# Representation Learning for Grounded Spatial Reasoning

**Michael Janner, Karthik Narasimhan, and Regina Barzilay**
Computer Science and Artificial Intelligence Laboratory
Massachusetts Institute of Technology
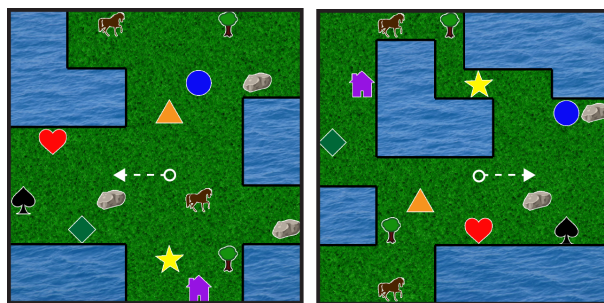`{janner, karthikn, regina}@csail.mit.edu`

## Abstract

The interpretation of spatial references is highly contextual, requiring joint inference over both language and the environment. We consider the task of spatial reasoning in a simulated environment, where an agent can act and receive rewards. The proposed model learns a representation of the world steered by instruction text. This design allows for precise alignment of local neighborhoods with corresponding verbalizations, while also handling global references in the instructions. We train our model with reinforcement learning using a variant of generalized value iteration. The model outperforms state-of-the-art approaches on several metrics, yielding a 45% reduction in goal localization error. [1]

## 1 Introduction

Understanding spatial references in natural language is essential for successful human-robot communication and autonomous navigation. This problem is challenging because interpretation of spatial references is highly context-dependent. For instance, the instruction *"Reach the cell above the westernmost rock"* translates into different goal locations in the two environments shown in Figure 1. Therefore, to enable generalization to new, unseen worlds, the model must jointly reason over the instruction text and environment configuration. Moreover, the richness and flexibility in verbalizing spatial references further complicates interpretation of such instructions.

[1]Code and dataset are available at `https://github.com/JannerM/spatial-reasoning`



*Reach the cell above the westernmost rock*

Figure 1: Sample 2D worlds and an instruction describing a goal location. The optimal path from a common start position, denoted by a white dashed line, varies considerably with changes in the map layout.

In this paper, we explore the problem of spatial reasoning in the context of interactive worlds. Specifically, we assume access to a simulated environment, in which an agent can take actions to interact with the world and is rewarded for reaching the location specified by the language instruction. This feedback is the only source of supervision the model uses for interpreting spatial references.

The key modeling task here is to induce a representation that closely ties environment observations and linguistic expressions. In prior work, this issue was addressed by learning representations for each modality and then combining them, for instance, with concatenation (Misra et al., 2017). While this approach captures high-level correspondences between instructions and maps, it does not encode de-

tailed, lower-level mappings between specific positions on the map and their descriptions. As our experiments demonstrate, combining the language and environment representations in a spatially localized manner yields significant performance gains on the task.

To this end, our model uses the instruction text to drive the learning of the environment representation. We start by converting the instruction text into a real-valued vector using a recurrent neural network with LSTM cells (Hochreiter and Schmidhuber, 1997). Using this vector as a kernel in a convolution operation, we obtain an instruction-conditioned representation of the state. This allows the model to reason about immediate local neighborhoods in references such as *"two cells to the left of the triangle"*. We further augment this design to handle global references that involve information concerning the entire map (e.g. *"the westernmost rock"*). This is achieved by predicting a global value map using an additional component of the instruction representation. The entire model is trained with reinforcement learning using the environmental reward signal as feedback.

We conducted our experiments using a 2D virtual world as shown in Figure 1. Overall, we created over 3,300 tasks across 200 maps, with instructions sourced from Mechanical Turk. We compare our model against two state-of-the-art systems adapted for our task (Misra et al., 2017; Schaul et al., 2015). The key findings of our experiments are threefold. First, our model can more precisely interpret instructions than baseline models and find the goal location, yielding a 45% reduction in Manhattan distance error over the closest competitor. Second, the model can robustly generalize across new, unseen map layouts. Finally, we demonstrate that factorizing the instruction representation enables the model to sustain high performance when handling both local and global references.

## 2   Related Work

**Spatial reasoning in text**   This topic has attracted both theoretical and practical interest. From the linguistic and cognitive perspectives, research has focused on the wide range of mechanisms speakers use to express spatial relations (Tenbrink, 2007; Viethen and Dale, 2008; Byrne and Johnson-Laird, 1989;

Li and Gleitman, 2002). The practical implications of this research are related to autonomous navigation (Moratz and Tenbrink, 2006; Levit and Roy, 2007; Tellex et al., 2011) and human-robot interaction (Skubic et al., 2004).

Previous computational approaches include techniques such as proximity fields (Kelleher et al., 2006), spatial templates (Levit and Roy, 2007) and geometrically defined mappings (Moratz and Tenbrink, 2006; Kollar et al., 2010). More recent work in robotics has integrated text containing position information with spatial models of the environment to obtain accurate maps for navigation (Walter et al., 2013; Hemachandra et al., 2014). Most of these approaches typically assume access to detailed geometry or other forms of domain knowledge. In contrast to these knowledge-rich approaches, we are learning spatial reference via interaction with the environment, acquiring knowledge of the environment in the process.

**Instruction following**   Spatial reasoning is a common element in many papers on instruction following (MacMahon et al., 2006; Vogel and Jurafsky, 2010; Chen and Mooney, 2011; Artzi and Zettlemoyer, 2013; Kim and Mooney, 2013; Andreas and Klein, 2015). As a source of supervision, these methods assume access to demonstrations, which specify the path corresponding with provided instructions. In our setup, the agent is only driven by the final rewards when the goal is achieved. This weaker source of supervision motivates development of new techniques not considered in prior work.

More recently, Misra et al. (2017) proposed a neural architecture for jointly mapping instructions and visual observations (pixels) to actions in the environment. Their model separately induces text and environment representations, which are concatenated into a single vector that is used to output an action policy. While this representation captures coarse correspondences between the modalities, it doesn't encode mappings at the level of local neighborhoods, negatively impacting performance on our task.

**Universal value functions**   The idea of generalized value functions has been explored before in Schaul et al. (2015). The technique, termed UVFA,

presents a clever trick of factorizing the value function over states and goals using singular value decomposition (SVD) and then learning a regression model to predict the low-rank vectors. This results in quick and effective generalization to all goals in the same state space. However, their work stops short of exploring generalization over map layouts, which our model is designed to handle. Furthermore, our setup also involves specifying goals using natural language instructions, which is different from the coordinate-style specification used in that work.

## 3   General Framework

**Task setup**   We model our task as a Markov Decision Process (MDP), where an autonomous agent is placed in an interactive environment with the capability to choose actions that can affect the world. A goal is described in text, and rewards are available to the agent correspondingly. The MDP can be represented by the tuple $\langle S, A, X, T, R \rangle$, where $S$ is the set of all possible state configurations, $A$ is the set of actions available to the agent, $X$ is the set of all goal specifications[2] in natural language, $T(s'|s, a, x)$ is the transition distribution, and $R(s, x)$ is the reward function. A state $s \in S$ includes information such as the locations of different entities along with the agent's own position. In this work, $T$ is deterministic in the environments considered; however, our methods also apply in the stochastic case.

**Text instructions**   Prior work has investigated human usage of different types of referring expressions to describe spatial relations (Levinson, 2003; Viethen and Dale, 2008). In order to build a robust instruction following system, we examine several categories of spatial expressions that exhibit the wide range of natural language goal descriptions. Specifically, we consider instructions that utilize objects/entities present in the environment to describe a goal location. These instructions can be categorized into three groups:

(a) Text referring to a specific entity. (e.g. *Go to the circle.*)

(b) Text specifying a location using a single refer-

ent entity. (e.g. *Reach the cell above the westernmost rock.*)

(c) Text specifying a location using multiple referent entities. (e.g. *Move to the goal two squares to the left of the heart and top right of house.*)

These three categories exemplify an increasing level of complexity, with the last one having multiple levels of indirection.

In each category, we have both local and global references to objects. Local references require an understanding of spatial prepositional phrases such as 'above', 'in between' and 'next to' in order to determine the precise goal location. This comprehension is invariant to the global position of the object landmark(s) provided in the instruction. A global reference, on the other hand, contains superlatives such as 'easternmost' and 'topmost', which require reasoning over the entire map. For example, in the case of (a) above, a local reference would describe a unique object[3] (e.g. *Go to the circle*), whereas a global reference might require comparing the positions of all objects of a specific type (e.g. *Go to the northernmost tree*).

A point to note is that we do not assume any access to mapping from instructions to objects or entities in the worlds or a knowledge of spatial ontology – the system has to learn this entirely through feedback from the environment.

**Generalized Value Iteration**   Learning to reach the goal while maximizing cumulative reward can be done by using a *value function $V(s)$* (Sutton and Barto, 1998) which represents the agent's notion of expected future reward from state $s$. A popular algorithm to learn an optimal value function is Value Iteration (VI) (Bellman, 1957), which uses the technique of dynamic programming.

In the standard Bellman equation, the value function is dependent solely on state. Schaul et al. (2015) proposed a value function $V(s, g)$ describing the expected reward from being in state $s$ given goal $g$, capturing that state values are goal-dependent and that a single environment can offer many such goals. We also make use of such a *generalized value function*, although our goals are not observed directly as
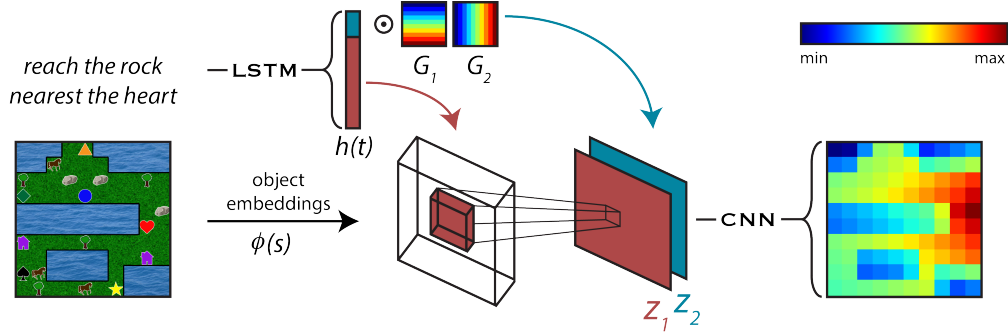
---

Figure 2: A schematic depiction of our model. Text instructions are represented as a vector $h(t)$ and states as embeddings $\phi(s)$. A portion of the text representation is used as a convolutional kernel on $\phi(s)$, giving a *text-conditioned* local state representation $z_1$. The remaining components are used as coefficients in a linear combination of gradient functions to give a global map-level representation $z_2$. $z_1$ and $z_2$ are concatenated and input to a convolutional neural network to predict the final value map.

coordinate locations or states themselves but rather described in natural language. With $x$ denoting a textual description of a goal, our VI update equations are:

$$Q(s, a, x) = R(s, x) + \gamma \sum_{s' \in S} T(s'|s, a, x)V(s', x)$$

$$V(s, x) = \max_a Q(s, a, x) \qquad (1)$$

where $Q$ is the action-value function, tracking the value of choosing action $a$ in state $s$. Once an optimal value function is learned, a straightforward action policy is:

$$\pi(s, x) = \arg\max_a Q(s, a, x) \qquad (2)$$

## 4   Model

Generalization over both environment configurations and text instructions requires a model that meets two desiderata. First, it must have a flexible representation of goals, one which can encode both the local structure and global spatial attributes inherent to natural language instructions. Second, it must be compositional, in order to learn a generalizable representation of the language even though each unique instruction will only be observed with a single map during training. Namely, the learned representation for a given instruction should still be useful even if the objects on a map are rearranged or the layout is changed entirely.

To that end, our model combines the textual instructions with the map in a spatially localized manner, as opposed to prior work which joins goal representations and environment observations via simpler functions like an inner product (Schaul et al., 2015). While our approach can more effectively learn local relations specified by language, it cannot naturally capture descriptions at the global environment level. To address this problem, we also use the language representation to predict coefficients for a basis set of gradient functions which can be combined to encode global spatial relations.

More formally, inputs to our model (see Figure 2) consist of an environment observation $s$ and textual description of a goal $x$. For simplicity, we will assume $s$ to be a 2D matrix, although the model can easily be extended to other input representations. We first convert $s$ to a 3D tensor by projecting each cell to a low-dimensional embedding ($\phi$) as a function of the objects contained in that cell. In parallel, the text instruction $x$ is passed through an LSTM recurrent neural network (Hochreiter and Schmidhuber, 1997) to obtain a continuous vector representation $h(x)$. This vector is then split into *local* and *global* components $h(x) = [h_1(x); h_2(x)]$. The local component, $h_2(x)$, is reshaped into a kernel to perform a convolution operation on the state embedding $\phi(s)$ (similar to Chen et al. (2015)):

$$z_1 = \psi_1(\phi(s); h_2(x)) \qquad (3)$$

Meanwhile, the three-element global component

**Algorithm 1** Training Procedure
___
 1: Initialize experience memory $\mathcal{D}$
 2: Initialize model parameters $\Theta$
 3: **for** $epoch=1,M$ **do**
 4:     Sample instruction $x \in X$ and associated environment $E$
 5:     Predict value map $\hat{V}(s,x;\Theta)$ for all $s \in E$
 6:     Choose start state $s_0$ randomly
 7:     **for** $t=1,N$ **do**
 8:         Select $a_t=\mathrm{argmax}_a \sum_s T(s|s_{t-1},a)\hat{V}(s,x;\Theta)$
 9:         Observe next state $s_t$ and reward $r_t$
10:     Store trajectory $(\boldsymbol{s}=s_0,s_1,...,\boldsymbol{r}=r_0,r_1,...)$ in $\mathcal{D}$
11:     **for** $j=1,J$ **do**
12:         Sample random trajectory $(\boldsymbol{s},\boldsymbol{r})$ from $\mathcal{D}$
13:         Perform gradient descent step on loss $\mathcal{L}(\theta)$
___

$h_1(x)$ is used to form the coefficients for a vertical and horizontal gradient along with a corresponding bias term.[4] The gradients, denoted $G_1$ and $G_2$ in Figure 2, are matrices of the same dimensionality as the state observation with values increasing down the rows and along the columns, respectively. The axis-aligned gradients are weighted by the elements of $h_1(x)$ and summed to give a final global gradient spanning the entire 2D space, analogous to how steerable filters can be constructed for any orientation using a small set of basis filters (Freeman and Adelson, 1991):

$$z_2 = h_1(x)[1] \cdot G_1 + h_1(x)[2] \cdot G_2 + h_1(x)[3] \cdot J \quad (4)$$

in which $J$ is the all-ones matrix also of the same dimensionality as the observed map.

Finally, the local and global information maps are concatenated into a single tensor, which is then processed by a convolutional neural network (CNN) with parameters $\theta$ to approximate the generalized value function:

$$\hat{V}(s,x) = \psi_2([z_1;z_2];\theta) \quad (5)$$

for every state $s$ in the map.

**Reinforcement Learning** Given our model's $\hat{V}(s,x)$ predictions, the resulting policy (Equation 2) can be enacted, giving a continuous trajectory of states $\{s_t, s_{t+1}, \ldots\}$ on a single map and their associated rewards $\{r_t, r_{t+1}, \ldots\}$ at each timestep

___
[4]Note that we are referring to gradient filters here, not the gradient calculated during backpropagation in deep learning.

$t$. We stored entire trajectories (as opposed to state transition pairs) in a replay memory $\mathcal{D}$ as described in Mnih et al. (2015). The model is trained to produce an accurate value estimate by minimizing the following objective:

$$\mathcal{L}(\Theta) = \mathrm{E}_{s \sim \mathcal{D}} \left[ \hat{V}(s,x;\Theta) - \left( R(s,x) + \gamma \max_a \sum_{s'} T(s'|s,a)\hat{V}(s',x;\Theta^-) \right) \right]^2 \quad (6)$$

where $s$ is a state sampled from $\mathcal{D}$, $\gamma$ is the discount factor, $\Theta$ is the set of parameters of the entire model, and $\Theta^-$ is the set of parameters of a target network copied periodically from our model. The complete training procedure is shown in Algorithm 1.

## 5 Experimental Setup

**Puddle world navigation data** In order to study generalization across a wide variety of environmental conditions and linguistic inputs, we develop an extension of the *puddle world* reinforcement learning benchmark (Sutton, 1996; Mankowitz et al., 2016). States in a $10 \times 10$ grid are first filled with either grass or water cells, such that the grass forms one connected component. We then populate the grass region with six unique objects which appear only once per map (triangle, star, diamond, circle, heart, and spade) and four non-unique objects (rock, tree, horse, and house) which can appear any number of times on a given map. See Figure 1 for an example visualization.

| Split | Local | Global |
|-------|-------|--------|
| Train | 1566  | 1071   |
| Test  | 399   | 272    |

Table 1: Overall statistics of our dataset.

Goal positions are chosen uniformly at random from the set of grass cells, encouraging the use of spatial references to describe goal locations which do not themselves contain a unique object. We used the Mechanical Turk crowdsourcing platform (Buhrmester et al., 2011) to collect natural language descriptions of these goals. Human annotators were asked to describe the positions of these

- *Reach the horse below the rock and to the left of the green diamond*

- *Move to the square two below and one left of the star*

- *Go to the cell above the bottommost horse*

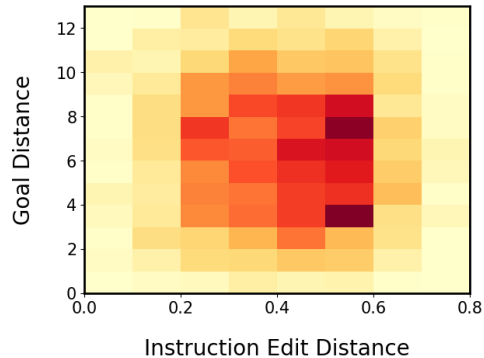Figure 3: Example goal annotations collected with Mechanical Turk.



Figure 4: A heatmap showing the normalized instruction edit distance and goal Manhattan distance corresponding to the most similar instructions between the train and test set. For each instruction in the test set, we find the five most similar instructions in the training set. Even for those instructions which are similar, the goal locations they describe can be far apart.

goals using surrounding objects. At the end of each trial, we asked the same participants to provide goal locations given their own text instructions. This helped filter out a majority of instructions that were ambiguous or ill-specified. Table 1 provides some statistics on the data, and Figure 3 shows example annotations. In total, we collected 3308 instructions, ranging from 2 to 43 words in length, describing over 200 maps. There are 361 unique words in the annotated instructions. We do not perform any pre-processing on the raw annotations.

It is plausible that a model designed to handle only local references could not handle global ones (consider our own model without the global gradient maps). For clearer interpretation of results, we evaluate our model in two modes: trained and tested on local and global data separately, or as a combined dataset. While local instructions were obtained easily, the global instructions were collected by designing a task in which only nonunique objects were presented to the annotators.[5] This precluded simple instructions like *"go left of the ⟨object⟩"* because there would always be more than one of each object type. Therefore, we obtained text with global properties (e.g. *middle rock*, *leftmost tree*) to sufficiently pinpoint an object. On average, we collected 31 unique local instructions and 10 unique global instructions per map.

To quantify the diversity of our dataset, we find the five nearest instructions in the training set for every instruction in the test set, as measured by edit distance (using the word as a unit) normalized by test instruction length. For each of these pairs, we also measure the Manhattan distance between their corresponding goal locations. Figure 4, which visu-

alizes this analysis, underscores the difficulty of this task; even when two instructions are highly similar, they might correspond to entirely different target locations. This is the case in the example in Figure 1, which has a distance of four between the references goals.

**Baselines** We compare our model to the following baselines:

**UVFA (text)** is a variant of the model described in (Schaul et al., 2015) adapted for our task. The original model made use of two MLPs to learn low dimensional embeddings of states and goals which were then combined via dot product to give value estimates. Goals were represented either as $(x, y)$ coordinates or as states themselves. As our goals are not observed directly but described in text, we replace the goal MLP with the same LSTM as in our model. The state MLP has an identical architecture to that of the UVFA: two hidden layers of dimension 128 and ReLU activations. For consistency with the UVFA, we represent states as binary vectors denoting the presence of each type of object at every position.

**CNN + LSTM** is a variant of the model described in Misra et al. (2017), who developed it for a language-grounded block manipulation task. It first convolves the map layout to a low-dimensional rep-

---

[5]The other objects were added back into the map after collecting the instruction.
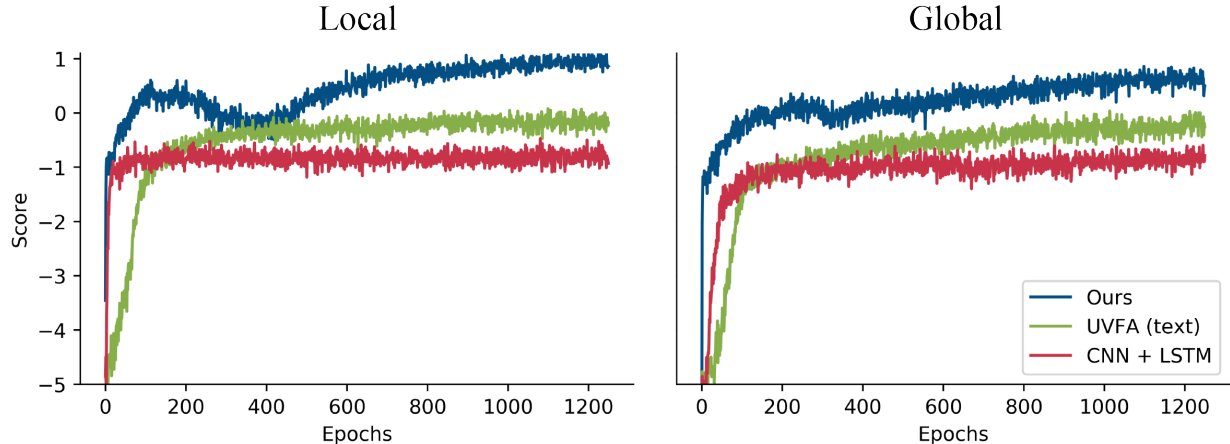
Figure 5: Reward achieved by our model and the two baselines on the training environments during reinforcement learning on both local and global instructions. Each epoch corresponds to simulation on 500 goals, with a goal simulation terminating either when the agent reaches the goal state or has taken 75 actions.

resentation (as opposed to the MLP of the UVFA) and concatenates this to the LSTM's instruction embedding (as opposed to a dot product). These concatenated representations are then input to a two-layer MLP.

We also perform analysis to study the representational power of our model, introducing two more comparison models:

**UVFA (pos)** is the original UVFA model from (Schaul et al., 2015), which we evaluate on our modified puddle worlds to determine the difficulty of environment generalization independently from instruction interpretation.

**Our model (w/o gradient)** is an ablation of our model without the global gradient maps, which allows us to determine the gradients' role in representation-building.

In additional to our reinforcement learning experiments, we train these models in a supervised setting to isolate the effects of architecture choices from other concerns inherent to reinforcement learning algorithms. For this purpose, we constructed a dataset of ground-truth value maps for all human-annotated goals using value iteration. We use the models to predict value maps for the entire grid and minimize the mean squared error (MSE) compared to the ground truth values:

$$\mathcal{L}'(\Theta) = \sum_{(s,x)} [\hat{V}(s,x;\Theta) - \bar{V}(s,x)]^2 \quad (7)$$

### 5.1 Implementation details

Our model implementation uses an LSTM with a learnable 15-dimensional embedding layer, 30 hidden units, 8-dimensional embeddings $\phi(s)$, and a 3x3 kernel applied to the embeddings, giving a dimension of 72 for $h_2(t)$. The final CNN has layers of $\{3, 6, 12, 6, 3, 1\}$ channels, all with 3x3 kernels and padding of length 1 such that the output value map prediction is equal in size to the input observation. For each map, a reward of 3 is given for reaching the correct goal specified by human annotation and a reward of $-1$ is given for falling in a puddle cell. The only terminal state is when the agent is at the goal. Rewards are discounted by a factor of 0.95. We use Adam optimization (Kingma and Ba, 2015) for training all models.

## 6 Results

We present empirical results on two different datasets - our annotated puddle world and an existing block navigation task (Bisk et al., 2016).

### 6.1 Puddle world navigation

**Comparison with the state-of-the-art** We first investigate the ability of our model to learn solely from environment simulation. Figure 5 shows the discounted reward achieved by our model as well as the two baselines for both instruction types. In both experiments, our model is the only one of the

| | Local | | Global | | Combined | |
|---|---|---|---|---|---|---|
| | Policy Quality | Distance | Policy Quality | Distance | Policy Quality | Distance |
| UVFA (text) | 0.56 | 4.71 | 0.62 | 6.28 | 0.61 | 5.06 |
| CNN + LSTM | 0.80 | 5.73 | 0.82 | 6.13 | 0.82 | 5.67 |
| Our model | **0.87** | **2.18** | **0.90** | **3.35** | **0.89** | **3.04** |

Table 2: Performance of models trained via reinforcement learning on a held-out set of environments and instructions. Policy quality is the true expected normalized reward and distance denotes the Manhattan distance from goal location prediction to true goal position. We show results from training on the local and global instructions both separately and jointly.

| | Local | | | Global | | |
|---|---|---|---|---|---|---|
| | MSE | Policy Quality | Distance | MSE | Policy Quality | Distance |
| UVFA (pos) | 1.30 | 0.48 | 4.96 | 1.04 | 0.52 | 5.39 |
| UVFA (text) | 3.23 | 0.57 | 4.97 | 1.9 | 0.62 | 5.31 |
| CNN + LSTM | 0.42 | 0.86 | 4.08 | 0.43 | 0.83 | 4.18 |
| Our model (w/o gradient) | **0.25** | **0.94** | 2.39 | 0.61 | 0.87 | 5.15 |
| Our model | **0.25** | **0.94** | **2.34** | **0.41** | **0.89** | **3.81** |

Table 3: Performance on a test set of environments and instructions after supervised training. Lower is better for MSE and Manhattan distance; higher is better for policy quality. The gradient basis significantly improves the reconstruction error and goal localization of our model on global instructions, and expectedly does not affect its performance on local instructions.

three to achieve an average nonnegative reward after convergence (0.88 for local instructions and 0.49 for global instructions), signifying that the baselines do not fully learn how to navigate through these environments.

Following Schaul et al. (2015), we also evaluated our model using the metric of *policy quality*. This is defined as the expected discounted reward achieved by following a softmax policy of the value predictions. Policy quality is normalized such that an optimal policy has a score of 1 and a uniform random policy has a score of 0. Intuitively, policy quality is the true normalized expectation of score over all maps in the dataset, instructions per map, and start states per map-instruction pair. Our model outperforms both baselines on this metric as well on the test maps (Table 2). We also note that the performance of the baselines flip with respect to each other as compared to their performance on the training maps (Figure 5). While the UVFA variant learned a better policy on the train set, it did not generalize to new environments as well as the CNN + LSTM.

Finally, given the nature of our environments, we can use the predicted value maps to infer a goal location by taking the position of the maximum value. We use the Manhattan distance from this predicted position to the actual goal location as a third metric. The accuracy of our model's goal predictions is more than twice that of the baselines on local references and roughly $45\%$ better on global references.

**Analysis of learned representations** For the representation analysis in a supervised setting, we compared the predicted value maps of all models against
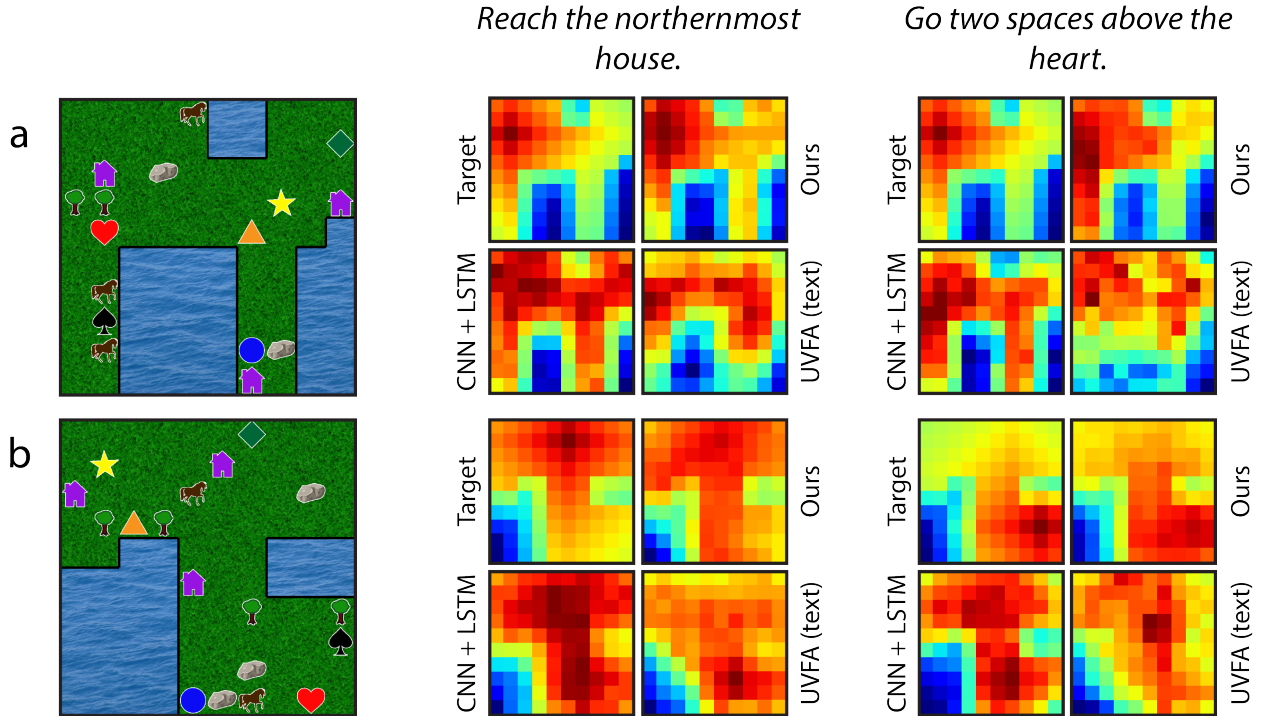
Figure 6: Value map predictions for two environments paired with two instructions each. Despite the difference in instructions, with one being global and the other local in nature and sharing no objects in their descriptions, they refer to the same goal location in the environment in (a). However, in (b), the descriptions correspond to different locations on the map. The vertical axis considers variance in goal location for the same instruction, depending on the map configuration.

the unseen test split of maps. Table 3 shows the results of this study. As expected, our model without the global gradient performs no differently from the full model on local references, but has higher MSE and average distances to true goal than the full model on global references. We also note that UVFA (pos) performs much worse than both CNN+LSTM and our model, showing the difficulty of environment generalization even when the goals are observed directly. (The original UVFA paper (Schaul et al., 2015) demonstrated effective generalization over goal states within a *single* environment.)

Surprisingly, our model trained via reinforcement learning has more precise goal location predictions (as measured via Manhattan distance) than when trained on true state values in a supervised manner. However, the MSE of the value predictions are much higher in the RL setting (*e.g.*, 0.80 vs 0.25 for supervised on local instructions). This shows that despite the comparative stability of the supervised setting, minimization of value prediction error does not nec-

essarily lend itself to the best policy or goal localization. Conversely, having a higher MSE does not always imply a worse policy, as seen also in the performance of the two UVFA variants in Table 3.

**Generalization** One of the criteria laid out for our model was its ability to construct language representations and produce accurate value maps, independent of layouts and linguistic variation. Figure 6 provides examples of two layouts, each with two different instructions. In the first map (top), we have both instructions referring to the same location. Our model is able to mimic the optimal value map accurately, while the other baselines are not as precise, either producing a large field of possible goal locations (*CNN+LSTM*) or completely missing the goal (*UVFA-text*).

On the vertical axis, we observe generalization across different maps with the same instructions. Our model is able to precisely identify the goals in each scenario in spite of significant variation in their

*reach cell one to the right and two down from horse located to the right of the heart*



*locate the cell that is filled w/a rock to the left of the teal spade and above the purple house which is above and to the left of the blue circle*
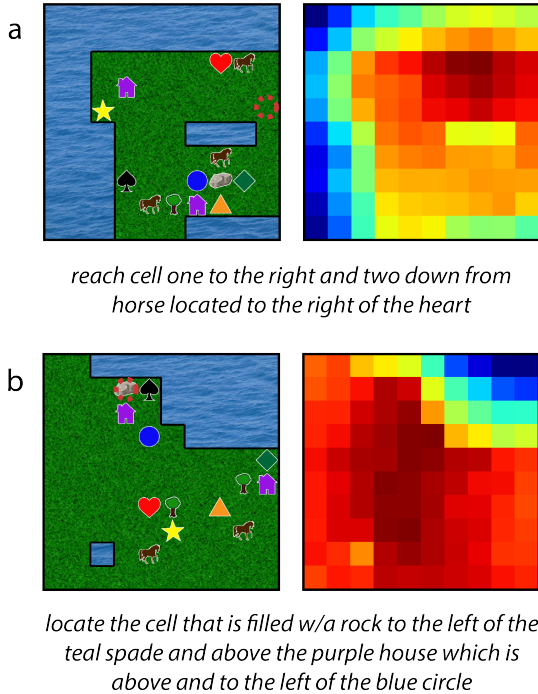
Figure 7: Examples of failure cases for our model. Multiple levels of indirection in (a) and a long instruction filled with redundant information in (b) make the instruction difficult to interpret. Intended goal locations are outlined in red for clarity.

locations. This proves harder for the other representations.

Although our model is compositional in the sense that it transfers knowledge of spatial references between different environments, some types of instructions do prove challenging. We identify two of the poorest predictions in Figure 7. We see that multiple levels of indirection (as in 7a, which references a location relative to an object relative to another object) or unnecessarily long instructions (as in 7b, which uniquely identifies a position by the eighth token but then proceeds with redundant information) are still a challenge.

**Learning curve** Due to the manual effort that comes with constructing a dataset of human annotations, it is also important to consider the sample-efficiency of a model. Figure 8 shows the quality policy and prediction error on local instructions as a function of training set size. We observe that our model reaches 0.90 policy quality with only 400
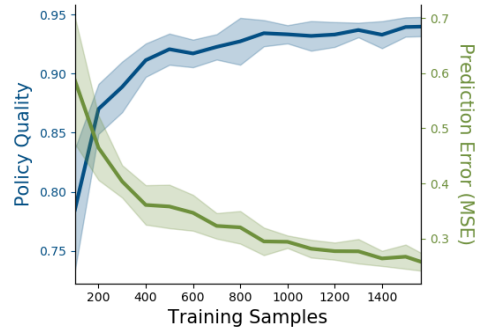


Figure 8: Effect of training set size on held-out predictions. The curves show the mean of ten training runs and the shaded regions show standard deviation. Our model's policy quality is greater than 0.90 with as few as 400 training goal annotations.

samples, demonstrating efficient generalization capability.

| | Language Grounding Dataset | |
|---|---|---|
| | Policy Quality | Distance |
| UVFA (text) | 0.15 | 4.61 |
| CNN + LSTM | 0.17 | 4.11 |
| Our model | **0.74** | **3.94** |

Table 4: The performance of our model and two baselines on the ISI Language Grounding dataset (Bisk et al., 2016). Our model once again outperforms the baselines, although all models have a lower policy quality on this dataset than on our own.

### 6.2 ISI Grounding Dataset

We also evaluate our model on the ISI Language Grounding dataset (Bisk et al., 2016), which contains human-annotated instructions describing how to arrange blocks identified by numbers and logos. Although it does not contain variable environment maps as in our dataset, it has a larger action space and vocabulary. The caveat is that the task as posed in the original dataset is not compatible with our model. For a policy to be derived from a value map with the same dimension as the state observation, it is implicitly assumed that there is a single con-
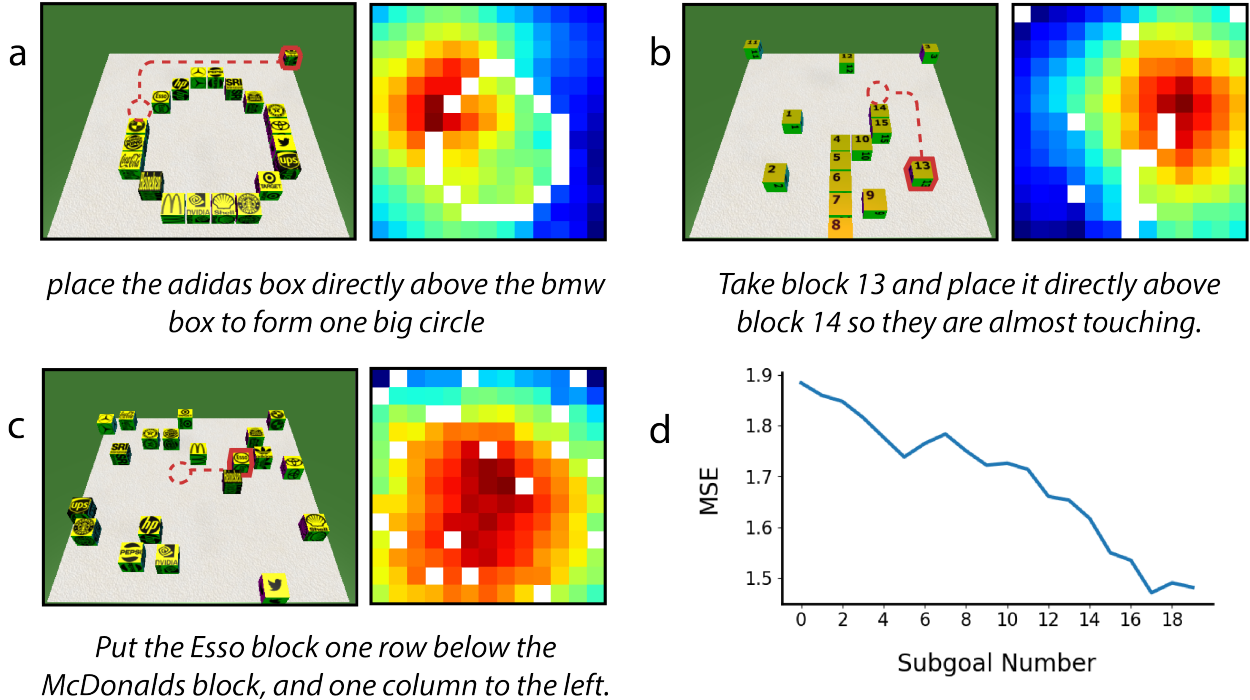
*place the adidas box directly above the bmw box to form one big circle*

*Take block 13 and place it directly above block 14 so they are almost touching.*

*Put the Esso block one row below the McDonalds block, and one column to the left.*

Figure 9: (a-c) Visualizations of tasks from the ISI Language Grounding dataset (Bisk et al., 2016) and our model's value map predictions. The agentive block and goal location are outlined in red for visibility. (d) The MSE of the value map prediction as a function of a subgoal's ordering in an overall task. The model performs better on subgoals later in a task despite the subgoals being treated completely independently during both training and testing.

trollable agent, whereas the ISI set allows multiple blocks to be moved. We therefore modify the ISI setup using an oracle to determine which block is given agency during each step. This allows us to retain the linguistic variability of the dataset while overcoming the mismatch in task setup. The states are discretized to a $13 \times 13$ map and the instructions are lemmatized.

Performance on the modified ISI dataset is reported in Table 4 and representative visualizations are shown in Figure 9. Our model outperforms both baselines by a greater margin in policy quality than on our own dataset.

Misra et al. (2017) also use this dataset and report results in part by determining the minimum distance between an agent and a goal during an evaluation lasting $N$ steps. This evaluation metric is therefore dependent on this timeout parameter $N$. Because we discretized the state space so as to be able to represent it as a grid of embeddings, the notion of a single step has been changed and direct comparison limited

to $N$ steps is ill-defined.[6] Hence, due to modifications in the task setup, we cannot compare directly to the results in Misra et al. (2017).

**Understanding grounding evaluation** An interesting finding in our analysis was that the difficulty of the language interpretation task is a function of the stage in task execution (Figure 9(d)). In the ISI Language Grounding set (Bisk et al., 2016), each individual instruction (describing where to move a particular block) is a subgoal in a larger task (such as constructing a circle with all of the blocks). The value maps predicted for subgoals occurring later in a task are more accurate than those occurring early in the task. It is likely that the language plays a less crucial role in specifying the subgoal position in the

---

[6]When a model is available and the states are not overwhelmingly high-dimensional, policy quality is a useful metric that is independent of this type of parameter. As such, it is our default metric here. However, estimating policy quality for environments substantially larger than those investigated here is a challenge in itself.

final steps of a task. As shown in Figure 9(a), it may be possible to narrow down candidate subgoal positions just by looking at a nearly-constructed high-level shape. In contrast, this would not be possible early in a task because most of the blocks will be randomly positioned. This finding is consistent with a result from Branavan et al. (2011), who reported that strategy game manuals were useful early in the game but became less essential further into play. It appears to be part of a larger trend that the marginal benefit of language in such grounding tasks can vary predictably between individual instructions.

# 7 Conclusions

We have described a novel approach for grounded spatial reasoning. Combining the language representation in a spatially localized manner allows for increased precision of goal identification and improved performance on unseen environment configurations. Alongside our models, we present *Puddle World Navigation*, a new grounding dataset for testing the generalization capacity of instruction-following algorithms in varied environments.

## Acknowledgement

## References

Jacob Andreas and Dan Klein. 2015. Alignment-based compositional semantics for instruction following. In *EMNLP*.

Yoav Artzi and Luke Zettlemoyer. 2013. Weakly supervised learning of semantic parsers for mapping instructions to actions. *TACL*, 1(1):49–62.

Richard Bellman. 1957. *Dynamic Programming*. Princeton University Press, Princeton, NJ, USA, 1 edition.

Yonatan Bisk, Deniz Yuret, and Daniel Marcu. 2016. Natural language communication with robots. In *NAACL HLT*, pages 751–761.

SRK Branavan, David Silver, and Regina Barzilay. 2011. Learning to win by reading manuals in a monte-carlo framework. In *ACL HLT*, pages 268–277.

Michael Buhrmester, Tracy Kwang, and Samuel D. Gosling. 2011. Amazon's mechanical turk. *Perspectives on Psychological Science*, 6(1):3–5. PMID: 26162106.

Ruth MJ Byrne and Philip N Johnson-Laird. 1989. Spatial reasoning. *Journal of memory and language*, 28(5):564–575.

David L Chen and Raymond J Mooney. 2011. Learning to interpret natural language navigation instructions from observations. *San Francisco, CA*, pages 859–865.

Kan Chen, Jiang Wang, Liang-Chieh Chen, Haoyuan Gao, Wei Xu, and Ram Nevatia. 2015. Abc-cnn: An attention based convolutional neural network for visual question answering. *arXiv preprint arXiv:1511.05960*.

William T. Freeman and Edward H. Adelson. 1991. The design and use of steerable filters. *IEEE TPAMI*, 13(9):891–906, September.

Sachithra Hemachandra, Matthew R Walter, Stefanie Tellex, and Seth Teller. 2014. Learning spatial-semantic representations from natural language descriptions and scene classifications. In *ICRA*, pages 2623–2630. IEEE.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation*, 9(8):1735–1780.

John D Kelleher, Geert-Jan M Kruijff, and Fintan J Costello. 2006. Proximity in context: an empirically grounded computational model of proximity for processing topological spatial expressions. In *ACL*, pages 745–752.

Joohyun Kim and Raymond J Mooney. 2013. Adapting discriminative reranking to grounded language learning. In *ACL (1)*, pages 218–227.

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *ICLR*.

Thomas Kollar, Stefanie Tellex, Deb Roy, and Nicholas Roy. 2010. Toward understanding natural language directions. In *Human-Robot Interaction*, pages 259–266. IEEE.

Stephen C Levinson. 2003. *Space in language and cognition: Explorations in cognitive diversity*, volume 5. Cambridge University Press.

Michael Levit and Deb Roy. 2007. Interpretation of spatial language in a map navigation task. *IEEE Transactions on Systems, Man, and Cybernetics, Part B (Cybernetics)*, 37(3):667–679.

Peggy Li and Lila Gleitman. 2002. Turning the tables: Language and spatial reasoning. *Cognition*, 83(3):265–294.

Matt MacMahon, Brian Stankiewicz, and Benjamin Kuipers. 2006. Walk the talk: Connecting language, knowledge, and action in route instructions. *AAAI*, 2(6):4.

Daniel J. Mankowitz, Timothy Arthur Mann, and Shie Mannor. 2016. Iterative hierarchical optimization for misspecified problems (IHOMP). *CoRR*, abs/1602.03348.

Dipendra K Misra, John Langford, and Yoav Artzi. 2017. Mapping instructions and visual observations to actions with reinforcement learning. *EMNLP*.

Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Andrei A. Rusu, Joel Veness, Marc G. Bellemare, Alex Graves, Martin Riedmiller, Andreas K. Fidjeland, Georg Ostrovski, Stig Petersen, Charles Beattie, Amir Sadik, Ioannis Antonoglou, Helen King, Dharshan Kumaran, Daan Wierstra, Shane Legg, and Demis Hassabis. 2015. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 02.

Reinhard Moratz and Thora Tenbrink. 2006. Spatial reference in linguistic human-robot interaction: Iterative, empirically supported development of a model of projective relations. *Spatial cognition and computation*, 6(1):63–107.

Tom Schaul, Daniel Horgan, Karol Gregor, and David Silver. 2015. Universal value function approximators. In *ICML*, pages 1312–1320.

M. Skubic, D. Perzanowski, S. Blisard, A. Schultz, W. Adams, M. Bugajska, and D. Brock. 2004. Spatial language for human-robot dialogs. *IEEE Transactions on Systems, Man, and Cybernetics, Part C (Applications and Reviews)*, 34(2):154–167, May.

Richard S Sutton and Andrew G Barto. 1998. *Introduction to reinforcement learning*. MIT Press.

Richard S. Sutton. 1996. Generalization in reinforcement learning: Successful examples using sparse coarse coding. In *NIPS*.

Stefanie Tellex, Thomas Kollar, Steven Dickerson, Matthew R Walter, Ashis Gopal Banerjee, Seth J Teller, and Nicholas Roy. 2011. Understanding natural language commands for robotic navigation and mobile manipulation. In *AAAI*.

Thora Tenbrink. 2007. *Space, time, and the use of language: An investigation of relationships*, volume 36. Walter de Gruyter.

Jette Viethen and Robert Dale. 2008. The use of spatial relations in referring expression generation. In *Proceedings of the Fifth International Natural Language Generation Conference*, pages 59–67. Association for Computational Linguistics.

Adam Vogel and Dan Jurafsky. 2010. Learning to follow navigational directions. In *ACL*, pages 806–814.

Matthew R Walter, Sachithra Hemachandra, Bianca Homberg, Stefanie Tellex, and Seth Teller. 2013. Learning semantic maps from natural language descriptions. Robotics: Science and Systems.