
Learning compositional tasks from language instructions

Anonymous Author(s)

Affiliation
Address
email

Abstract

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

13 1 Introduction

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

21 Prior work studying the linguistic “systematicity” of
22 neural networks have focused on sequence mapping tasks in a supervised learning setting (Lake and Baroni,
23 2018; Lake, 2019; Andreas, 2019). In this work,
24 we are interested in compositional generalization of a reinforcement learning agent following natural language
25 instructions in an embodied environment. In particular, we explore the hypothesis that a language-
26 conditioned reinforcement learning agent with a compositional inductive bias in its behavioral policy will
27 exhibit systematic generalization to unobserved natural language instructions.
28 In this work,
29 we are interested in compositional generalization of a reinforcement learning agent following natural language
30 instructions in an embodied environment. In particular, we explore the hypothesis that a language-
31 conditioned reinforcement learning agent with a compositional inductive bias in its behavioral policy will
32 exhibit systematic generalization to unobserved natural language instructions.
33 In this work,
34 we are interested in compositional generalization of a reinforcement learning agent following natural language
35 instructions in an embodied environment. In particular, we explore the hypothesis that a language-
36 conditioned reinforcement learning agent with a compositional inductive bias in its behavioral policy will
37 exhibit systematic generalization to unobserved natural language instructions.

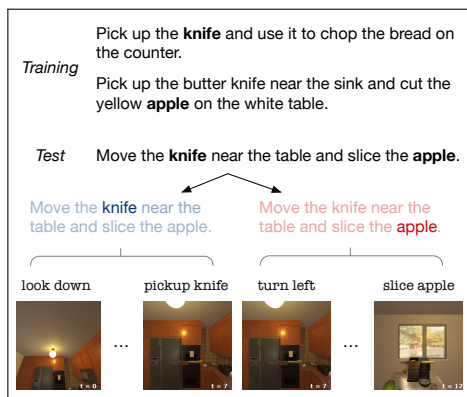


Figure 1: Zero-shot generalization to an unseen task of slicing an apple. The test task is composed of known primitive subtasks – *picking up a knife* and *slicing the apple* – each of which were encountered in training tasks. Our agent learns to decompose a natural language task description into subtasks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

37 There has been a flurry of recent work on embodied learning tasks such as question answering (Gupta
38 et al., 2017), navigation (Anderson et al., 2018) and object interaction (Shridhar et al., 2020; Carvalho
39 et al., 2020) in embodied settings. In particular, the ALFRED task (Shridhar et al., 2020) studies
40 agents that exploit detailed natural language instructions to generalize to novel instructions in novel
41 environments at test time. Such existing benchmarks offer limited flexibility to study systematic
42 generalization since (i) the benchmarks were not built for this purpose and it is unclear to what extent
43 systematic generalization skills are required to solve the tasks and (ii) the tasks demand challenging
44 reasoning skills such as visual recognition and planning over large number of time-steps which makes
45 it difficult to study compositional generalization ability in isolation.

46 In this work we develop a set of tasks in the AI2Thor virtual home environment (Kolve et al., 2017)
47 which test the compositionality of embodied agents. In order to make progress in systematic gener-
48 alization, we make two simplifying assumptions: we assume access to an oracle object recognizer
49 and we study generalization in a single kitchen layout. This allows us to study the degree to which a
50 policy captures the systematicity in language by studying its zero-shot generalization performance on
51 held out natural language instructions.

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

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

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

71 2 Related work

72 **Compositional generalization** Prior work has studied compositional generalization in sequence
73 mapping tasks. Benchmarks such as SCAN (Lake and Baroni, 2018) and gSCAN (Ruis et al., 2020)
74 study translating synthetic text descriptions to an action sequence (e.g. jump twice \rightarrow JUMP JUMP).
75 gSCAN couples SCAN instances with entities in a grid environment and solving a task requires
76 grounding the text and entities similar to our work. Prior approaches for these benchmarks impose
77 compositional inductive biases in models by augmenting models with memory (Lake, 2019) and
78 data augmentation (Andreas, 2019). In this work we use attention mechanisms and introduce a novel
79 policy parameterization to impose compositional inductive biases.

80 **Text based embodied control** Advances in photo-realistic simulation environments such as Deep-
81 Mind Lab (Beattie et al., 2016) and AI2Thor (Kolve et al., 2017) have driven recent progress in
82 embodied agents that learn from text instructions. Chaplot et al. (2018) consider a simple navigation
83 task where an agent has to move to an object specified by a set of attributes such as shape and
84 color. They propose the gated attention model to generalize compositionally in the attribute space.
85 Hill et al. (2019) consider systematic generalization in 2D and 3D environments with synthetic
86 text instructions. Compared to these work, we consider object interaction tasks in a photo realistic
87 simulated environment with human-authored language instructions.

88 ALFRED (Shridhar et al., 2020) couples tasks in the AI2Thor environment with detailed text
89 descriptions of tasks. In contrast, we consider a simplified setup of learning compositional skills from

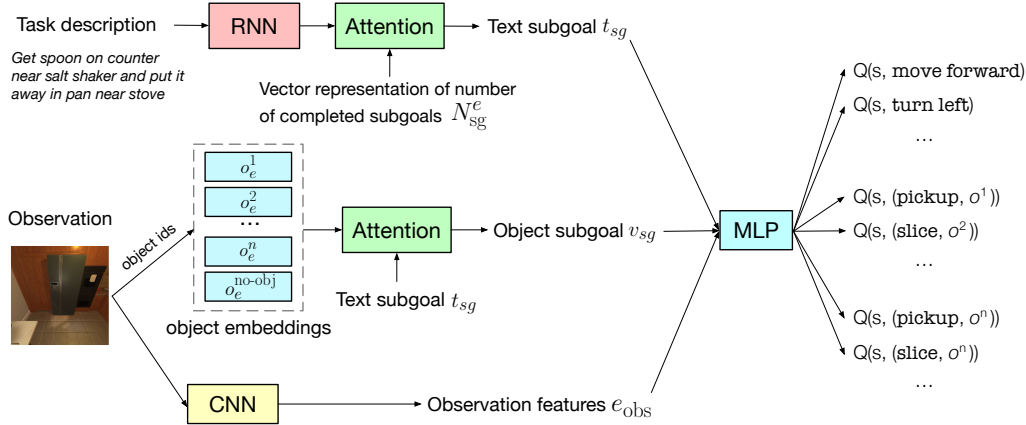


Figure 2: Approach Overview: We perform attention over the text instruction to construct an embedding t_{sg} that represents the current subgoal. The text embedding subgoal t_{sg} attends to scene object embeddings to construct an object subgoal representation v_{sg} . An MLP takes t_{sg} , v_{sg} and 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.

90 high-level task descriptions. We further do not assume access to expert task demonstrations. These
 91 assumptions allows us to focus on compositional generalization to zero-shot tasks, which is not the
 92 main goal of the ALFRED benchmark. However, the approach presented here can potentially be
 93 applicable to ALFRED when combined with learning from demonstrations.

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

101 Finally, generalization to unseen instructions has been considered in prior work such as Oh et al.
 102 (2017); Lynch and Sermanet (2020), although compositional generalization is not their main focus.

103 3 Problem

104 We consider an embodied agent acting in a kitchen environment to solve basic tasks from language
 105 instructions (See Fig. 4 for an example task). At the beginning of an episode the agent receives a text
 106 instruction τ . Our goal is to learn a policy $\pi(a|s, \tau)$; $a \in \mathcal{A}$, $s \in \mathcal{S}$ that predicts actions in order to
 107 complete tasks. The agent state s is partially observable – it receives an egocentric observation obs
 108 of the scene. We further assume that an oracle object recognition model provides the object ids for
 109 objects in the egocentric observation.

110 The action space consists of navigation and object interaction actions $\mathcal{A} = \mathcal{A}_{nav} \cup \mathcal{A}_{int}$. There are
 111 8 navigation actions $\mathcal{A}_{nav} = \{move\ forward, move\ back, move\ left, move\ right, turn\ left, turn\ right,$
 112 $look\ up, look\ down\}$. Interaction actions $\mathcal{A}_{int} = \mathcal{B} \times \mathcal{I}$ are specified using an interaction $b \in \mathcal{B}$ and
 113 an object id $o \in \mathcal{I}$ where $\mathcal{B} = \{pickup, place, slice, toggle\ on, toggle\ off, turn\ on, turn\ off\}$ and \mathcal{I} is a
 114 pre-defined set of identifiers of objects that are available to the agent for interaction in the current
 115 observation.

116 The agent receives a positive reward for successfully completing a task. It also receives a small
 117 negative reward for every time-step. In addition, we also assume that every correct object interaction
 118 receives a positive reward. In addition to providing a denser learning signal, the rewards are also used
 119 to identify subgoals as described in section 4.1. In practice such dense rewards may be unavailable,
 120 but this is outside the scope of our study and left as future work.

121 **4 Approach**

122 We approach the problem by considering a task τ to be composed of subgoals g_1, \dots, g_n , where each
 123 subgoal g_i involves navigating to a particular object and interacting with it. For example, the task
 124 *place an apple on the table* involves finding the apple and picking it up, followed by navigating to the
 125 table and putting down the apple, which can be considered to be the two subgoals for executing the
 126 task. Each object interaction required to complete the task thus corresponds to a subgoal. Since every
 127 subgoal completion receives a positive reward, the number of subgoals completed at every time-step
 128 N_{sg} is known to the agent. The subgoals themselves are not known to the agent – we use attention on
 129 the text instruction to compute a latent subgoal representation.

130 **4.1 Text subgoal inference**

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

$$t_{sg} = \text{Attention}(\text{query} = N_{sg}^e, \text{keys} = \text{values} = H) = \sum_{h \in H} \frac{\exp((Qs)^\top (Kh))}{\sum_{h' \in H} \exp((Qs)^\top (Kh'))} Vh \quad (1)$$

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

139 **4.2 Cross-modal reasoning**

140 Given the text subgoal t_{sg} , we use an attention mechanism to reason about objects in the scene
 141 within some distance to the agent. This helps the agent understand if objects of interest relevant
 142 to the subgoal are present nearby. Let the set of nearby scene objects be $O = \{o^1, \dots, o^n\}$, where
 143 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
 144 embedding table that produces object embeddings $O_e = \{o_e^1, \dots, o_e^n\}$. The cross-modal attention
 145 is given by Eq. (2) where the text subgoal attends to the scene object embeddings (Q', K', V' are
 146 learnable parameter matrices). We augment the scene objects embeddings O_e with an additional
 147 learned embedding o_e^{no-obj} which is expected to absorb any probability mass not assigned to scene
 148 objects $O'_e = O_e \cup o_e^{no-obj}$. The attention produces an *object subgoal* embedding v_{sg} .

$$v_{sg} = \text{Attention}(\text{query} = t_{sg}, \text{keys} = \text{values} = O'_e) = \sum_{o_e \in O'_e} \frac{\exp((Q't_{sg})^\top (K'o_e))}{\sum_{o'_e \in O'_e} \exp((Q't_{sg})^\top (K'o'_e))} V'o_e \quad (2)$$

149 **4.3 Policy learning**

150 We use a deep Q-learning algorithm to train a policy (Mnih et al., 2013), where a neural network
 151 is trained to approximate the state-action value function $Q(s, a)$. Given the current observation,
 152 text subgoal and object subgoal, the state-action value for a navigation action $a \in \mathcal{A}_{nav}$ is given by
 153 Eq. (3), where f_{nav} is an MLP (multi-layer perceptron) and $e_{obs} = f_{CNN}(\text{obs})$ is a feature vector of
 154 the observation image computed using a CNN encoder.

$$Q_{nav}(s, a) = f_{nav}(a|e_{obs}, t_{sg}, v_{sg}) \quad (3)$$

155 The state-action values for interaction actions $a = (b, o) \in \mathcal{B} \times \mathcal{I}$ can be analogously modeled as in
 156 Eq. (4). We found it helpful to decompose the state-action value in an additive fashion over an action
 157 score f_{int}^a and an object score f_{int}^o as in Eq. (5). Intuitively, f_{int}^a learns to model action preferences,
 158 whereas f_{int}^o learns to ground text goals to physical objects. In addition to sharing parameters across
 159 actions and objects, this decomposition allows us to model state-action values of object interactions
 160 not experienced during training, as long as the specific interaction and the object were encountered.

Task type	Task descriptions
<i>pick up pot</i>	Go to the stove and pick up the pot. Pick up the pot on the bottom right burner on the stove. Take the cooking pot from the stove.
<i>place spoon in pan</i>	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. Move spoon from the counter and into the pan.
<i>slice bread with knife</i>	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. Pick up the sharp knife with a yellow handle, and slice the bread on the white 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.

161 Unless specified otherwise we use the decomposed value function $Q_{\text{int}}^{\text{add}}$ in our experiments.

$$Q_{\text{int}}^{\text{full}}(s, a) = f_{\text{int}}(a|e_{\text{obs}}, t_{\text{sg}}, v_{\text{sg}}) \quad (4)$$

$$Q_{\text{int}}^{\text{add}}(s, a) = f_{\text{int}}^a(b|e_{\text{obs}}, t_{\text{sg}}, v_{\text{sg}}) + f_{\text{int}}^o(o|t_{\text{sg}}) \text{ where } a = (b, o) \in \mathcal{B} \times \mathcal{I} \quad (5)$$

162 In summary, the state-action value function is modeled as in Eq. (6).

$$Q(s, a) = \begin{cases} Q_{\text{nav}}(s, a); & a \in \mathcal{A}_{\text{nav}} \\ Q_{\text{int}}^{\text{add}}(s, a); & a \in \mathcal{A}_{\text{int}} \end{cases} \quad (6)$$

163 The overall model (see Fig. 2 for an illustration) including parameters of the subgoal inference (Eq.
164 1) and cross-modal reasoning (Eq. 2) components, as well as the MLPs in Eqs. (3) and (5) are trained
165 end-to-end using a double-DQN algorithm (Van Hasselt et al., 2016). Once the model has been
166 trained we construct a greedy policy by choosing actions with the highest state-action values for
167 inference.

168 5 Experiments

169 5.1 Tasks

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

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

178 We use Amazon Mechanical Turk to collect natural language descriptions of tasks for training and
179 evaluation. A turker is shown key observation frames during the execution of a particular task and is
180 asked to describe in a sentence how they would describe the task to a robot. Turkers were instructed
181 to do their best to correctly identify task relevant objects. But often descriptions from the turkers
182 incorrectly identify objects such as identifying a potato as an avocado. Such descriptions were
183 manually fixed so that the correct object identities are mentioned in the instructions. We collected
184 5 natural language descriptions each for 35 tasks that include *pickup*, *place* and *slice* tasks. The
185 descriptions consist of 170 unique tokens and have an average length of 12 tokens. Table 1 shows
186 example descriptions collected for some tasks. See appendix B for instructions given to Turkers in
187 the data collection process.

188 The pickup tasks are used for evaluating multi-task and zero-shot generalization with seen and unseen
189 descriptions of tasks. We use 10 *pickup* tasks - *pickup X* where $X \in \{\text{apple, bread, tomato, potato,}$

	place tasks	slice tasks
Training tasks	<i>place apple in plate</i> <i>place butterknife in plate</i> <i>place spoon in plate</i> <i>place butterknife in pan</i> <i>place potato in pan</i> <i>place spoon in pan</i> <i>place apple in pot</i> <i>place butterknife in pot</i> <i>place potato in pot</i>	<i>slice apple with knife</i> <i>slice tomato with knife</i> <i>slice bread with knife</i> <i>slice apple with butterknife</i> <i>slice potato with butterknife</i> <i>slice bread with butterknife</i>
Test tasks (obj-obj setting)	<i>place potato in plate</i> <i>place apple in pan</i> <i>place spoon in pot</i>	<i>slice potato with knife</i> <i>slice tomato with butterknife</i>
Test tasks (task-obj setting)	<i>place knife in plate</i> <i>place knife in pan</i> <i>place knife in pot</i>	<i>slice lettuce with knife</i> <i>slice lettuce with butterknife</i>

Table 2: Task types used for training and testing on place and slice tasks. The *obj-obj* setting considers test tasks 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).

190 *lettuce, spoon, bread, butter knife, plate, pot*}. These tasks are used for evaluating generalization to
191 seen and unseen descriptions of known short-horizon tasks. They are also used in generalization to
192 longer horizon tasks as described later in this section.

193 The *place* and *slice* tasks are used for evaluating generalization to longer-horizon unseen tasks.
194 Table 2 shows tasks used for training and evaluation. In addition to multitask generalization, we use
195 these tasks to study zero-shot compositional generalization to unseen task descriptions. The unseen
196 descriptions can correspond to tasks that were encountered during training, similar to the *pickup* tasks.
197 A more challenging generalization scenario is to generalize to text descriptions of unseen tasks.

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

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

208 The training tasks in Table 2 were designed such that each object argument appears in multiple tasks.
209 Furthermore, when choosing object arguments for a given task type, we prioritized objects that appear
210 in as many tasks as possible. For instance, in the pickup and place tasks setup, the objects were plate,
211 pan, pot, spoon, etc. where each object appears in at least three of the training tasks. This ensures that
212 there are enough occurrences of each object type for the agent to understand and ground the object
213 type. It also helps the agent disentangle the notion of an object versus a task in a given instruction.

214 5.2 Baselines

215 We compare the proposed approach against the following baselines.

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

219 **Gated Attention** This architecture (Chaplot et al., 2018) combines the instruction representation
220 with the visual observation using a gated attention operation. The fused representation is fed to an
221 MLP which models the state-action values. All models and baselines are trained using the DDQN
222 Q-learning algorithm.

Tasks Descriptions	Training tasks		Test tasks		
	seen	unseen	unseen obj-obj	unseen task-obj	
Model	RNN	0.65	0.65	0.26	0.13
	Gated Attention	0.92	0.85	0.66	0.34
	Ours				
	(a) Q_{int}^{add} (no cross modal)	0.89	0.76	0.84	0.77
	(b) Q_{int}^{full} + cross modal	0.93	0.85	0.44	0.34
(c) Q_{int}^{add} + 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.

223 5.3 Hyperparameters

224 Word embeddings and the RNN have representation size 32. Objects are represented by embeddings
 225 of size 32 from an embedding table. The CNN observation features have size 512 and the CNN
 226 encoder has 1.7M parameters, which constitutes 90% of the overall model parameters. The MLPs in
 227 Eqs. (3) and (4) are single hidden layer MLPs with 256 hidden units and ReLU activation.

228 5.4 Results

229 **Short-horizon tasks** We first consider pickup tasks that involve a single object interaction. In these
 230 tasks the agent has to identify the object reference mentioned in the text description and then find and
 231 pick up the relevant object. We train and evaluate on 10 pickup task types. Four text descriptions of
 232 each task type are part of the training set and the remaining descriptions (i.e., 1 per task type) are part
 233 of the test set. Identifying the correct subgoal for these tasks involves paying attention to the verbs
 234 and nouns in the task description as well as the overall context. On the training and test descriptions,
 235 our agent trained from scratch achieves success rates of 0.9, 0.92 respectively.

236 **Longer-horizon Tasks** We now consider tasks that involve two
 237 subgoals, which includes the place and slice tasks in Table 2. Jointly
 238 learning text grounding and subgoal inference for long horizon tasks
 239 can be challenging. We thus consider a curriculum learning strategy
 240 where an agent is gradually trained on tasks of increasingly longer
 241 horizon. The agent is first pre-trained on the pickup tasks as de-
 242 scribed in the previous section, and then fine-tuned on the training
 243 tasks in Table 2. Fig. 3 compares the learning progress of agents
 244 trained from scratch and an agent that has been pre-trained on the
 245 pickup tasks. The pre-trained agent learns twice as fast compared
 246 to the agent trained from scratch and achieves perfect success rate
 247 on training tasks.

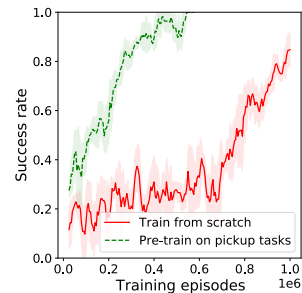


Figure 3: Learning progress of agent trained from scratch and agent pre-trained on pickup tasks.

248 **Generalization** Table 3 shows the average task completion success rate of models under different
 249 generalization scenarios. The RNN and Gated Attention baselines are limited by the fact that the text
 250 instruction is represented using the same encoding across all time-steps, which has limited ability to
 251 capture compositional information. The inductive bias of Gated Attention enables better performance,
 252 but it has difficulty generalizing to unseen tasks. The attention based model outperforms these
 253 baselines, which indicates that the attention mechanism helps exploit compositional information in
 254 the instruction better than a fixed encoding.

255 In addition to better performance, the attention model has the advantage of being more interpretable.
 256 Fig. 4 shows the agent’s actions and the attention pattern over time for an example task. The agent
 257 learns to identify object references in the instruction and uses attention as a sub-goal representation.
 258 This mimics a hierarchical policy where a high-level controller provides a sub-goal and a low-level
 259 controller executes it Jiang et al. (2019b). The agent further learns to ground object references in the

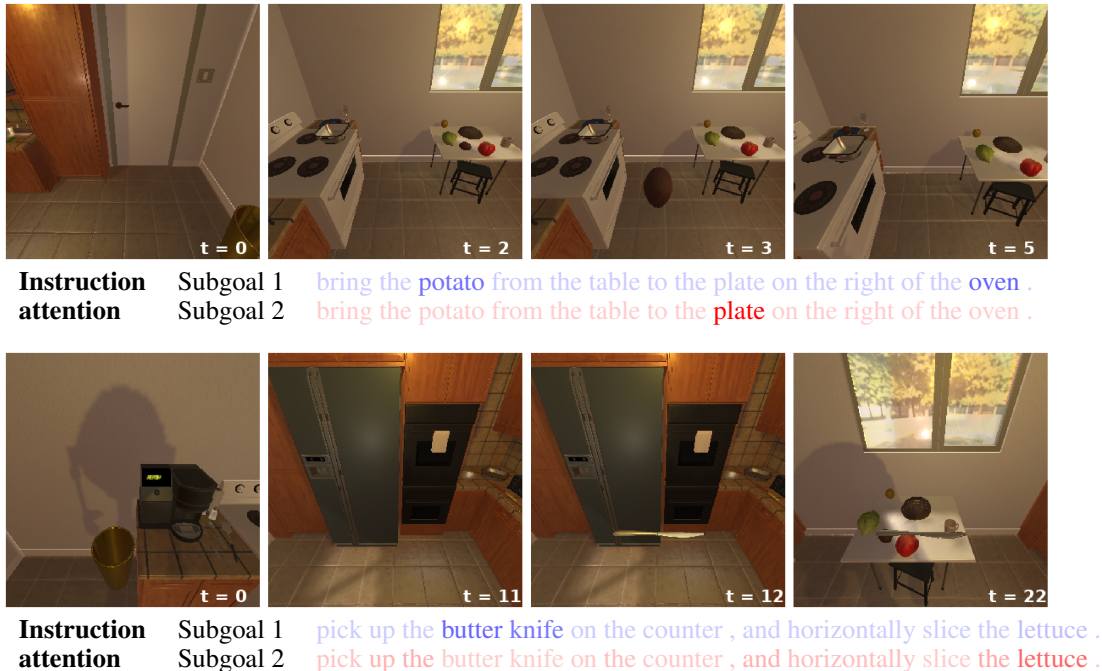


Figure 4: Agent’s observation at different time-steps while performing a place task and a slice task. The attention distribution in the text goal inference component while executing each subgoal is also given below the agent observations.

260 text instruction to objects in the scene. Notably, these attention patterns and grounding are learned
 261 from the reward signal alone without any other supervision. More example of agent trajectories are
 262 given in appendix A.

263 5.5 Ablations

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

266 **Cross modal reasoning** We examine model performance without the cross modal reasoning com-
 267 ponent. In this case the MLPs in Eqs. (3) and (5) only receive the text subgoal and observation
 268 encoding as inputs and the visual subgoal v_{sg} is omitted. From rows (a) and (c) in table Table 3 it
 269 is clear that the cross-modal reasoning components helps ground text in scene objects and enables
 270 better generalization across all settings.

271 **Interaction Q-values** We examine the benefit of decomposing the value function approximation
 272 of interaction actions in an additive fashion in Q_{int}^{add} (Eq. (5)). We compare it against Q_{int}^{full} (Eq. (4)),
 273 which treats each (interaction, object) pair as a separate atomic action. Comparing rows (b), (c) in
 274 Table 3 we see that the additive decomposition is crucial for generalization to unseen tasks.

275 6 Conclusion

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

281 Compared to existing benchmarks such as ALFRED we made simplifying assumptions such as
 282 oracle visual recognition, relatively short horizon tasks and generalization within single kitchen

283 layout which allows us to focus on compositional generalization in embodied settings. However, the
284 ideas presented here can potentially be combined with curriculum learning and learning from human
285 demonstrations to perform complex tasks that require planning over hundreds of time-steps such as
286 in the ALFRED setting, and we leave this to future work.

287 References

- 288 Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid,
289 Stephen Gould, and Anton Van Den Hengel. 2018. Vision-and-language navigation: Interpreting
290 visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE*
291 *Conference on Computer Vision and Pattern Recognition*, pages 3674–3683.
- 292 Jacob Andreas. 2019. Good-enough compositional data augmentation. *arXiv preprint*
293 *arXiv:1904.09545*.
- 294 Jacob Andreas, Dan Klein, and Sergey Levine. 2017. <http://arxiv.org/abs/1611.01796> Modular Multi-
295 task Reinforcement Learning with Policy Sketches. *arXiv:1611.01796 [cs]*. ArXiv: 1611.01796.
- 296 Charles Beattie, Joel Z Leibo, Denis Teplyashin, Tom Ward, Marcus Wainwright, Heinrich Küttler,
297 Andrew Lefrancq, Simon Green, Víctor Valdés, Amir Sadik, et al. 2016. Deepmind lab. *arXiv*
298 *preprint arXiv:1612.03801*.
- 299 Wilka Carvalho, Anthony Liang, Kimin Lee, Sungryull Sohn, Honglak Lee, Richard L Lewis,
300 and Satinder Singh. 2020. Reinforcement learning for sparse-reward object-interaction tasks in
301 first-person simulated 3d environments. *arXiv preprint arXiv:2010.15195*.
- 302 Devendra Singh Chaplot, Kanthashree Mysore Sathyendra, Rama Kumar Pasumarthi, Dheeraj
303 Rajagopal, and Ruslan Salakhutdinov. 2018. Gated-attention architectures for task-oriented
304 language grounding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 32.
- 305 Noam Chomsky. 2009. *Syntactic structures*. De Gruyter Mouton.
- 306 Saurabh Gupta, James Davidson, Sergey Levine, Rahul Sukthankar, and Jitendra Malik. 2017.
307 Cognitive mapping and planning for visual navigation. In *Proceedings of the IEEE Conference on*
308 *Computer Vision and Pattern Recognition*, pages 2616–2625.
- 309 Felix Hill, Andrew Lampinen, Rosalia Schneider, Stephen Clark, Matthew Botvinick, James L
310 McClelland, and Adam Santoro. 2019. Environmental drivers of systematicity and generalization
311 in a situated agent. *arXiv preprint arXiv:1910.00571*.
- 312 Yiding Jiang, Shixiang Gu, Kevin Murphy, and Chelsea Finn. 2019a. Language as an abstraction for
313 hierarchical deep reinforcement learning. *arXiv preprint arXiv:1906.07343*.
- 314 YiDing Jiang, Shixiang (Shane) Gu, Kevin P Murphy, and Chelsea Finn. 2019b.
315 [http://papers.nips.cc/paper/9139-language-as-an-abstraction-for-hierarchical-deep-](http://papers.nips.cc/paper/9139-language-as-an-abstraction-for-hierarchical-deep-reinforcement-learning.pdf)
316 [reinforcement-learning.pdf](http://papers.nips.cc/paper/9139-language-as-an-abstraction-for-hierarchical-deep-reinforcement-learning.pdf) Language as an Abstraction for Hierarchical Deep Reinforcement
317 Learning. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d̄textquotesingle AlchÃ©-Buc,
318 E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems 32*, pages
319 9419–9431. Curran Associates, Inc.
- 320 Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and
321 Ross Girshick. 2017. Clevr: A diagnostic dataset for compositional language and elementary visual
322 reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*,
323 pages 2901–2910.
- 324 Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Daniel
325 Gordon, Yuke Zhu, Abhinav Gupta, and Ali Farhadi. 2017. AI2-THOR: An Interactive 3D
326 Environment for Visual AI. *arXiv*.
- 327 Brenden Lake and Marco Baroni. 2018. Generalization without systematicity: On the compositional
328 skills of sequence-to-sequence recurrent networks. In *International Conference on Machine*
329 *Learning*, pages 2873–2882. PMLR.

- 330 Brenden M Lake. 2019. Compositional generalization through meta sequence-to-sequence learning.
331 *arXiv preprint arXiv:1906.05381*.
- 332 Corey Lynch and Pierre Sermanet. 2020. <http://arxiv.org/abs/2005.07648> Grounding Language in
333 Play. *arXiv:2005.07648 [cs]*. ArXiv: 2005.07648.
- 334 Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan
335 Wierstra, and Martin Riedmiller. 2013. Playing atari with deep reinforcement learning. *arXiv*
336 *preprint arXiv:1312.5602*.
- 337 Junhyuk Oh, Satinder Singh, Honglak Lee, and Pushmeet Kohli. 2017.
338 <http://arxiv.org/abs/1706.05064> Zero-Shot Task Generalization with Multi-Task Deep Re-
339 inforcement Learning. *arXiv:1706.05064 [cs]*. ArXiv: 1706.05064.
- 340 Laura Ruis, Jacob Andreas, Marco Baroni, Diane Bouchacourt, and Brenden M. Lake. 2020.
341 <http://arxiv.org/abs/2003.05161> A Benchmark for Systematic Generalization in Grounded Lan-
342 guage Understanding. *arXiv:2003.05161 [cs]*. ArXiv: 2003.05161.
- 343 Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi,
344 Luke Zettlemoyer, and Dieter Fox. 2020. <http://arxiv.org/abs/1912.01734> ALFRED: A Benchmark
345 for Interpreting Grounded Instructions for Everyday Tasks. *arXiv:1912.01734 [cs]*. ArXiv:
346 1912.01734.
- 347 Hado Van Hasselt, Arthur Guez, and David Silver. 2016. Deep reinforcement learning with double
348 q-learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 30.
- 349 Yuke Zhu, Daniel Gordon, Eric Kolve, Dieter Fox, Li Fei-Fei, Abhinav Gupta, Roozbeh Mottaghi, and
350 Ali Farhadi. 2017. Visual semantic planning using deep successor representations. In *Proceedings*
351 *of the IEEE international conference on computer vision*, pages 483–492.