 277 tasks. LOReL results in Table 10 refer to the performance on the 6 seen instructions in the LOReL 278 evaluation dataset, same as ones reported in the original paper. 279

**Flat Baseline.** We implement a non-hierarchical baseline using Transformers, the details of which are in section 3.2.1.

# 4.3 How does the performance of LISA compare with non-hierarchical baselines in a low-data regime?

We consider three levels from the BabyAI environment and the LOReL Sawyer environment. From 284 the BabyAI environment, we consider the GoToSeq, SynthSeq and BossLevel tasks since they are 285 challenging and require performing several sub-tasks one after the other. Since these levels contain 286 instructions that are compositional in nature, when we train on limited data, the algorithm must be 287 able to learn the skills which form these complex instructions to generalize well to unseen instructions 288 at test time. Our results are given in Table 10. We train the models on a random sample of 1k, 10k 289 and 100k trajectories on the BabyAI dataset and 50k trajectories on the LOReL dataset. We use 290 more data from the LOReL dataset because of the sub-optimal nature of the trajectories. On all 291 the environments, our method is competitive to or outperforms the strong non-hierarchical decision 292 transformer baseline. The gap grows larger as we reduce the number of trajectories we train on, 293 indicating that our method is able to leverage the common sub-task structures better and glean more 294 information from limited data. As expected, with larger amounts of training data, it becomes hard to 295 beat the flat baseline since the Transformer sees more compositions during training and can generalize 296 297 better at test time [8]. As mentioned above, we evaluate on the same 6 seen instructions the original

<sup>\*</sup>We optimized a language-conditioned BC model following the details in the appendix of the LOReL paper to the best of our abilities but could not get better performance.



Figure 4: LISA Skill Heat map on LOReL. The sparsity and the bright spots show that specific options correspond to specific language tokens and by extension, skills

LOReL paper did. We also evaluated the performance on varying the language instructions, similar to the original paper with results in appendix section E.

We were pleasantly surprised that LISA is 2x better than the flat baseline on LOReL tasks, reaching 40% success rate despite the sub-optimal nature of the data. One explanation for this is that the discrete skill codes are able to capture *different ways of doing the same task*, thereby allowing LISA to *learn an implicit multi-modal policy*. This is not possible with the flat version as it has no way to compartmentalize these noisy trajectories, and perhaps tends to overfit on this noisy data, leading to

305 performance degradation.

## 306 4.4 What skills does LISA learn? Are they diverse?

To answer this question, we analyse the skills produced 307 308 by LISA and the language tokens corresponding to each 309 skill. We plot a heat map in Figure 4 corresponding to the correlation between the language tokens and skill codes. 310 Here, we plot the map corresponding to the LOReL dataset. 311 From the figure, we can see that certain skill codes cor-312 respond very strongly to certain language tokens and by 313 extension, tasks. We also see the sparse nature of the heat 314 315 maps which indicates that each skill corresponds to distinct language tokens. We also plot word clouds corresponding 316 to four different options in the LOReL environment in 317 figure 5 and we notice that different options are triggered 318 by different language tokens. From the figure, it is clear 319 that the skill on the top left corner corresponds to *close* 320 321 the drawer and the skill on the top right corresponds to turn faucet left. Similar word clouds and heat maps for the 322 BabyAI environments are in the appendix section **B.3**. 323



Figure 5: Word clouds on LOReL: We show the most correlated words for 4 different learnt skill codes on LOReL. We can see that the codes represent interpretable and distinguishable skills. For e.g, the code on the top left corresponds to closing the drawer. (note that container is a synonym for drawer in the LOREL dataset)

#### 324 4.5 Do the skills learned by LISA correspond to interpretable behavior?

We have seen that the different skills 325 correspond to different language to-326 kens, but do the policies conditioned 327 on these skills behave according to the 328 language tokens? To understand this, 329 we fix the skill code for the entire tra-330 jectory and run the policy i.e. we are 331 shutting off the skill predictor and al-332 ways predicting the same skill for the 333 entire trajectory. As we can see from 334 335 the word cloud and the corresponding GIF in Figure 3, the behaviour for 336 skill code 14 is exactly what we can 337 infer from the language tokens in the 338



Figure 6: **MI between language and skill codes:** We show the Mutual Information over training iterations for various settings of LISA on the BabyAI BossLevel environment

word cloud – *close the drawer*. More such images and GIFs can be found in the appendix section **B.5**.

#### **4.6** Why do the skills learned by LISA have such a strong correlation to language?

As mentioned in section 3.2, the *commitment loss* from VQ acts as a way to increase the MI between 341 the language and the skill codes during training. This allows the codes to be highly correlated with 342 language without any reconstruction losses. To analyze this, we plot the MI between the options 343 and the language during training on the BabyAI BossLevel with 1k trajectories and the plot can be 344 seen in figure 6. The plots show the MI increasing over training for a wide range of settings as we 345 vary the number of skills and the horizon. In the ablation studies below, we report the success rate 346 corresponding to each of these curves and we notice that there's almost a direct correlation with 347 increasing MI and task performance. This is very encouraging since it clearly shows that the skills 348 are encoding language and that directly impacts the performance of the behavior cloning policy. 349

#### 4.7 Can we use the learned skills to perform new composition tasks?

To test our composition performance, Table 2: LISA Composition Results: We show our performance 351 we evaluate on LOReL composition 352 on the LOReL Sawyer environment compared to baselines tasks using images. To this end, we Method Success Rate (in %) 353 handcraft 15 composition instructions, 354 Flat  $13.66 \pm 1.10$ some of which are from the LOReL 355 LOReL Planner  $18.38 \pm 2.18$ training data and some of which are 356 LISA (Ours)  $\textbf{21.06} \pm \textbf{2.13}$ 

<sup>357</sup> unseen. We have listed these instruc-

tions in the appendix table 4 but one such example is "*pull the handle and move black mug down*". As we can see, over 10 different runs, our performance is nearly 2x that of the non-hierarchical baseline. We also compare with the original LOReL planner on these composition tasks and we notice that we perform slightly better despite them having access to a reward function and a dynamics model trained on a lot more data. We set the time horizon to 40 from the usual 20 for all the methods while performing these experiments because of the compositional nature of the tasks.

Note that results in Table 10 show compositionality performance on the BabyAI dataset as we train with 0.1%-10% of the data. We evaluate on the gym environment generating any possible language instruction from the BabyAI grammar for a level and come across several unseen compositions at evaluation time.

# **368 5** Limitations and Future Work

We present LISA, a hierarchical imitation learning framework that can be used to learn interpretable skill abstractions from language-conditioned expert demonstrations. We showed that the skills are diverse and can be used to solve long-range language tasks and that our method outperforms a strong non-hierarchical baseline in the low-data regime.

However, there are several limitations to LISA and plenty of scope for future work. One limitation of LISA is that there are several hyperparameters to tune that may affect performance like the number of options and the horizon for each option. It certainly helps to have a good idea of the task to decide these hyperparameters even though the ablations show that the method is fairly robust to these choices. Its also useful to learn the horizon for each skill by learning a termination condition and we leave this for future work.

Although our method has been evaluated on the language-conditioned imitation learning setting, its not difficult to modify this method to make it work for image goals or demos, and in the RL setting as well. Its interesting to see if the vector quantization trick can be used to learn goal-conditioned skills in a more general framework.

# 383 References

9.5

Felix Agakov and David Barber. Variational information maximization for neural coding. In Nikhil Ranjan
 Pal, Nik Kasabov, Rajani K. Mudi, Srimanta Pal, and Swapan Kumar Parui, editors, *Neural Information Processing*, pages 543–548, Berlin, Heidelberg, 2004. Springer Berlin Heidelberg. ISBN 978-3-540-30499-

- 387
- [2] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen
   Gould, and Anton van den Hengel. Vision-and-language navigation: Interpreting visually-grounded
   navigation instructions in real environments. In *Proceedings of the IEEE Conference on Computer Vision* and Pattern Recognition (CVPR), 2018. 6

- [3] Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen
   Gould, and Anton van den Hengel. Vision-and-language navigation: Interpreting visually-grounded
   navigation instructions in real environments, 2018. 3
- [4] Dilip Arumugam, Siddharth Karamcheti, Nakul Gopalan, Lawson Wong, and Stefanie Tellex. Accurately and efficiently interpreting human-robot instructions of varying granularities. *Robotics: Science and Systems XIII*, Jul 2017. doi: 10.15607/rss.2017.xiii.056. URL http://dx.doi.org/10.15607/RSS.
   2017.XIII.056.3
- [5] Dilip Arumugam, Siddharth Karamcheti, Nakul Gopalan, Edward C. Williams, Mina Rhee, Lawson L.
   Wong, and Stefanie Tellex. Grounding natural language instructions to semantic goal representations for abstraction and generalization. *Auton. Robots*, 43(2):449–468, feb 2019. ISSN 0929-5593. doi: 10.1007/s10514-018-9792-8. URL https://doi.org/10.1007/s10514-018-9792-8. 3
- [6] Valts Blukis, Nataly Brukhim, Andrew Bennett, Ross A. Knepper, and Yoav Artzi. Following high-level
   navigation instructions on a simulated quadcopter with imitation learning, 2018. 3
- [7] Malik Boudiaf, Jérôme Rony, Imtiaz Masud Ziko, Eric Granger, Marco Pedersoli, Pablo Piantanida, and
   Ismail Ben Ayed. A unifying mutual information view of metric learning: cross-entropy vs. pairwise losses.
   In *European Conference on Computer Vision*, pages 548–564. Springer, 2020. 5
- [8] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind
  Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss,
  Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens
  Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack
  Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language
  models are few-shot learners. *CoRR*, abs/2005.14165, 2020. URL https://arxiv.org/abs/2005.
  14165. 7
- [9] Devendra Singh Chaplot, Kanthashree Mysore Sathyendra, Rama Kumar Pasumarthi, Dheeraj Rajagopal,
   and Ruslan Salakhutdinov. Gated-attention architectures for task-oriented language grounding, 2018. 3
- [10] Howard Chen, Alane Suhr, Dipendra Misra, Noah Snavely, and Yoav Artzi. Touchdown: Natural language
   navigation and spatial reasoning in visual street environments, 2020. 3, 6
- [11] Lili Chen, Kevin Lu, Aravind Rajeswaran, Kimin Lee, Aditya Grover, Michael Laskin, Pieter Abbeel,
   Aravind Srinivas, and Igor Mordatch. Decision transformer: Reinforcement learning via sequence modeling,
   2021. 6
- [12] Valerie Chen, Abhinav Gupta, and Kenneth Marino. Ask your humans: Using human instructions to
   improve generalization in reinforcement learning, 2021. 6
- [13] Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan Saharia,
   Thien Huu Nguyen, and Yoshua Bengio. Babyai: A platform to study the sample efficiency of grounded
   language learning, 2019. 2, 6, 7, 18
- Yan Duan, Marcin Andrychowicz, Bradly Stadie, Openai Jonathan Ho, Jonas Schneider, Ilya Sutskever,
   Pieter Abbeel, and Wojciech Zaremba. One-Shot imitation learning. In I Guyon, U V Luxburg, S Bengio,
   H Wallach, R Fergus, S Vishwanathan, and R Garnett, editors, *Advances in Neural Information Processing Systems 30*, pages 1087–1098. Curran Associates, Inc., 2017. 2
- [15] Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning
   skills without a reward function, 2018. 3
- [16] Chelsea Finn, Tianhe Yu, Tianhao Zhang, Pieter Abbeel, and Sergey Levine. One-shot visual imitation
   learning via meta-learning, 2017. 2
- [17] Daniel Fried, Ronghang Hu, Volkan Cirik, Anna Rohrbach, Jacob Andreas, Louis-Philippe Morency,
   Taylor Berg-Kirkpatrick, Kate Saenko, Dan Klein, and Trevor Darrell. Speaker-follower models for
   vision-and-language navigation, 2018. 3
- [18] Divyansh Garg, Shuvam Chakraborty, Chris Cundy, Jiaming Song, and Stefano Ermon. Iq-learn: Inverse
   soft-q learning for imitation. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.
   URL https://openreview.net/forum?id=Aeo-xqtb5p. 2
- [19] R.M. Gray and D.L. Neuhoff. Quantization. *IEEE Transactions on Information Theory*, 44(6):2325–2383,
   1998. doi: 10.1109/18.720541. 5

- [20] Karl Moritz Hermann, Felix Hill, Simon Green, Fumin Wang, Ryan Faulkner, Hubert Soyer, David
   Szepesvari, Wojciech Marian Czarnecki, Max Jaderberg, Denis Teplyashin, Marcus Wainwright, Chris
   Apps, Demis Hassabis, and Phil Blunsom. Grounded language learning in a simulated 3d world, 2017. 3
- 446 [21] Jonathan Ho and Stefano Ermon. Generative adversarial imitation learning, 2016. 2
- Hengyuan Hu, Denis Yarats, Qucheng Gong, Yuandong Tian, and Mike Lewis. Hierarchical decision
   making by generating and following natural language instructions, 2019. 3
- Yiding Jiang, Shixiang Gu, Kevin Murphy, and Chelsea Finn. Language as an abstraction for hierarchical deep reinforcement learning, 2019. 3
- [24] Dmitry Kalashnikov, Jacob Varley, Yevgen Chebotar, Benjamin Swanson, Rico Jonschkowski, Chelsea
   Finn, Sergey Levine, and Karol Hausman. Mt-opt: Continuous multi-task robotic reinforcement learning
   at scale, 2021. 2
- [25] Alexander C. Li, Carlos Florensa, Ignasi Clavera, and Pieter Abbeel. Sub-policy adaptation for hierarchical
   reinforcement learning, 2020. 3
- [26] Yunzhu Li, Jiaming Song, and Stefano Ermon. Infogail: Interpretable imitation learning from visual demonstrations, 2017. 3
- [27] Jelena Luketina, Nantas Nardelli, Gregory Farquhar, Jakob Foerster, Jacob Andreas, Edward Grefenstette,
   Shimon Whiteson, and Tim Rocktäschel. A survey of reinforcement learning informed by natural language,
   2019. 3
- [28] Corey Lynch and Pierre Sermanet. Language conditioned imitation learning over unstructured data.
   *Robotics: Science and Systems*, 2021. URL https://arxiv.org/abs/2005.07648.3
- 463 [29] Corey Lynch, Mohi Khansari, Ted Xiao, Vikash Kumar, Jonathan Tompson, Sergey Levine, and Pierre
   464 Sermanet. Learning latent plans from play. *Conference on Robot Learning (CoRL)*, 2019. URL https:
   465 //arxiv.org/abs/1903.01973. 3
- 466 [30] Corey Lynch, Mohi Khansari, Ted Xiao, Vikash Kumar, Jonathan Tompson, Sergey Levine, and Pierre
   467 Sermanet. Learning latent plans from play. *Conference on Robot Learning (CoRL)*, 2019. URL https:
   468 //arxiv.org/abs/1903.01973. 7
- (31) Chih-Yao Ma, Jiasen Lu, Zuxuan Wu, Ghassan AlRegib, Zsolt Kira, Richard Socher, and Caiming Xiong.
   Self-monitoring navigation agent via auxiliary progress estimation, 2019. 3
- [32] Matt MacMahon, Brian Stankiewicz, and Benjamin Kuipers. Walk the talk: Connecting language,
   knowledge, and action in route instructions. In *Proceedings of the 21st National Conference on Artificial Intelligence Volume 2*, AAAI'06, page 1475–1482. AAAI Press, 2006. ISBN 9781577352815. 3
- 474 [33] Oier Mees, Lukas Hermann, Erick Rosete-Beas, and Wolfram Burgard. Calvin a benchmark for language 475 conditioned policy learning for long-horizon robot manipulation tasks. *arXiv preprint arXiv:2112.03227*,
   476 2021. 6
- [34] Dipendra Misra, John Langford, and Yoav Artzi. Mapping instructions and visual observations to actions
   with reinforcement learning, 2017. 3
- [35] Dipendra Misra, Andrew Bennett, Valts Blukis, Eyvind Niklasson, Max Shatkhin, and Yoav Artzi. Mapping
   instructions to actions in 3d environments with visual goal prediction, 2019. 3
- [36] Ofir Nachum, Shixiang Gu, Honglak Lee, and Sergey Levine. Data-efficient hierarchical reinforcement
   learning, 2018. 3
- [37] Ashvin Nair, Vitchyr Pong, Murtaza Dalal, Shikhar Bahl, Steven Lin, and Sergey Levine. Visual reinforce ment learning with imagined goals, 2018. 2
- [38] Suraj Nair, Eric Mitchell, Kevin Chen, Brian Ichter, Silvio Savarese, and Chelsea Finn. Learning language conditioned robot behavior from offline data and crowd-sourced annotation, 2021. 2, 3, 6, 7, 19
- [39] Junhyuk Oh, Satinder Singh, Honglak Lee, and Pushmeet Kohli. Zero-shot task generalization with
   multi-task deep reinforcement learning, 2017. 3
- [40] Chris Paxton, Yonatan Bisk, Jesse Thomason, Arunkumar Byravan, and Dieter Fox. Prospection: Inter pretable plans from language by predicting the future, 2019. 3

- [41] Dean A Pomerleau. Efficient training of artificial neural networks for autonomous navigation. *Neural computation*, 3(1):88–97, 1991. ISSN 0899-7667. 2
- [42] Stéphane Ross and Drew Bagnell. Efficient reductions for imitation learning. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 661–668, 2010.
- [43] Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and structured
   prediction to no-regret online learning. In *Proceedings of the fourteenth international conference on artificial intelligence and statistics*, pages 627–635, 2011. 2
- [44] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert:
   smaller, faster, cheaper and lighter, 2020. 6
- [45] Archit Sharma, Shixiang Gu, Sergey Levine, Vikash Kumar, and Karol Hausman. Dynamics-aware
   unsupervised discovery of skills, 2020. 3
- [46] Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke
   Zettlemoyer, and Dieter Fox. Alfred: A benchmark for interpreting grounded instructions for everyday
   tasks, 2020. 3, 6
- [47] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. Cliport: What and where pathways for robotic manipula tion, 2021. 2, 3
- [48] Mohit Shridhar, Xingdi Yuan, Marc-Alexandre Côté, Yonatan Bisk, Adam Trischler, and Matthew
   Hausknecht. Alfworld: Aligning text and embodied environments for interactive learning, 2021. 6
- [49] Simon Stepputtis, Joseph Campbell, Mariano Phielipp, Stefan Lee, Chitta Baral, and Heni Ben Amor.
   Language-conditioned imitation learning for robot manipulation tasks, 2020. 3
- [50] Alane Suhr and Yoav Artzi. Situated mapping of sequential instructions to actions with single-step reward
   observation, 2018. 3
- [51] Richard Sutton, Doina Precup, and Satinder Singh. Between mdps and semi-mdps: A framework for
   temporal abstraction in reinforcement learning. 1999. 3, 4
- [52] Hao Tan, Licheng Yu, and Mohit Bansal. Learning to navigate unseen environments: Back translation with
   environmental dropout, 2019. 3
- 517 [53] Aaron van den Oord, Oriol Vinyals, and Koray Kavukcuoglu. Neural discrete representation learning,
   518 2018. 5
- [54] Jesse Vig. A multiscale visualization of attention in the transformer model. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics: System Demonstrations*, pages 37–42,
   Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-3007. URL
   https://www.aclweb.org/anthology/P19-3007. 2
- 523 [55] Sida I. Wang, Percy Liang, and Christopher D. Manning. Learning language games through interaction,
   2016. 3
- [56] Danfei Xu, Suraj Nair, Yuke Zhu, Julian Gao, Animesh Garg, Li Fei-Fei, and Silvio Savarese. Neural task
   programming: Learning to generalize across hierarchical tasks. In 2018 IEEE International Conference on Robotics and Automation (ICRA), pages 3795–3802. IEEE, 2018. 2
- [57] Tianhe Yu, Chelsea Finn, Annie Xie, Sudeep Dasari, Tianhao Zhang, Pieter Abbeel, and Sergey Levine.
   One-shot imitation from observing humans via domain-adaptive meta-learning, 2018. 2
- [58] Jesse Zhang, Haonan Yu, and Wei Xu. Hierarchical reinforcement learning by discovering intrinsic options,
   2021. 3

# 532 Checklist

533 1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contribu tions and scope? [Yes]
- (b) Did you describe the limitations of your work? [Yes] See section 5
- (c) Did you discuss any potential negative societal impacts of your work? [Yes] See appendix section
   A

539	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]	
540	2. If you are including theoretical results	
541	(a) Did you state the full set of assumptions of all theoretical results? $[N/A]$	
542	(b) Did you include complete proofs of all theoretical results? $[N/A]$	
543	3. If you ran experiments	
544 545 546	(a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] We will provide a zip file of code with the supplementary material	
547 548	<ul> <li>(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)?</li> <li>[Yes] See appendix section D</li> </ul>	
549 550	(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] Tables 10 and 2 are both over 3 seeds.	
551 552 553	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] This information is in the Training details section of section 5	
554	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets	
555	(a) If your work uses existing assets, did you cite the creators? [Yes]	
556	(b) Did you mention the license of the assets? [N/A]	
557	(c) Did you include any new assets either in the supplemental material or as a URL? $[N/A]$	
558 559	(d) Did you discuss whether and how consent was obtained from people whose data you're us- ing/curating? [N/A] The datasets we use are not collected with the help of humans	
560 561 562	(e) Did you discuss whether the data you are using/curating contains personally identifiable in- formation or offensive content? [N/A] The datasets we use are not collected with the help of humans	
563	5. If you used crowdsourcing or conducted research with human subjects	
564 565	<ul> <li>(a) Did you include the full text of instructions given to participants and screenshots, if applicable?</li> <li>[N/A]</li> </ul>	
566 567	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]	
568 569	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]	

# 570 Appendix

# 571 A Broader Societal Impact

We introduce a new method for language-conditioned imitation learning to perform complex navigation and manipulation tasks. Our intention is for this algorithm to be used in a real-world setting where humans can provide natural language instructions to robots that can carry them out. However, we must ensure that the language commands that we provide to these agents must be well aligned with the objectives of humans and must ensure that we are aware of what actions the agent could take given said command.

# 577 **B** More visualizations

# 578 B.1 Generating heat maps and word clouds

To generate heat maps and word clouds, for each evaluation instruction, we run the model and record all the 579 skill codes used in the trajectory generated. We now tokenize the instruction and for each skill code used in the 580 trajectory, record *all* the tokens from the language instruction. Once we have this mapping from skills to tokens, 581 we can plot heat maps and word clouds. This is the best we can do since we don't know exactly which tokens in 582 the instruction correspond to the skills chosen. Therefore, these plots can tend to be a little noisy but we still see 583 some clear patterns. Especially in BabyAI, since the vocabulary is small, we see that several skills correspond 584 to the same tokens because many instructions contain the same tokens. But in LOReL because each task uses 585 almost completely different words, we can see a very sparse heat map with clear correlations. 586

For reader viewability and aiding the interpretability on the LISA skills, we show below the unnormalized heatmaps showing the skill-word correlations, the column normalized heatmaps showing word frequencies for each skill as well as the row normalized showing the skill frequencies for each word.

### 590 B.2 BabyAI



Figure 7: Skill Heat map on BabyAI BossLevel



Figure 8: Word Freq. for each skill on BabyAI BossLevel (column normalized)



Figure 9: Skill Freq. for each word on BabyAI BossLevel (row normalized)

# 591 B.3 WordClouds

<sup>592</sup> Due to the small vocabulary in BabyAI environment, its hard to generate clean word clouds, nevertheless, we

<sup>593</sup> hope they help with interpreting LISA skills.



Figure 10: Word Cloud on BabyAI BossLevel for z = 1



Figure 11: Word Cloud on BabyAI BossLevel for z = 13



Figure 12: Word Cloud on BabyAI BossLevel for z = 37

# 594 B.4 LOReL Sawyer



Figure 13: Skill Heat map on LOReL Sawyer



Figure 14: Word Freq. for each skill on LOReL Sawyer (column normalized)

##cl ##ock ##fi ##ti ##ti ##tice ##wise ##wise ##wise ##wise ##zzle and away block bring cabinet camera clockwise																				- 0.8
dose closed doser container cup cup oupboard dark drawer drawer from from tront																				- 0.6
handra left left lighter mog open pull pull ngrit rep right strift																				- 0.4
slide spin tap towards translate translate un valve white	ò	1	2	3	i i	5	6	 8	9	10	'n	12	13	14	15	16	17	18	19	- 0.0

Figure 15: Skill Freq. for each word on LOReL Sawyer (row normalized)





Figure 16: Behavior and language corresponding to skill code 4: "turn faucet left"



Figure 17: Behavior and language corresponding to skill code 15: "move white mug right"



Figure 18: LOReL Composition task: "close drawer and turn faucet left"

# 596 C Datasets

# 597 C.1 BabyAI Dataset



Figure 19: BabyAI BossLevel

The BabyAI dataset [13] contains 19 levels of increasing difficulty where each level is set in a grid world where

an agent has a partially observed state of a square of side 7 around it. The agent must learn to perform various

- tasks of arbitrary difficulty such as moving objects between rooms, opening doors or closing them etc. all with a partially observed state and a language instruction.
- Each level comes with 1 million language conditioned trajectories, and we use a small subset of these for our

training. We evaluate our model on the environment provided with each level that generates a new language instruction and grid randomly.

We have provided details about the levels we evaluated on below. More details can be found in the original paper.

# 606 C.1.1 GoToSeq

- 607 Sequencing of go-to-object commands.
- Example command: "go to a box and go to the purple door, then go to the grey door"
- 609 Demo length:  $72.7 \pm 52.2$

### 610 C.1.2 SynthSeq

- Example command: "put a purple key next to the yellow key and put a purple ball next to the red box on your
- 612 left after you put a blue key behind you next to a grey door'
- 613 Demo length:  $81.8 \pm 61.3$

#### 614 C.1.3 BossLevel

- Example command: "pick up a key and pick up a purple key, then open a door and pick up the yellow ball"
- 616 Demo length:  $84.3 \pm 64.5$

### 617 C.2 LOReL Sawyer Dataset



Figure 20: LOReL Sawyer Environment

- This dataset [38] consists of *pseudo-expert* trajectories collected from a RL buffer of a a random policy and
- has been labeled with post-hoc crowd-sourced language instructions. Therefore, the trajectories complete the

620 language instruction provided but may not necessarily be optimal. The Sawyer dataset contains 50k language 621 conditioned trajectories on a simulated environment with a Sawyer robot of demo length 20.

We evaluate on the same set of instructions the original paper does for 10, which can be found in the appendix of

- the original paper. These consist of the following 6 tasks and rephrasals of these tasks where they change only the noun, only the verb, both noun and verb and rewrite the entire task (human provided). This comes to a total
- of 77 instructions for all 6 tasks combined. An example is shown below and the full list of instructions can be found in the original paper.
- 627 1. Close drawer
- 628 2. Open drawer
- 629 3. Turn faucet left
- 630 4. Turn faucet right
- 5. Move black mug right
- 632 6. Move white mug down

Table 3: LOReL Example rephrasals for the instruction "close drawer"					
Seen	Unseen Verb	Unseen Noun	Unseen Verb + Noun	Human Provided	
close drawer	shut drawer	close container	shut container	push the drawer shut	

For the composition instructions, we took these evaluation instructions from the original paper and combined

them to form 12 new composition instructions as shown below.

Table 4: LOReL Composition tasks Instructions open drawer and move black mug right pull the handle and move black mug down move white mug right move black mug down close drawer and turn faucet right close drawer and turn faucet left turn faucet left and move white mug down turn faucet right and close drawer move white mug down and turn faucet left open drawer and turn faucet counterclockwise slide the drawer closed and then shift white mug down turn faucet left and move white mug down move white mug down and move black mug right turn faucet right and open cabinet move black mug right and turn faucet right

We included the instructions *"move white mug right"* and *"move black mug down"* as composition tasks here in the hope that we may have skills corresponding to colors like black and white or directions like right and down

that can be composed to form these instructions but we did not observe such behaviour.

# 638 **D** Training details

We plan to release our code on acceptance. Here we include all hyper-parameters we used. We implement our models in PyTorch. Our original flat baseline implementation borrows from Decision Transformer codebase which uses GPT2 to learn sequential behavior. However, we decided to start from scratch in order to implement LISA to make our code modular and easily support hierarchy. We use 1 layer Transformer networks for both the skill predictor and the policy network in our experiments for the main paper. We tried using large number of layers but found them to be too computationally expensive without significant performance improvements. In BabyAI and LOReL results we train all models for three seeds.

For BossLevel environment we use 50 skill codes, for other environments we used the settings detailed in the table below:

### 648 D.1 LISA

Hyperparameter	BabyAI	LORL
Transformer Layers	1	1
Transformer Embedding Dim	128	128
Transformer Heads	4	4
Skill Code Dim	16	16
Number of Skills	20	20
Dropout	0.1	0.1
Batch Size	128	128
Policy Learning Rate	1e - 4	1e - 4
Skill Predictor Learning Rate	1e - 5	1e - 5
Language Model Learning Rate	1e - 6	1e - 6
VQ Loss Weight	0.25	0.25
Horizon	10	10
VQ EMA Update	0.99	0.99
Optimizer	Adam	Adam
	-	

#### Table 5: LISA Hyperparameters

# 649 D.2 Baselines

# Table 6: Flat Baseline Hyperparameters

Hyperparameter	BabyAI	LORL
Transformer Layers	2	2
Transformer Embedding Dim	128	128
Transformer Heads	4	4
Dropout	0.1	0.1
Batch Size	128	128
Policy Learning Rate	1e - 4	1e - 4
Language Model Learning Rate	1e - 6	1e - 6
Optimizer	Adam	Adam

For the original baseline for BabyAI, we used the code from the original repository. For the LOReL baseline, we 650 used the numbers from the paper for LOReL Images. For LOReL States BC baseline, we implemented it based 651 on the appendix section of the paper. We ran the LOReL planner from the original repository for the composition 652 instructions. 653

#### **D.3** Ablations 654

As mentioned in the paper, all our ablations were performed on BabyAI BossLevel with 1k trajectories over a 655 single seed for the sake of time. Unless otherwise specified, we use the following settings. We use a 1-layer, 656 4-head transformer for both the policy and the skill predictor. We use 50 options and a horizon of 10. We use a 657 batch size of 128 and train for 2500 iterations. We use a learning rate of 1e-6 for the language model and 1e-4 658 for the other parameters of the model. We use 2500 warm-up steps for the DT policy. Training was done on 659 660 Titan RTX GPUs.

#### Ε **Detailed LOReL Sawyer results** 661

We provide details results on the LOReL evaluation instructions below for LISA and the flat baseline in the same 662 format as the original paper. The results are averaged over 10 runs. The time horizon used was 20 steps. 663

Table 7: Task-wise success rates (in %) on LOReL Sawyer.

Task	Flat	LISA
close drawer	10	100
open drawer	60	20
turn faucet left	0	0
turn faucet right	0	30
move black mug right	20	60
move white mug down	0	30

Table 8: Rephrasal-wise success rates (in %) on LOReL Sawyer.

Rephrasal Type	Flat	LISA
seen	15	40
unseen noun	13.33	33.33
unseen verb	28.33	30
unseen noun+verb	6.7	20
human	26.98	27.35

#### More experiments F 664

#### F.1 Ablation Studies 665

For the sake of time, all our ablations were performed with a 1-layer, 4-head transformer for the skill predictor 666 and for the policy. All our ablations are on the BabyAI-BossLevel environment with 1k expert trajectories. 667

Our first experiment varies the horizon of the skills. The table below shows the results on BabyAI BossLevel 668

for 4 different values of the horizon. We see that the method is fairly robust to the different choices of horizon 669

unless we choose a very small horizon. For this case, we notice that a horizon of 5 performs best, but this could 670 vary with different tasks. 671

Table 9: Ablation on horizon. We fixed the number of options to be 50 for these experiments

Horizon	1	5	10	50
Success Rate (in %)	32	52	47	47

672 We also tried varying the number of skills the skill-predictor can choose from and found that this hyperparameter

is fairly robust as well unless we choose an extremely high or low value. We suspect using more skills worsens 673

performance because it leads to a harder optimization problem and the options don't clearly correspond to 674 specific language skills.

675

Table 10: Ablation on number of options. We fixed the horizon to be 10 for these experiments

Number of Options	10	20	50	100
Success Rate (in %)	47	47	47	43

#### F.2 Can we leverage the interpretability of the skills produced by LISA for manual planning? 676

Since our skills are so distinct and interpretable, its tempting to try and manually plan over the skills based on 677 the language tokens they encode. In the LOReL environment, using the same composition tasks as section above, 678

we observe the word clouds of options and simply plan by running the fixed option code corresponding to a 679 task for a certain horizon and then switch to the next option corresponding to the next task. This means we are 680 using a manual (human) skill predictor as opposed to our trained skill predictor. While this doesn't work as 681 well because skills are a function of both language and trajectory and we can only interpret the language part as 682 humans, it still shows how interpretable our skills are as humans can simply observe the language tokens and 683 plan over them to complete tasks. We show a successful example and a failure case below for the instruction 684 "close drawer and turn faucet right" in figure 21. We first observe that the two skills we want to compose are 685 Z = 14 and Z = 2 as shown by the word clouds. We then run skill 14 for 20 steps and skill 2 for 20 steps. In 686 the failure case, the agent closes the drawer but then pulls it open again when trying to turn the faucet to the 687 688 right.



Z = 14





Successful Manual Planning



Unsuccessful Manual Planning

Figure 21: We show a successful manual planning and unsuccessful manual planning example for the instruction "close the drawer and turn the faucet right"

#### 689 F.3 Do the skills learned transfer effectively to similar tasks?

We want to ask the question whether we can use the skills learned on one task as a initialization point for a similar task or even freeze the learned skills for the new task. To this end, we set up experiments where we

trained LISA on the BabyAI GoTo task with 1k trajectories and tried to transfer the learned options to the GoToSeq task with 1k trajectories. Similarly we trained LISA on GoToSeq with 1k trajectories and tried to

GoToSeq task with 1k trajectories. Similarly we trained LISA on GoToSeq with 1k trajectories transfer to BossLevel with 1k trajectories. The results are shown in figures 22 and 23 respectively.

As we can see from the GoToSeq experiment in figure 22, there is no major difference between the three methods.

<sup>696</sup> We notice that we can achieve good performance even by holding the learned option codes from GoTo frozen.

<sup>697</sup> This is because the skills in GoTo and GoToSeq are very similar except that GoToSeq composes these skills as

tasks. We also notice that finetuning doesn't make a big difference – once again probably because the skills are

699 similar for both environments.



Figure 22: Transferring skills on the BabyAI GoToSeq environment



Figure 23: Transferring skills on the BabyAI BossLevel environment

In the BossLevel case in 23, however, we do notice that the frozen skills perform slightly worse than the other 700 two methods. This is because the BossLevel contains more skills than those from GoToSeq. We also notice 701 that the performance of finetuning and starting from scratch is nearly the same. This could be because the 702 meta-controller needs to adapt to use the new skills in the BossLevel environment anyway and there is no benefit 703 from loading learned options here.

#### 705 F.4 State-based skill-predictor

We have already spoken in section 3.2 about the fact that our method using two transformers doesn't necessarily 706 mean its more compute-heavy than the non-hierarchical counterpart. But to test whether we really need two 707 transformers, we perform an ablation study that replaces the skill predictor to be just a state-based selector 708 MLP as opposed to a trajectory-based transformer. Our results show that the performance is only slightly worse 709 when using a state-based skill predictor in this case, but once again this may not be the case in more complex 710 environments. However, we notice that the skills collapsed in this case and the model tends to use fewer skills 711 than normal as shown in figure 24. This is expected because the skill predictor is now predicting skills with 712 much less information. 713

Table 11: Comparing state-based MLP skill predictor vs trajectory-based transformer skill predictor. We fixed the number of options to be 50 and horizon as 10 for these experiments

Skill Predictor Architecture	Success Rate
State-based MLP	46%
Trajectory-based Transformer	47%



Figure 24: State-based skill predictor heat map shows that the model tends to use fewer options compared to figure 7

#### 714 F.5 Continuous skill codes

We also compare to the non-quantized counterpart where we learn skills from a continuous distribution as opposed to a categorical distribution. We expect this to perform better because the skill predictor has access to a larger number of skill codes to choose from and this is what we observe in table 12. However, this comes at the price of interpretability and its harder to interpret and choose continuous skill codes than discrete codes. We also observe that on the LOReL with states environment, using discrete codes performs better than using continuous codes (table 13). This could be because learning a multi-modal policy with discrete skills is an easier optimization problem than learning one with continuous skills (see the end of section 4.3).

Table 12: Ablation on Quantization on BabyAI BossLevel. We fix number of options to 50 and horizon to 10

Skill codes	Success Rate (in %)
Continuous	51
Discrete	47

Table 13: Ablation on Quantization on LOReL with states on seen tasks. We fix number of options to 20 and horizon to 10

Skill codes	Success Rate (in %)
Continuous	60.0
Discrete	66.7