# Learning Compositional Tasks from Language Instructions

Lajanugen Logeswaran<sup>1</sup>, Wilka Carvalho<sup>2</sup>, Honglak Lee<sup>1,2</sup>

<sup>1</sup>LG AI Research <sup>2</sup>University of Michigan, Ann Arbor llajan@lgresearch.ai

#### Abstract

The ability to combine learned knowledge and skills to solve novel tasks is a key aspect of generalization in humans that allows us to understand and perform tasks described by novel language utterances. While progress has been made in supervised learning settings, no work has yet studied compositional generalization of a reinforcement learning agent following natural language instructions in an embodied environment. We develop a set of tasks in a photo-realistic simulated kitchen environment that allow us to study the degree to which a behavioral policy captures the systematicity in language by studying its zero-shot generalization performance on held out natural language instructions. We show that our agent which leverages a novel additive action-value decomposition in tandem with attention-based subgoal prediction is able to exploit composition in text instructions to generalize to unseen tasks.

#### Introduction

Human language is characterized by systematic compositionality: one can combine known components – such as words or phrases – to produce novel linguistic combinations (Chomsky 2009). This is a key aspect of generalization in humans and enables us to understand and perform tasks specified by novel language utterances over familiar words or phrases. If you know what a "laptop" and a "fridge" are, you can easily understand how to perform the task "place the laptop in the fridge" even if you have never placed a laptop in a fridge.

Prior work studying the linguistic "systematicity" of neural networks have focused on sequence mapping tasks in a supervised learning setting (Lake and Baroni 2018; Lake 2019; Andreas 2019). In this work, we are interested in compositional generalization of a reinforcement learning agent following natural language instructions in an embodied environment. In particular, we explore the hypothesis that a language-conditioned reinforcement learning agent with a compositional inductive bias in its behavioral policy will exhibit systematic generalization to unobserved natural language instructions.

There has been a flurry of recent work on embodied learning tasks such as question answering (Gupta et al. 2017),



Figure 1: Zero-shot generalization to an unseen task of slicing an apple. The test task is composed of known primitive subtasks – *picking up a knife* and *slicing the apple* – each of which were encountered in training tasks. Our agent learns to decompose a natural language task description into subtasks using attention and executes them using low-level actions.

vision-language navigation (Anderson et al. 2018) and object interaction (Shridhar et al. 2020; Carvalho et al. 2020; Singh et al. 2020; Corona et al. 2020; Blukis et al. 2022; Min et al. 2021) in embodied settings. In particular, the AL-FRED task (Shridhar et al. 2020) studies agents that exploit detailed natural language instructions to generalize to novel instructions in novel environments at test time. Such existing benchmarks offer limited flexibility to study systematic generalization since (i) the benchmarks were not built for this purpose and it is unclear to what extent systematic generalization skills are required to solve the tasks and (ii) the tasks demand challenging reasoning skills such as visual recognition, planning over large number of time-steps and exploration in unseen layouts which makes it difficult to study compositional generalization ability in isolation.

In this work we develop a set of tasks in the AI2Thor virtual home environment (Kolve et al. 2017) which test the compositionality of embodied agents. In order to make

Copyright © 2023, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

progress in systematic generalization, we make two simplifying assumptions: we assume access to an oracle object recognizer and we study generalization in a single kitchen layout. This allows us to study the degree to which a policy captures the systematicity in language by studying its zero-shot generalization performance on held out natural language instructions.

Despite these simplifications, agents still need to understand the instruction to figure out the sequence of object interactions that need to be performed and act over many timesteps with limited guidance. In order to succesfully generalize at test time, an agent needs to learn to ground natural language instructions to temporally extended goal-oriented behaviors or "skills" in a compositional manner to perform novel tasks that are compositions of the tasks presented at train time. In addition, the agent needs to ground the actions and objects referenced in the text to actions and objects in the environment. We leverage this setting to develop and study a policy with an inductive bias for compositionality and show that this enables systematic generalization in the context of combining behavioral skills learned purely from reward without expert demonstrations.

We present an attention-based agent that learns to predict subgoals from language instructions via a learned attention mechanism. Our agent uses these subgoals with a novel policy parametrization which decomposes the actionvalue function in an additive fashion that enables estimating the action-value for novel object-interactions composed of objects and interactions experienced during training.

We show evidence that this parametrization facilitates exploiting the compositional nature of text instructions by showing systematic generalization to both unseen task descriptions and unseen tasks. We present an example in Figure 1, where the agent is able to systematically generalize the behavior "pickup up the knife" to "move the knife" and "cut the yellow apple" to "slice the apple". Thanks to the additive inductive bias afforded by our action-value parametrization, it is able to compose these behaviors to perform the novel task "move the knife near the table and slice the apple" at test time.

## **Related Work**

Compositional Generalization Compositionality is an important aspect of generalization: decomposing a learning problem into components that the agent knows how to address. An agent with a compositional understanding of its environment and its own skills can more easily combine these skills to solve tasks never observed during training. Prior work has studied compositional generalization in sequence mapping tasks. Benchmarks such as SCAN (Lake and Baroni 2018) and gSCAN (Ruis et al. 2020) study translating synthetic text descriptions to an action sequence (e.g. jump twice  $\rightarrow$  JUMP JUMP). gSCAN couples SCAN instances with entities in a grid environment and solving a task requires grounding the text and entities similar to our work. Prior approaches for these benchmarks impose compositional inductive biases in models by augmenting models with memory (Lake 2019) and data augmentation (Andreas 2019; Shaw et al. 2020). Lake (2019) proposed a memory-augmented sequence-to-sequence model and a meta-learning algorithm to enable models to recognize the compositional structure of instances. Andreas (2019) introduce a data augmentation approach where they synthesize new examples by combining phrases from known examples to provide a compositional inductive bias to models. Shaw et al. (2020) demonstrate a grammar induction method and combine it with a pre-trained T5 model (Raffel et al. 2019) for semantic parsing tasks. In this work we use attention mechanisms and introduce a novel policy parameterization to impose compositional inductive biases.

Text-based Embodied Control There is rich literature on instruction following agents (Chen and Mooney 2011; Tellex et al. 2011; Mei, Bansal, and Walter 2016; Fried, Andreas, and Klein 2017; Suhr et al. 2019). Advances in photo-realistic simulation environments such as DeepMind Lab (Beattie et al. 2016) and AI2Thor (Kolve et al. 2017) have driven recent progress in embodied agents that learn from text instructions. Chaplot et al. (2018) consider a simple navigation task where an agent has to move to an object specified by a set of attributes such as shape and color. They propose the gated attention model to generalize compositionally in the attribute space. Hill et al. (2019) consider systematic generalization in 2D and 3D environments with synthetic text instructions. Compared to these work, we consider object interaction tasks in a photo realistic simulated environment with human-authored language instructions.

ALFRED (Shridhar et al. 2020) is a recently proposed benchmark which couples tasks in the AI2Thor environment with text instructions. The dataset comes with expert demonstrations and human-written text descriptions of tasks. Prior approaches to ALFRED train policies using imitation learning (Singh et al. 2020; Corona et al. 2020). More recently, policies based on spatial semantic map building have received more interest (Blukis et al. 2022; Min et al. 2021).

In contrast to the ALFRED learning setup, we consider a simplified scenario of learning compositional skills from high-level task descriptions. We further do not assume access to expert task demonstrations. These assumptions allow us to focus on compositional generalization to zero-shot tasks, which is not the main goal of the ALFRED benchmark. However, the approach presented here can potentially be applicable to ALFRED when combined with learning from demonstrations.

**Hierarchical Reinforcement Learning** Learning to directly map percepts to low-level action sequences can be challenging. An alternative hierarchical approach is to first come up with a sequence of subtasks, which can be considered as high-level actions (Andreas, Klein, and Levine 2017; Zhu et al. 2017). Each of those subtasks can then be realized using low-level actions. Our policy has an implicit hierarchical structure where latent subgoals are represented as text embeddings using attention. Language was used as an abstraction for the high-level policy in Jiang et al. (2019a) for object rearrangement tasks based on the CLEVR engine (Johnson et al. 2017).

Finally, generalization to unseen instructions has been



Figure 2: Approach Overview. We perform attention over the text instruction to construct an embedding  $t_{sg}$  that represents the current subgoal. The text embedding subgoal  $t_{sg}$  attends to scene object embeddings to construct an object subgoal representation  $v_{sg}$ . An MLP takes  $t_{sg}$ ,  $v_{sg}$  and observation features  $e_{obs}$  as input and predicts state-action values Q(s, a). The entire model is trained end-to-end using Q-learning. See text for details.

considered in prior work such as Oh et al. (2017); Lynch and Sermanet (2020), although compositional generalization is not their main focus.

### Problem

We consider an embodied agent acting in a kitchen environment to solve basic tasks from language instructions (See Fig. 5 for an example task). At the beginning of an episode the agent receives a text instruction  $\tau$ . Our goal is to learn a policy  $\pi(a|s, \tau)$ ;  $a \in A, s \in S$  that predicts actions in order to complete tasks. The agent state s is partially observable – it receives an egocentric observation *obs* of the scene. We further assume that an oracle object recognition model provides the object ids for objects in the egocentric observation.

The action space consists of navigation and object interaction actions  $\mathcal{A} = \mathcal{A}_{nav} \cup \mathcal{A}_{int}$ . There are 8 navigation actions  $\mathcal{A}_{nav} = \{move \ forward, \ move \ back, \ move \ left, \ move \ right, \ turn \ left, \ turn \ right, \ look \ up, \ look \ down\}$ . Interaction actions  $\mathcal{A}_{int} = \mathcal{B} \times \mathcal{I}$  are specified using an interaction  $b \in \mathcal{B}$  and an object id  $o \in \mathcal{I}$  where  $\mathcal{B} = \{pickup, \ place, \ slice, \ tog$  $gle \ on, \ toggle \ off, \ turn \ on, \ turn \ off\}$  and  $\mathcal{I}$  is a pre-defined set of identifiers of objects that are available to the agent for interaction in the current observation.

The agent receives a positive reward for successfully completing a task. It also receives a small negative reward for every time-step. In addition, we also assume that every correct object interaction receives a positive reward. In addition to providing a denser learning signal, the rewards are also used to identify subgoals as described in the next section. In practice such dense rewards may be unavailable, but this is outside the scope of our study and left as future work.

#### Approach

We approach the problem by considering a task  $\tau$  to be composed of subgoals  $g_1, ..., g_n$ , where each subgoal  $g_i$  involves

navigating to a particular object and interacting with it. For example, the task *place an apple on the table* involves finding the apple and picking it up, followed by navigating to the table and putting down the apple, which can be considered to be the two subgoals for executing the task. Each object interaction required to complete the task thus corresponds to a subgoal. Since every subgoal completion receives a positive reward, the number of subgoals completed at every time-step  $N_{sg}$  is known to the agent. The subgoals themselves are not known to the agent – we use attention on the text instruction to compute a latent subgoal representation.

### **Text Subgoal Inference**

Given instruction  $\tau$  composed of the tokens  $(w_1, ..., w_n)$ , we obtain the corresponding token embeddings  $E = (e_1, ..., e_n)$  and use an RNN to encode the instruction to obtain a sequence of contextualized token representations  $H = (h_1, ..., h_n)$ . We compute a *text subgoal*  $t_{sg}$  for a given time-step by computing attention on the instruction using  $N_{sg}^e$  as query where  $N_{sg}^e$  is a vector representation of  $N_{sg}$ . This is shown in Eq. (1) (Q, K, V) are learnable parameter matrices).

$$t_{sq} = \text{Attention}(\text{query} = N_{sq}^e, \text{keys} = \text{values} = H)$$
 (1)

$$= \sum_{h \in H} \frac{\exp((Qs)^{\top}(Kh))}{\sum_{h' \in H} \exp((Qs)^{\top}(Kh'))} Vh$$
(2)

We expect the attention to focus on words in the instruction relevant to executing the current subgoal. For instance, if the agent is expected to interact with an apple, the attention module could learn to focus on the word 'apple'.

#### Cross-modal Reasoning

Given the text subgoal  $t_{sg}$ , we use an attention mechanism to reason about objects in the scene within some distance from

Task type	Task descriptions
pick up pot	Go to the stove and pick up the pot. Pick up the pot on the bottom right burner on the stove.
place spoon in pan	get spoon on counter near salt shaker and put it away in pan near stove. Pick up the spoon from the table near the salt shaker and move it to the pan on the counter by the sink.
place knife in plate	Pick up a knife from the left side of the kitchen, near the trashcan, and move it to the plate. get knife near sink and put on opposite side of stove on dark plate
slice bread with knife	Pick the knife and slice the bread. Take the knife with the yellow handle from the counter by the sink and use it to cut horizontal slices out of the loaf of bread on the white table.
slice lettuce with butterknife	Using the silver butter knife you have in your hand, cut up the lettuce into thin slices. Pick up the butter knife on the counter and cut the lettuce on the table.

Table 1: Example task types and corresponding task descriptions. Note that the task descriptions are used for training and testing agents. The task types are not known to the agents.

the agent. This helps the agent understand if objects of interest relevant to the subgoal are present nearby. Let the set of nearby scene objects be  $O = \{o^1, ..., o^n\}$ , where the  $o^i \in \mathcal{I}$ are object ids provided by an oracle. The  $o^i$ 's can thus be treated as indexes into an embedding table that produces object embeddings  $O_e = \{o_e^1, ..., o_e^n\}$ . The cross-modal attention is given by Eq. (4) where the text subgoal attends to the scene object embeddings (Q', K', V') are learnable parameter matrices). We augment the scene objects embeddings  $O_e$ with an additional learned embedding  $o_e^{no-obj}$  which is expected to absorb any probability mass not assigned to scene objects  $O'_e = O_e \cup o_e^{no-obj}$ . The attention produces an *object* subgoal embedding  $v_{sq}$ .

$$v_{sg} = \text{Attention}(\text{query} = t_{sg}, \text{keys} = \text{values} = O'_e)$$
 (3)

$$= \sum_{o_e \in O'_e} \frac{\exp((Q't_{sg})^\top (K'o_e))}{\sum\limits_{o'_e \in O'_e} \exp((Q't_{sg})^\top (K'o'_e))} V'o_e \tag{4}$$

# **Policy Learning**

We use a deep Q-learning algorithm to train a policy (Mnih et al. 2013), where a neural network is trained to approximate the state-action value function Q(s, a). Given the current observation, text subgoal and object subgoal, the state-action value for a navigation action  $a \in \mathcal{A}_{nav}$  is given by Eq. (5), where  $f_{nav}$  is an MLP (multi-layer perceptron) and  $e_{obs} = f_{CNN}(obs)$  is a feature vector of the observation image computed using a CNN encoder.

$$Q_{\text{nav}}(s,a) = f_{\text{nav}}(a|e_{\text{obs}}, t_{sg}, v_{sg})$$
(5)

The state-action values for interaction actions  $a = (b, o) \in \mathcal{B} \times \mathcal{I}$  can be analogously modeled as in Eq. (6). We found it helpful to decompose the state-action value in an additive fashion over an action score  $f_{int}^a$  and an object score  $f_{int}^o$  as in Eq. (7). Intuitively,  $f_{int}^a$  learns to model action preferences, whereas  $f_{int}^o$  learns to ground text goals to physical objects. In addition to sharing parameters across actions and objects, this decomposition allows us to model state-action

values of object interactions not experienced during training, as long as the specific interaction and the object were encountered. Unless specified otherwise we use the decomposed value function  $Q_{\rm int}^{\rm add}$  in our experiments.

$$Q_{\text{int}}^{\text{full}}(s,a) = f_{\text{int}}(a|e_{\text{obs}}, t_{sg}, v_{sg}) \tag{6}$$

$$Q_{\rm int}^{\rm add}(s,a) = f_{\rm int}^{a}(b|e_{\rm obs}, t_{sg}, v_{sg}) + f_{\rm int}^{o}(o|t_{sg})$$
(7)

where  $a = (b, o) \in \mathcal{B} \times \mathcal{I}$ .

In summary, the state-action value function is modeled as in Eq. (8).

$$Q(s,a) = \begin{cases} Q_{\text{nav}}(s,a); & a \in \mathcal{A}_{\text{nav}} \\ Q_{\text{int}}^{\text{add}}(s,a); & a \in \mathcal{A}_{\text{int}} \end{cases}$$
(8)

The overall model (see Fig. 2 for an illustration) including parameters of the subgoal inference (Eq. 1) and crossmodal reasoning (Eq. 4) components, as well as the MLPs in Eqs. (5) and (7) are trained end-to-end using a double-DQN algorithm (Van Hasselt, Guez, and Silver 2016). Once the model has been trained we construct a greedy policy by choosing actions with the highest state-action values for inference.

### Experiments

#### **Tasks and Dataset**

We use the AI2Thor (Kolve et al. 2017) environment as a testbed for our experiments. While there exist prior benchmarks that couple language instructions with embodied environments such as ALFRED (Shridhar et al. 2020), they were not designed to study compositional generalization. We thus construct a new task setup that allows us to flexibly vary tasks and object arguments. We consider the following task types in our experiments,

- *pickup x*: Find and pick up object x
- *place x in y*: Find and pick up object *x*, followed by navigating towards *y* and placing it.
- *slice x with y*: Secure cutting instrument *y*, find object *x* and perform the slice action on it.

We use Amazon Mechanical Turk to collect natural language descriptions of tasks for training and evaluation. A turker is shown key observation frames during the execution of a particular task and is asked to describe in a sentence how they would describe the task to a robot. Turkers were instructed to do their best to correctly identify task relevant objects. But often descriptions from the turkers incorrectly identify objects such as identifying a potato as an avocado. Such descriptions were manually fixed so that the correct object identities are mentioned in the instructions. We collected 5 natural language descriptions each for 35 task types that include *pickup*, *place* and *slice* tasks. The descriptions consist of 170 unique tokens and have an average length of 12 tokens. Table 1 shows example descriptions collected for some task types.

**Short Horizon Tasks** The pickup tasks are used for evaluating multi-task and zero-shot generalization with seen and unseen descriptions of task types. We use 10 *pickup* task types - *pickup X* where  $X \in \{apple, bread, tomato, potato, lettuce, spoon, bread, butter knife, plate, pot\}$ . These tasks are used for evaluating generalization to seen and unseen descriptions of known *short-horizon* task types. They are also used in generalization to longer horizon tasks as described later in this section.

**Longer Horizon Tasks** The *place* and *slice* tasks are used for evaluating generalization to *longer-horizon* tasks. Table 2 shows tasks used for training and evaluation. In addition to multitask generalization, we use these tasks to study zero-shot compositional generalization to unseen task descriptions. The unseen descriptions can correspond to task types that were encountered during training (i.e., seen task types), similar to the *pickup* tasks. A more challenging generalization scenario is to generalize to text descriptions of task types not encountered during training (i.e., unseen task types).

We consider two types of tasks in the latter scenario. The obj-obj setting examines the ability of the agent to generalize to tasks composed of unseen combinations of objects. For instance, in the test task type *place potato in plate*, the relevant objects *potato*, *plate* were encountered during training in task types such as *place potato in pan* and *place apple in plate*.

The task-obj setting is a harder generalization problem where the agent is expected to generalize to unseen combinations of tasks and objects. For the test task type *slice lettuce with knife*, the object *lettuce* was never observed in the context of a *slice* task during training. However, the agent has access to *pickup* tasks and is expected to learn to interact with lettuce by using the *pickup lettuce* task. This can be challenging because the agent was only taught how to pick up lettuce, and did not learn to associate lettuce with slice tasks.

The seen task types in Table 2 were designed such that each object argument appears in multiple task types. Furthermore, when choosing object arguments for a given task type, we prioritized objects that appear in as many tasks as possible. For instance, in the pickup and place tasks setup, the objects were plate, pan, pot, spoon, etc. where each ob-

#### Seen task types

place apple in plate place butterknife in plate place spoon in plate place butterknife in pan place potato in pan place spoon in pan place apple in pot place butterknife in pot place potato in pot slice apple with knife slice tomato with knife slice bread with knife slice apple with butterknife slice potato with butterknife slice bread with butterknife

#### Unseen task types (obj-obj setting)

place potato in plate	slice potato with knife
place apple in pan	slice tomato with butterknife
place spoon in pot	

Unseen task types (task-obj setting)

place knife in plate	slice lettuce with knife
place knife in pan	slice lettuce with butterknife
place knife in pot	

Table 2: Task types used for training and testing on place and slice tasks. *Seen* and *Unseen* correspond to task types whose descriptions were respectively seen and not seen during training. The obj-obj setting considers test task types composed of unseen combinations of objects. The task-obj setting considers generalization to unseen combinations of tasks and objects (e.g., learning to slice lettuce when taught how to slice objects and how to pickup lettuce). Note that all pickup task types are treated as seen and omitted in this table for brevity.

ject appears in at least three of the training tasks. This ensures that there are enough occurrences of each object type for the agent to understand and ground the object type. It also helps the agent disentangle the notion of an object versus a task in a given instruction.

### Baselines

We compare the proposed approach against the following baselines.

**RNN** In this baseline we replace the attentional model with an RNN that produces an embedding of the text instruction. While this model can potentially work for unseen instructions, we examine if the encoding effectively captures the compositional information present in the instructions.

**Gated Attention** We consider the gated attention architecture from Chaplot et al. (2018). This architecture combines the instruction representation with the visual observation using a gated attention operation. The fused representation is fed to an MLP which models the state-action values.

All models and baselines are trained using the DDQN (Double deep Q-learning) algorithm (Van Hasselt, Guez, and Silver 2016).

### Hyperparameters

Word embeddings and the RNN have representation size 32. Objects are represented by embeddings of size 32 from an

	Seen tasks		Unseen tasks	
Model	Seen descriptions	Unseen descriptions	Unseen descriptions obj-obj	Unseen descriptions task-obj
RNN	0.65	0.65	0.26	0.13
Gated Attention	0.92	0.85	0.66	0.34
Ours				
(a) $Q_{\text{nav}} + Q_{\text{int}}^{\text{add}}$ (no cross modal)	0.89	0.76	0.84	0.77
(b) $Q_{\text{nav}} + Q_{\text{int}}^{\text{full}}$ (with cross modal)	0.93	0.85	0.44	0.34
(c) $Q_{\text{nav}} + Q_{\text{int}}^{\text{add}}$ (with cross modal)	0.95	0.87	0.94	0.91

Table 3: Task success rates of models under different generalization settings. Models are evaluated on seen/unseen descriptions of seen tasks and on unseen descriptions of unseen tasks. For unseen tasks, we further evaluate under unseen combinations of objects as well as unseen combinations of tasks and objects. Best numbers are boldfaced.

pick up the silver butter **knife** closest to the edge of the counter .

the **spoon** is between the spatula and the fork on the left countertop; pick it up.

move to the table, pick up the tomato.

pick up the lettuce from the table.

pick up the **potato** between the lettuce and the tomato .

Figure 3: Visualizing task attention for pickup tasks. Words in darker shades received higher attention probabilities.

embedding table. The CNN observation features have size 512 and the CNN encoder has 1.7M parameters, which constitues 90% of the overall model parameters. The MLPs in Equations (5) and (6) are single hidden layer MLPs with 256 hidden units and ReLU activation.

### Results

**Short Horizon Tasks** We first consider pickup tasks that involve a single object interaction. In these tasks the agent has to identify the object reference mentioned in the text description and then find and pick up the relevant object. We train and evaluate on 10 pickup task types. Four text descriptions of each task type are part of the training set and the remaining descriptions (i.e., 1 per task type) are part of the test set.

On the training and test descriptions, our agent trained from scratch achieves success rates of 0.9, 0.92 respectively. Identifying the correct subgoal for these tasks involves paying attention to the verbs and nouns in the task description as well as the overall context. Figure 3 visualizes the task attention in the subgoal inference component for a subset of test task descriptions, from which it is clear that the agent learns to focus on the relevant parts of the instruction.

**Longer Horizon Tasks** We now consider tasks that involve two subgoals, which includes the place and slice tasks in Table 2. Jointly learning text grounding and subgoal inference for long horizon tasks can be challenging. We thus consider a curriculum learning strategy where an agent is gradually trained on tasks of increasingly longer horizon.



Figure 4: Learning progress of agent trained from scratch and agent pre-trained on pickup tasks.

The agent is first pre-trained on the pickup tasks as described in the previous section, and then fine-tuned on descriptions of seen tasks in Table 2. Figure 4 compares the learning progress of agents trained from scratch and an agent that has been pre-trained on the pickup tasks. The pre-trained agent learns twice as fast compared to the agent trained from scratch and achieves perfect success rate on training task descriptions.

**Generalization** Table 3 shows the average task completion success rate of models under different generalization scenarios. We report performance on seen/unseen descriptions of seen task types as well as (unseen) descriptions of unseen task types from Table 2. The RNN and Gated Attention baselines are limited by the fact that the text instruction is represented using the same encoding across all time-steps, which has limited ability to capture compositional information. The inductive bias of Gated Attention enables better performance, but it has difficulty generalizing to unseen tasks. The attention based model outperforms these baselines, which indicates that the attention mechanism helps exploit compositional information in the instruction better than a fixed encoding.

In addition to better performance, the attention model has the advantage of being more interpretable. Figure 5 shows



**Instruction attention** Subgoal 1 Subgoal 2 bring the potato from the table to the plate on the right of the oven . bring the potato from the table to the plate on the right of the oven .



Instruction attention

- Subgoal 1 Subgoal 2
- pick up the butter knife on the counter , and horizontally slice the lettuce
   pick up the butter knife on the counter , and horizontally slice the lettuce

Figure 5: Agent's observation at different time-steps while performing a place task and a slice task. The attention distribution in the text goal inference component while executing each subgoal is also given below the agent observations.

the agent's actions and the attention pattern over time for example tasks. The agent learns to identify object references in the instruction and uses attention as a sub-goal representation. This mimics a hierarchical policy where a high-level controller provides a sub-goal and a low-level controller executes it (Jiang et al. 2019b). The agent further learns to ground object references in the text instruction to objects in the scene. Notably, these attention patterns and grounding are learned from the reward signal alone without any other supervision.

### Ablations

We perform ablations to study the impact of cross-modal reasoning and decomposing the value function in an additive fashion.

**Cross-modal Reasoning** We examine model performance without the cross modal reasoning component. In this case the MLPs in Equations (5) and (7) only receive the text subgoal and observation encoding as inputs and the visual subgoal  $v_{sg}$  is omitted. This is shown as  $Q_{nav} + Q_{int}^{add}$  (no cross modal) in Table 3. From rows (a) and (c) in table Table 3 it is clear that the cross-modal reasoning components helps ground text in scene objects and enables better generalization across all settings.

**Interaction Q-values** We examine the benefit of decomposing the value function approximation of interaction ac-

tions in an additive fashion in  $Q_{int}^{add}$  (Equation (7)). We compare it against  $Q_{int}^{full}$  (Equation (6)), which treats each (interaction, object) pair as a separate atomic action. Comparing rows (b), (c) in Table 3 we see that the additive decomposition is crucial for generalization to unseen tasks.

# Conclusion

In this work we proposed attention based agents that can exploit the compositional nature of language instructions to generalize to unseen tasks. The policy mimics a hierarchical process where a text embedding obtained via attention represents the subgoal to be executed and the policy network executes the low level actions. The proposed method performs strongly against baselines on a testbed we created based on a photorealistic simulated environment and provides some interpretability.

Compared to existing benchmarks such as ALFRED we made simplifying assumptions such as oracle visual recognition, relatively short horizon tasks and generalization within single kitchen layout which allows us to focus on compositional generalization in embodied settings. However, the ideas presented here can potentially be combined with curriculum learning and learning from human demonstrations to perform complex tasks that require planning over hundreds of time-steps such as in the ALFRED setting, and we leave this to future work.

# References

Anderson, P.; Wu, Q.; Teney, D.; Bruce, J.; Johnson, M.; Sünderhauf, N.; Reid, I.; Gould, S.; and Van Den Hengel, A. 2018. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 3674–3683.

Andreas, J. 2019. Good-enough compositional data augmentation. *arXiv preprint arXiv:1904.09545*.

Andreas, J.; Klein, D.; and Levine, S. 2017. Modular Multitask Reinforcement Learning with Policy Sketches. *arXiv:1611.01796 [cs]*. ArXiv: 1611.01796.

Beattie, C.; Leibo, J. Z.; Teplyashin, D.; Ward, T.; Wainwright, M.; Küttler, H.; Lefrancq, A.; Green, S.; Valdés, V.; Sadik, A.; et al. 2016. Deepmind lab. *arXiv preprint arXiv:1612.03801*.

Blukis, V.; Paxton, C.; Fox, D.; Garg, A.; and Artzi, Y. 2022. A persistent spatial semantic representation for high-level natural language instruction execution. In *Conference on Robot Learning*, 706–717. PMLR.

Carvalho, W.; Liang, A.; Lee, K.; Sohn, S.; Lee, H.; Lewis, R. L.; and Singh, S. 2020. Reinforcement Learning for Sparse-Reward Object-Interaction Tasks in Firstperson Simulated 3D Environments. *arXiv preprint arXiv:2010.15195*.

Chaplot, D. S.; Sathyendra, K. M.; Pasumarthi, R. K.; Rajagopal, D.; and Salakhutdinov, R. 2018. Gated-attention architectures for task-oriented language grounding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.

Chen, D.; and Mooney, R. 2011. Learning to interpret natural language navigation instructions from observations. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 25.

Chomsky, N. 2009. *Syntactic structures*. De Gruyter Mouton.

Corona, R.; Fried, D.; Devin, C.; Klein, D.; and Darrell, T. 2020. Modular Networks for Compositional Instruction Following. *arXiv preprint arXiv:2010.12764*.

Fried, D.; Andreas, J.; and Klein, D. 2017. Unified pragmatic models for generating and following instructions. *arXiv preprint arXiv:1711.04987*.

Gupta, S.; Davidson, J.; Levine, S.; Sukthankar, R.; and Malik, J. 2017. Cognitive mapping and planning for visual navigation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2616–2625.

Hill, F.; Lampinen, A.; Schneider, R.; Clark, S.; Botvinick, M.; McClelland, J. L.; and Santoro, A. 2019. Environmental drivers of systematicity and generalization in a situated agent. *arXiv preprint arXiv:1910.00571*.

Jiang, Y.; Gu, S.; Murphy, K.; and Finn, C. 2019a. Language as an abstraction for hierarchical deep reinforcement learning. *arXiv preprint arXiv:1906.07343*.

Jiang, Y.; Gu, S. S.; Murphy, K. P.; and Finn, C. 2019b. Language as an Abstraction for Hierarchical Deep Reinforcement Learning. In *Advances in Neural Information Processing Systems 32*, 9419–9431. Curran Associates, Inc. Johnson, J.; Hariharan, B.; Van Der Maaten, L.; Fei-Fei, L.; Lawrence Zitnick, C.; and Girshick, R. 2017. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2901–2910.

Kolve, E.; Mottaghi, R.; Han, W.; VanderBilt, E.; Weihs, L.; Herrasti, A.; Gordon, D.; Zhu, Y.; Gupta, A.; and Farhadi, A. 2017. AI2-THOR: An Interactive 3D Environment for Visual AI. *arXiv*.

Lake, B.; and Baroni, M. 2018. Generalization without systematicity: On the compositional skills of sequence-to-sequence recurrent networks. In *International Conference on Machine Learning*, 2873–2882. PMLR.

Lake, B. M. 2019. Compositional generalization through meta sequence-to-sequence learning. *arXiv preprint arXiv:1906.05381*.

Lynch, C.; and Sermanet, P. 2020. Grounding Language in Play. *arXiv:2005.07648 [cs]*. ArXiv: 2005.07648.

Mei, H.; Bansal, M.; and Walter, M. R. 2016. Listen, attend, and walk: Neural mapping of navigational instructions to action sequences. In *Thirtieth AAAI Conference on Artificial Intelligence*.

Min, S. Y.; Chaplot, D. S.; Ravikumar, P.; Bisk, Y.; and Salakhutdinov, R. 2021. FILM: Following Instructions in Language with Modular Methods. *arXiv preprint arXiv:2110.07342*.

Mnih, V.; Kavukcuoglu, K.; Silver, D.; Graves, A.; Antonoglou, I.; Wierstra, D.; and Riedmiller, M. 2013. Playing atari with deep reinforcement learning. *arXiv preprint arXiv:1312.5602*.

Oh, J.; Singh, S.; Lee, H.; and Kohli, P. 2017. Zero-Shot Task Generalization with Multi-Task Deep Reinforcement Learning. *arXiv:1706.05064 [cs]*. ArXiv: 1706.05064.

Raffel, C.; Shazeer, N.; Roberts, A.; Lee, K.; Narang, S.; Matena, M.; Zhou, Y.; Li, W.; and Liu, P. J. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*.

Ruis, L.; Andreas, J.; Baroni, M.; Bouchacourt, D.; and Lake, B. M. 2020. A Benchmark for Systematic Generalization in Grounded Language Understanding. *arXiv:2003.05161 [cs]*. ArXiv: 2003.05161.

Shaw, P.; Chang, M.-W.; Pasupat, P.; and Toutanova, K. 2020. Compositional generalization and natural language variation: Can a semantic parsing approach handle both? *arXiv preprint arXiv:2010.12725*.

Shridhar, M.; Thomason, J.; Gordon, D.; Bisk, Y.; Han, W.; Mottaghi, R.; Zettlemoyer, L.; and Fox, D. 2020. ALFRED: A Benchmark for Interpreting Grounded Instructions for Everyday Tasks. *arXiv:1912.01734 [cs]*. ArXiv: 1912.01734.

Singh, K. P.; Bhambri, S.; Kim, B.; Mottaghi, R.; and Choi, J. 2020. Moca: A modular object-centric approach for interactive instruction following. *arXiv preprint arXiv:2012.03208*.

Suhr, A.; Yan, C.; Schluger, J.; Yu, S.; Khader, H.; Mouallem, M.; Zhang, I.; and Artzi, Y. 2019. Executing instructions in situated collaborative interactions. *arXiv* preprint arXiv:1910.03655.

Tellex, S.; Kollar, T.; Dickerson, S.; Walter, M.; Banerjee, A.; Teller, S.; and Roy, N. 2011. Understanding natural language commands for robotic navigation and mobile manipulation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 25.

Van Hasselt, H.; Guez, A.; and Silver, D. 2016. Deep reinforcement learning with double q-learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30.

Zhu, Y.; Gordon, D.; Kolve, E.; Fox, D.; Fei-Fei, L.; Gupta, A.; Mottaghi, R.; and Farhadi, A. 2017. Visual semantic planning using deep successor representations. In *Proceedings of the IEEE international conference on computer vision*, 483–492.