

CraftAssist Instruction Parsing: Semantic Parsing for a Minecraft Assistant

Yacine Jernite*

Kavya Srinet*

Jonathan Gray

Arthur Szlam

Facebook AI Research

Abstract

We propose a large scale semantic parsing dataset focused on instruction-driven communication with an agent in Minecraft. We describe the data collection process which yields additional 35K human generated instructions with their semantic annotations. We report the performance of three baseline models and find that while a dataset of this size helps us train a usable instruction parser, it still poses interesting generalization challenges which we hope will help develop better and more robust models.

1 Introduction

Semantic parsing is used as a component for natural language understanding in human-robot interaction systems (Tellex et al., 2011; Matuszek et al., 2013), and for virtual assistants (Kollar et al., 2018). Recently, researchers have shown success with deep learning methods for semantic parsing, e.g. (Dong and Lapata, 2016; Jia and Liang, 2016; Zhong et al., 2017). However, to fully utilize powerful neural network approaches, it is necessary to have large numbers of training examples. In the space of human-robot (or human-assistant) interaction, the publicly available semantic parsing datasets are small. Furthermore, it can be difficult to reproduce the end-to-end results (from utterance to action) because of the wide variety of robot setups and proprietary nature of personal assistants.

In this work, we introduce a new semantic parsing dataset for human-bot interactions. Our “robot” or “assistant” is embodied in the sandbox construction game Minecraft¹, a popular multi-player open-world voxel-based crafting game. We also provide the associated platform for executing the logical forms in game.

¹<https://minecraft.net/en-us/>. We limit ourselves to creative mode for this work

Situating the assistant in Minecraft has several benefits for studying task oriented natural language understanding (NLU). Compared to physical robots, Minecraft allows less technical overhead irrelevant to NLU, such as difficulties with hardware and large scale data collection. On the other hand, our bot has all the basic in-game capabilities of a player, including movement and placing or removing voxels. Thus Minecraft preserves many of the NLU elements of physical robots, such as discussions of navigation and spatial object reference.

Furthermore, working in Minecraft may enable large scale human interaction because of its large player base, in the tens of millions. Although Minecraft’s simulation of physics is simplified, the task space is complex. There are many atomic objects in Minecraft, such as animals and block-types, that require no perceptual modeling. For researchers interested in the interactions between perception and language, collections of voxels making up a “house” or a “hill” are not atomic objects and the assistant cannot apprehend them without a perceptual system.

Our contributions in the paper are as follows:

Grammar: We develop a set of action primitives and grammar over these primitives that comprise a mid-level interface to Minecraft, for machine learning agents. See Section 3.

Data: Using a collection of language templates to convert logical forms over the primitives into pseudo-natural language, we build a dataset of language instructions with logical form annotation by having crowd-sourced workers rephrase the language outputs, as in (Wang et al., 2015). We also collect a test set of crowd-sourced annotations of commands generated independent of our grammar. In addition to the natural language commands and the associated logical forms, we also make available the code to execute these in the

game, allowing the reproduction of end-to-end results. See Section 4.

Models: We show the results of several neural semantic parsing models trained on our data. See Section 5 and 6. We also will provide access to an interactive bot using these models for parsing².

2 Related Work

There have been a number of datasets of natural language paired with logical forms to evaluate semantic parsing approaches, e.g. (Price, 1990; Tang and Mooney, 2001; Cai and Yates, 2013; Wang et al., 2015; Zhong et al., 2017). The dataset presented in this work is an order of magnitude larger than those in (Price, 1990; Tang and Mooney, 2001; Cai and Yates, 2013) and is similar in scale to (Wang et al., 2015; Zhong et al., 2017). We use the data collection strategy in (Wang et al., 2015) to build the pairings between logical forms and natural language: first building the grammar, then generating from the grammar via templates, and then using crowd-sourced workers to rephrase the templated generations. However, we also collect a test set of “free” commands and use crowd-sourced workers to annotate these.

In addition to connecting natural language to logical forms, our dataset connects both of these to a dynamic environment. In (Tellex et al., 2011; Matuszek et al., 2013) semantic parsing has been used for interpreting natural language commands for robots. In our paper, the “robot” is embodied in the Minecraft game instead of in the physical world.

Semantic parsing in a voxel-world recalls (Wang et al., 2017), where the authors describe a method for building up a programming language from a small core via interactions with players.

We demonstrate the results of several neural parsing models on our dataset. In particular, we show the results of a reimplementation of (Dong and Lapata, 2016) adapted to our grammar. There have been several other papers proposing neural architectures for semantic parsing, for example (Jia and Liang, 2016; Zhong et al., 2017). In those papers, as in this one, the models are trained with full supervision of the mapping from natural language to logical forms, without considering the results of executing the logical form (in this case, the effect on the environment of executing the ac-

tions denoted by the logical form). There has been progress towards “weakly supervised” semantic parsing (Artzi and Zettlemoyer, 2013; Liang et al., 2016; Guu et al., 2017) where the logical forms are hidden variables, and the only supervision given is the result of executing the logical form. There are now approaches that have shown promise without even passing through (discrete) logical forms at all (Riedel et al., 2016; Neelakantan et al., 2016). We hope that the dataset introduced here, which has supervision at the level of the logical forms, but whose underlying grammar and environment can be used to generate essentially infinite weakly supervised or execution rewards, will also be useful for studying these models.

Minecraft, especially via the MALMO project (Johnson et al., 2016) has been used as a base environment for several machine learning papers. Often Minecraft is used as a testbed for reinforcement learning (Shu et al., 2017; Udagawa et al., 2016; Alaniz, 2018; Oh et al., 2016; Tessler et al., 2017). In these papers, the agent is trained to complete tasks by issuing low level actions (as opposed to our higher level primitives) and receiving a reward on success. Some of these papers(e.g. (Oh et al., 2017)) do consider simplified, templated language as a method for composable specifying tasks, but training an RL agent to execute the scripted primitives in our grammar is already nontrivial, and so the task space and language is more constrained than what we use here. Nevertheless, our work may be useful to researchers interested in RL- using our grammar and executing in game can supply (hard) tasks and descriptions. Another set of papers (Kitaev and Klein, 2017; Yi et al., 2018) have used Minecraft for visual question answering with logical forms. Our work extends these to interactions with the environment. Finally, (Allison et al., 2018) is a more focused study on how a human might interact with a Minecraft agent; our collection of free generations (see 4.2.2) includes annotated examples from similar studies of players interacting with a player pretending to be a bot.

3 A Natural Language Interface

We want to interpret natural language commands given to an agent with a pre-defined set of capabilities. We start by providing an overview of these capabilities and the action space that they entail, then define a grammar to capture this action space.

²Instructions will be available at <http://craftassist.io/acl2019demo>

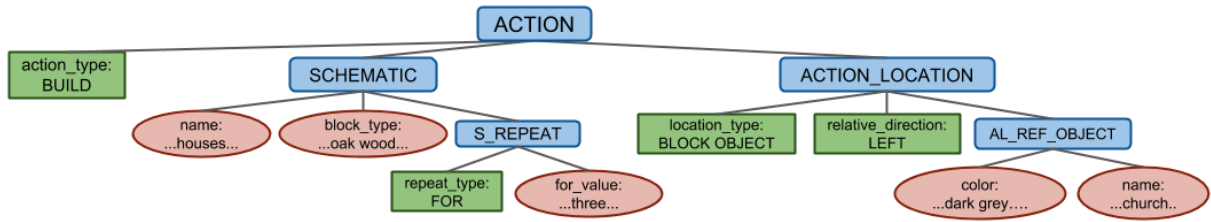


Figure 1: Parse tree for “Make three oak wood houses to the left of the dark grey church.”

3.1 Agent Action Space

The goal of the proposed agent is to help a player create structures and mechanisms in a voxelized world by moving around, and placing and removing blocks. To this end, the agent needs to be able to understand a number of high-level commands, which we present here.

Basic action commands First, we need commands corresponding to high level actions of the agent. For example, we may ask it to BUILD an object from a known schematic or to copy an existing structure at a given location, or to DESTROY one. Similarly, it might be useful to be able to ask the agent to DIG a hole of a given shape at a specified location, or on the contrary to FILL one up. The agent can also be asked to complete an already structure however it sees fit (this action is called FREEBUILD), or to SPAWN game mobs. Finally, we need to be able to direct the agent to MOVE to a location.

Teaching and querying the bot In order to understand most of the above commands, the agent needs to have an internal representation of the world. We want to be able to add to this representation by allowing the user to TAG existing objects with names or properties. This can be considered a basic version of the self-improvement capabilities in (Kollar et al., 2013; Thomason et al., 2015; Wang et al., 2016, 2017). Conversely, to query this internal state, we can ask the agent to ANSWER questions about the world. This part of the grammar is similar to the visual question-answering in (Yi et al., 2018)

Control commands Additionally, we want to be able to ask the agent to STOP or RESUME an action, or to UNDO the result of a recent command. Finally, the agent needs to be able to understand when a sentence does not correspond to any of the above mentioned actions, and map it to a NOOP command.

3.2 Parsing Grammar

All of the above commands are represented as trees encoding all of the information necessary for their execution. Figure 1 presents an example parse tree for the BUILD command “*Make three oak wood houses to the left of the dark grey church.*”

Internal nodes Each action has a set of possible arguments, which themselves have a recursive argument structure. Each of these action types and complex arguments corresponds to an internal node (blue rounded rectangles in Figure 1), with its children providing more specific information. For example, the BUILD action can specify a SCHEMATIC (what we want to build) and a LOCATION child (where we want to build it). In turn, the SCHEMATIC can specify a general category (house, bridge, temple, etc...), as well as a set of properties (size, color, building material, etc...), and in our case also has a REPEAT child subtree specifying how many we want to build. Similarly, the LOCATION can specify an absolute location, a distance, direction, and information about the location REFERENCE OBJECT stored in a child subtree.

One notable feature of this representation is that we do not know *a priori* which of a node’s possible children will be specified. For example, BUILD can have a SCHEMATIC and a LOCATION specified (“*Build a house over there.*”), just a SCHEMATIC (“*Build a house.*”), just a LOCATION (“*Build something next to the bridge.*”), or neither (“*Make something.*”).

The full grammar is specified in Figure 3. In addition to the various LOCATION, REFERENCE OBJECT, SCHEMATIC, and REPEAT nodes which can be found at various levels, another notable subtree is the action’s STOP CONDITION, which essentially allows the agent to understand “while” loops (for example: “dig down until you hit the bedrock” or “follow me”).

Leaf nodes Eventually, arguments have to be specified in terms of values which correspond to agent primitives. We call these nodes categorical leaves (green rectangles in Figure 1). The root of the tree has a categorical leaf child which specifies the **action type**, BUILD in our example. There are also nodes specifying the **repeat type** in the REPEAT sub-tree (“make three houses” corresponds to executing a FOR loop), the LOCATION TYPE (the location is given in reference to the BLOCK OBJECT that is the “dark grey church”), and the **relative direction** to the reference, here LEFT.

However, there are limits to what we can represent with a pre-specified set of hard-coded primitives, especially if we want our agent to be able to learn new concepts or new values. Additionally, even when there is a pre-specified agent primitive, mapping some parts of the command to a specific value might be better left to an external module (e.g. mapping a number string to an integer value). For both of these reasons, we also have span leaves (red ovals in Figure 1). This way, a model can learn to generalize to e.g. colors or size descriptions that it has never seen before. The SCHEMATIC is specified by the command substring corresponding to its **name** (“houses”) and the requested **block type** (“oak wood”). The range of the for loop is specified by the REPEAT’s **for value** (“three”), and the REFERENCE OBJECT is denoted in the command by its generic **name** and specific **color** (“church” and “dark grey”).

4 The CAIP Dataset

This paper introduces the CraftAssist Instruction Parsing (CAIP) dataset of English-language commands and their associated “action trees”, as defined in Section 3 (see Appendix A for examples and a full grammar specification). CAIP is a composite dataset containing a combination of algorithmically generated commands and human-written natural language commands.

4.1 Generated Data

We start by algorithmically generating action trees (logical forms over the grammar) with associated surface forms through the use of templates. To that end, we first define a set of template objects, which link an atomic concept in the game world to several ways it can be described through language. For example the template object MOVE links the action type MOVE to the utterances

go, walk, move,... Likewise, the template object RelativeDirection links all of the direction primitives to their names. Some template objects also have purely linguistic functions in order to make the sentence more natural but without referring to any information relevant to the tree. For example, the object ALittle can be realized into *a bit, a little, somewhat, ...*

Then, we build recursive templates for each action as recursive sequences of templates and template objects. For each of these templates, we can then sample a game value and its corresponding string. By concatenating these, we obtain an action tree and its corresponding language description. Consider for example the template [Move, ALittle, RelativeDirection] made up of the template objects described above. One possible realization could be the description *go a little to the left* paired with an action tree specifying the action type as MOVE, and an ACTION LOCATION sub-tree with which a child relative direction categorical node which has value LEFT. Finally, in addition to the action-specific templates, we also generate training data for the NOOP action type by sampling dialogue lines from the Cornell Movie Dataset (Danescu-Niculescu-Mizil and Lee, 2011).

We wrote 3,900 templates in total. We can create a training example for a parsing model by choosing one of them at random, and then sampling a (description, tree) pair from it, which, given the variety and modularity of the template objects, yields virtually unlimited data (for practical reasons, we pre-generate a set of 800K training, 5K validation, and 5K test examples for our experiments). The complete list of templates and template objects is included in the Supplementary Material.

4.2 Collected Data

To supplement the generated data, natural language commands written by crowd-sourced workers were collected in a variety of settings.

4.2.1 Rephrases

While the template generations yield a great variety of language, they cannot cover all possible ways of phrasing a specific instruction. In order to supplement them, we asked crowd-sourced workers to rephrase some of the produced instructions into commands in alternate, natural English that does not change the meaning of the sentence.

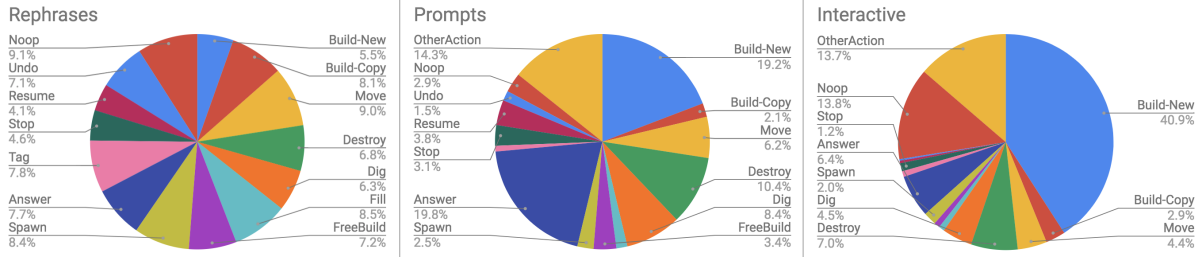


Figure 2: Frequency of each action type in the different data collection schemes described in Section 4.2. The BUILD action is divided into BUILD-NEW (a command to build a brand new structure, which may specify a SCHEMATIC) and BUILD-COPY (an command to duplicate an existing structure, which specifies a REFERENCE OBJECT).

This setup enables the collection of unique English commands whose action trees are already known. Note that a rephrased sentence will have the same action tree structure, but the positions of the words corresponding to span nodes may change. To account for this, words contained in a span range in the original sentence are highlighted in the task, and crowd-sourced workers are asked to highlight the corresponding words in their rephrased sentence. Then the action tree span values are substituted for the rephrased sentence to get the corresponding tree. This yields a total of 32K rephrases. We use 30K for training, 1K for validation, and 1K for testing.

4.2.2 Image and Text Prompts

We also presented crowd-sourced workers with a description of the capabilities of an assistant bot in a creative virtual environment (which matches the set of allowed actions in the grammar), and (optionally) some images of a bot in a game environment. They were then asked to provide examples of commands that they might issue to an in-game assistant. We refer to these instructions as “prompts” in the rest of this paper. The complete instructions shown to workers is included in appendix 19.

4.2.3 Interactive Gameplay

We asked crowd-sourced workers to play creative-mode Minecraft with an assistant bot, and they were instructed to use the in-game chat to direct the bot in whatever way they chose. The exact instructions are included in appendix B.2. Players in this setting had no prior knowledge of the bot’s capabilities or the parsing grammar.

4.2.4 Annotation Tool

Both prompts and interactive instructions come without a reference tree and need to be annotated. To facilitate this process, we designed a web-based tool which asks users a series of multiple-choice questions to determine the semantic content of a sentence. The responses to some questions will prompt other more specific questions, in a process that mirrors the hierarchical structure of the grammar. The responses are then processed to produce an action tree. This allows crowd-sourced workers to provide annotations with no knowledge of the specifics of the grammar described above. For each sentence annotated with the tool, three responses from distinct users were collected, and a sentence was included in the dataset only if at least two out of three responses matched exactly. This yields 1265 annotated prompts, and 817 annotated interactive instructions. A screenshot of the tool is included in Appendix 21.

4.3 Dataset Statistics

Action Frequencies Since the different data collection settings described in Section 4.2 imposed different constraints and biases on the crowd-sourced workers, the distribution of actions in each subset of data is therefore very different. For example, in the Interactive Gameplay scenario, workers were given no prior indication of the bot’s capabilities, and spent much of their time asking the bot to build things. The action frequencies of each subset are shown in Figure 2.

Grammar coverage Some crowd-sourced commands describe an action that is outside the scope of the grammar. To account for this, users of the action tree annotation tool are able to mark that a sentence is a command to perform an action that is not listed. The resulting action trees are

labelled OTHERACTION, and their frequency in each dataset is shown in Figure 2. Note that annotators that choose OTHERACTION still have the option to label other nodes in the action tree like LOCATION and REFERENCE OBJECT.

5 Baseline Models

In order to assess the challenges of the dataset, we implement several baseline models which read a sentence and output an Action Tree, including an adaptation of the Seq2Tree model of (Dong and Lapata, 2016) to our grammar.

Sentence Encoder All of our models rely on a sentence encoder. In this work, we use a bidirectional GRU encoder (Cho et al., 2014) which encodes a sentence of length T $\mathbf{s} = (w_1, \dots, w_T)$ into a sequence of T dimension d vectors:

$$f_{GRU}(\mathbf{s}) = (\mathbf{h}_1, \dots, \mathbf{h}_T) \in \mathbb{R}^{d \times T}$$

Multi-Headed Attention Our models also use multi-head attention over the sentence representation. We use the implementation of (Klein et al., 2017), with a residual connection. Given K matrices $\mathbf{M}^\alpha = (M_1^\alpha, \dots, M_K^\alpha) \in \mathbb{R}^{d \times d \times K}$, we define:

$$\alpha_n^k = \text{softmax}\left(\frac{\mathbf{x}^\top M_k^\alpha(\mathbf{h}_1, \dots, \mathbf{h}_T)}{\sqrt{d}}\right)$$

$$\mathbf{x}^\alpha = \sum_{k=1}^K \alpha_n^{k\top}(\mathbf{h}_1, \dots, \mathbf{h}_T)$$

$$\text{attn}(\mathbf{x}, (\mathbf{h}_1, \dots, \mathbf{h}_T); \mathbf{M}^\alpha) = \mathbf{x} + \mathbf{x}^\alpha$$

5.1 Node Predictions

The output tree is made up of internal, categorical, and span nodes. We denote each of these sets by \mathcal{I} , \mathcal{C} and \mathcal{S} respectively, and the full set of nodes as $\mathcal{N} = \mathcal{I} \cup \mathcal{C} \cup \mathcal{S}$. Given a sentence, our aim is to predict the state of each of the nodes $n \in \mathcal{N}$ in the corresponding Action Tree.

Each node in an Action Tree is either active or inactive. We denote the state of a node n by $a_n \in \{0, 1\}$. All the descendants of an inactive internal node $n \in \mathcal{I}$ are considered to be inactive. Additionally, each categorical node $n \in \mathcal{C}$ has a set of possible values C^n . Thus, in a specific Action Tree, each active categorical node has a category label $c_n \in \{1, \dots, |C^n|\}$. Finally, active span nodes $n \in \mathcal{S}$ for a sentence of length T have a start and end index $(s_n, e_n) \in \{1, \dots, T\}^2$.

We take the following approach to predicting the state of a tree. First, we compute a node representation \mathbf{r}^n for each node $n \in \mathcal{N}$ based on the input sentence \mathbf{s} :

$$(\mathbf{r}_1, \dots, \mathbf{r}_{|\mathcal{N}|}) = f_{REP}((h_1, \dots, h_T))$$

Then, we compute the probabilities of each of the labels as:

$$\forall n \in \mathcal{N}, \quad p(a_n) = \sigma(\langle \mathbf{r}_n, \mathbf{p}_n \rangle) \quad (1)$$

$$\forall n \in \mathcal{C}, \quad p(c_n) = \text{softmax}(M_n^c \mathbf{r}_n) \quad (2)$$

$$\forall n \in \mathcal{S}, \quad p(s_n) = \text{softmax}(\mathbf{r}_n^\top M_n^s(\mathbf{h}_1, \dots, \mathbf{h}_T))$$

$$p(e_n) = \text{softmax}(\mathbf{r}_n^\top M_n^e(\mathbf{h}_1, \dots, \mathbf{h}_T)) \quad (3)$$

where the following are model parameters:

$$\forall n \in \mathcal{N}, \quad \mathbf{p}_n \in \mathbb{R}^d$$

$$\forall n \in \mathcal{C}, \quad M_n^c \in \mathbb{R}^{d \times d}$$

$$\forall n \in \mathcal{S}, \quad (M_n^s, M_n^e)_n \in \mathbb{R}^{d \times d \times 2}$$

Our proposed baselines differ from each other by how we compute the node representations \mathbf{r}_n from the sentence. We present three implementations f_{REP} in Section 5.2.

5.2 Node Representation

Independent predictions Our first model computes \mathbf{r}_n independently for each node by attending over the sentence representation. More specifically, each node $n \in \mathcal{N}$ has a parameter $\mathbf{v}_n \in \mathbb{R}^n$. We compute \mathbf{r}_n by simply using \mathbf{v}_n to attend over the sequence encoding $(\mathbf{h}_1, \dots, \mathbf{h}_T)$ using K headed attention parameterized $\mathbf{M}^\nu \in \mathbb{R}^{d \times d \times K}$:

$$\mathbf{r}_n = \text{attn}(\mathbf{x}, (\mathbf{h}_1, \dots, \mathbf{h}_T); \mathbf{M}^\nu) \quad (4)$$

Seq2Tree We also implement the recurrent node representation function from Seq2Tree model of (Dong and Lapata, 2016). It uses a recurrent decoder to compute representations for the children of a node in sequence based on the previously predicted siblings and the parent’s representation. So let $n^p \in \mathcal{I}$ be an internal node, let (c_1, \dots, c_m) be its children. Let the recurrent hidden state of a node n be noted as \mathbf{g}^n , and \circ be the concatenation. We then compute:

$$\mathbf{r}_{c_t} = \text{attn}(\mathbf{v}_{c_t} + \mathbf{g}^{c_t-1}, (\mathbf{h}_1, \dots, \mathbf{h}_T); \mathbf{M}^\sigma) \quad (5)$$

$$\mathbf{g}^{c_t} = \begin{cases} f_{rec}(\mathbf{g}^{c_t-1}, \mathbf{v}'_{c_t} \circ \mathbf{g}^{n^p}), & \text{if } a_{c_t} = 1 \\ \mathbf{g}^{c_t-1}, & \text{else} \end{cases} \quad (6)$$

Where $\mathbf{M}^\nu \in \mathbb{R}^{d \times d \times K}$ is a tree-wise parameter (as in the independent prediction case), f_{rec} is the GRU recurrence function, and \mathbf{v}'_{c_t} is a node parameter (one per category for categorical nodes).

SentenceRec One possible limitation of the Seq2tree model predicted above is that the tree-side recurrent update do not directly depend on the input sentence. This can be addressed by a simple modification: we simply add the node representation \mathbf{r}_{c_t} to the input for the recurrent update:

$$\mathbf{r}_{c_t} = \text{attn}(\mathbf{v}_{c_t} + \mathbf{g}^{c_{t-1}}, (\mathbf{h}_1, \dots, \mathbf{h}_T); \mathbf{M}^\sigma) \quad (7)$$

$$\mathbf{g}^{c_t} = \begin{cases} f_{rec}(\mathbf{g}^{c_{t-1}}, (\mathbf{v}'_{c_t} + \mathbf{r}_{c_t}) \circ \mathbf{g}^{n^p}), & \text{if } a_{c_t} = 1 \\ \mathbf{g}^{c_{t-1}}, & \text{else} \end{cases} \quad (8)$$

We refer to this model as SentenceRec.

5.3 Sequential Prediction

We predict the state of the Action Tree given a sentence in a sequential manner, by predicting the state of the nodes ($\{a_n; \forall n \in \mathcal{N}\}$, $\{c_n; \forall n \in \mathcal{C}\}$, and $\{(s_n, e_n); \forall n \in \mathcal{S}\}$) in Depth First Search order. Additionally, since an inactive node’s descendant are all inactive, we can skip the sub-trees rooted at n if we predict $a_n = 0$. Let us thus note the parent of a node n as $\pi(n)$. Given Equations 1 to 3, the log-likelihood of a tree with states $(\mathbf{a}, \mathbf{c}, \mathbf{s}, \mathbf{e})$ given a sentence \mathbf{s} can be written as:

$$\begin{aligned} \mathcal{L} = & \sum_{n \in \mathcal{N}} a_{\pi(n)} \log(p(a_n)) \\ & + \sum_{n \in \mathcal{C}} a_n \log(p(c_n)) \\ & + \sum_{n \in \mathcal{S}} a_n \left(\log(p(s_n)) + \log(p(e_n)) \right) \quad (9) \end{aligned}$$

Not that since in all of our models the representation \mathbf{r}_n of a node n only depends on nodes that have been seen before it in a DFS search, this loss lends itself well to beam search prediction.

6 Experiments

Training Data We train our model jointly on the (virtually unlimited) template generations and set of 33K training rephrases. Early experiments showed that a model trained exclusively on templated generations failed to reach accuracies better than 40% on the validation rephrases. Training on rephrases did a little better (up to 65%) but

still trailed behind models trained on both (around 80%, see Tale 1).

The action types represented in all three test datasets (rephrases, prompts and interactive) are very different, as shown in Figure 2. In order to address both of these issues, we sample training examples evenly between templates and rephrases according to each of the test setting distributions (no replacement til all examples of a subset have been seen).

Modeling Details We use a 2-layer GRU sentence encoder and all hidden layers in our model have dimension $d = 256$. We use pre-trained word embeddings computed with FastText with subword information (Bojanowski et al., 2017), to which we concatenate free learnable dimensions (these are initialized to be 0, and we tried adding 0, 8, 32 and 64 free dimensions). All models are trained with Adagrad, using label smoothing, dropout, and word dropout for regularization. In all settings, we selected the model which reached the best accuracy on the validation rephrases to evaluate on the test sets. We provide our model and training code.

Overview of Results Table 1 presents tree-level accuracies for the proposed training settings. First, we notice that all models are able to reach near-perfect accuracy on generations from our templates, which means they can invert the generation process described in Section 4.1. The accuracy on the validation and test rephrased data is also high, up to 80.7% for the SentenceRec model. However, the worse performance on instructions from both prompts and interactive shows that our setting poses significant generalization challenges. In particular, all models have significant trouble with the prompts, which come from crowd-sourced workers asked to imagine general game commands and may not fit the exact Minecraft setting. Still, 86% of the annotations are valid under our grammar, and we hope that future work will be better able to address the domain shift to be able to predict those.

On the “interactive” commands, the models do a little better. In general, the SentenceRec seems to have a small edge over the base Seq2Tree model, but the main difference seems to be between the independent prediction and recurrent models. While the latter do much better when trained in-distribution (12% absolute gap), the for-

Sampling	Model	Temp. Valid.	Rep. Valid.	Rep. Test	Prompts	Interactive
Rephrases	Independent	0.979	0.810	0.806	0.171	0.307
	Seq2Tree	0.979	0.801	0.794	0.180	0.231
	SentenceRec	0.976	0.814	0.807	0.159	0.255
Prompts	Independent	0.976	0.804	0.773	0.184	0.370
	Seq2Tree	0.987	0.819	0.789	0.176	0.321
	SentenceRec	0.980	0.828	0.776	0.179	0.360
Interactive	Independent	0.976	0.782	0.709	0.179	0.337
	Seq2Tree	0.975	0.802	0.734	0.196	0.454
	SentenceRec	0.980	0.820	0.771	0.195	0.465

Table 1: Success of trained models over various training and test distributions. Each group of three rows corresponds to a distribution over top-level commands used during training. “Rephrases”, “Prompts”, and “Interactive” as in Figure 2. In the columns, “Temp” refers to the templates distribution, “Rep.” to rephrases (from the template distribution), and “Prompts” and “Interactive” as before.

Model	Node	Test P/R/F		
		Rephrases	Prompts	Interactive
Ind.	INT	97/95/96	75/77/76	75/83/79
	CAT	92/90/91	44/56/50	52/66/58
	SPAN	94/91/93	53/42/47	58/51/54
SRec	INT	97/96/97	75/77/76	86/84/85
	CAT	92/91/92	42/54/47	59/66/62
	SPAN	94/93/94	51/43/46	68/56/61

Table 2: Per-node Precision, Recall and F1 for models trained with interactive sampling for all node types

mer does seem to adapt better to the distribution shift when trained using the rephrases or prompts sampling.

Analysis Table 2 gives insights into model behaviors on CATegorical, INTernal and SPAN nodes. Accurate prediction of a categorical or span node depends on having predicted all of the internal nodes on the path to the root, which explains why CAT and SPAN P/R/F numbers are lower than INT. Additionally, both models have more trouble predicting span than categorical nodes.

We also computed confusion matrices for the best SentenceRec model (see Appendix C). For internal nodes, both models seem to have trouble identifying the scope of some LOCATION and REPEAT nodes: i.e. even when identifying that the command specifies a location, is it the location where the command needs to be executed, or the action where the command’s argument is located? There is also a confusion between SCHEMATIC and ACTION REFERENCE OBJECTS, which we assume comes from the difficulty of interpreting whether the speaker is asking the model to build an object it knows (SCHEMATIC) or another copy of an object in the world (REFERENCE OBJECT), a prediction which must rely on an understanding of

the context.

Finally, the internal parent being absent seems to account for most of the CAT and SPAN mistakes (aside from the action type, which is a child of the root). For action types, the model seems to have trouble recognizing questions and Fill requests mostly. The model also seems to often confuse Mobs (animated creatures in the game) with Objects, which is indeed difficult to disambiguate without some background knowledge. For spans, the model mostly makes the mistake of predicting a node as inactive when it is present. It should be noted that span mis-match are especially rare, except for the model sometimes confusing the depth and height of an object when both are present.

7 Conclusion

In this work, we have described a grammar over a control system for a Minecraft assistant. We then discussed the creation of a dataset of natural language utterances with associated logical forms from this grammar that can be executed in-game. Finally, we showed the results of using this new dataset to train several neural models for parsing natural language instructions. We find that the models we trained were able to fit the templated data nearly perfectly and the rephrased data with some accuracy, but struggled to adapt to the human-generated data. In our view, the problem of using the small number of annotated (grammar-free) human data with the infinite generations of our grammar to improve results on human-distribution to be an exciting area of research.

References

- Stephan Alaniz. 2018. Deep reinforcement learning with model learning and monte carlo tree search in minecraft. *arXiv preprint arXiv:1803.08456*.
- Fraser Allison, Ewa Luger, and Katja Hofmann. 2018. How players speak to an intelligent game character using natural language messages. *Transactions of the Digital Games Research Association*, 4(2).
- Yoav Artzi and Luke Zettlemoyer. 2013. Weakly supervised learning of semantic parsers for mapping instructions to actions. *Transactions of the Association for Computational Linguistics*, 1:49–62.
- Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching word vectors with subword information. *TACL*, 5:135–146.
- Qingqing Cai and Alexander Yates. 2013. Large-scale semantic parsing via schema matching and lexicon extension. In *Proceedings of the 51st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 423–433.
- Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Dzmitry Bahdanau, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. 2014. Learning phrase representations using RNN encoder-decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL*, pages 1724–1734.
- Cristian Danescu-Niculescu-Mizil and Lillian Lee. 2011. Chameleons in imagined conversations: A new approach to understanding coordination of linguistic style in dialogs. In *Proceedings of the Workshop on Cognitive Modeling and Computational Linguistics, ACL 2011*.
- Li Dong and Mirella Lapata. 2016. Language to logical form with neural attention. *arXiv preprint arXiv:1601.01280*.
- Kelvin Guu, Panupong Pasupat, Evan Zheran Liu, and Percy Liang. 2017. From language to programs: Bridging reinforcement learning and maximum marginal likelihood. *arXiv preprint arXiv:1704.07926*.
- Robin Jia and Percy Liang. 2016. Data recombination for neural semantic parsing. *arXiv preprint arXiv:1606.03622*.
- Matthew Johnson, Katja Hofmann, Tim Hutton, and David Bignell. 2016. The malmo platform for artificial intelligence experimentation. In *IJCAI*, pages 4246–4247.
- Nikita Kitaev and Dan Klein. 2017. Where is misty? interpreting spatial descriptors by modeling regions in space. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 157–166.
- Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander M. Rush. 2017. **Opennmt: Open-source toolkit for neural machine translation**. In *Proc. ACL*.
- Thomas Kollar, Danielle Berry, Lauren Stuart, Karolina Owczarzak, Tagyoung Chung, Lambert Mathias, Michael Kayser, Bradford Snow, and Spyros Matsoukas. 2018. The alexa meaning representation language. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 3 (Industry Papers)*, volume 3, pages 177–184.
- Thomas Kollar, Jayant Krishnamurthy, and Grant P Strimel. 2013. Toward interactive grounded language acquisition. In *Robotics: Science and systems*, volume 1, pages 721–732.
- Chen Liang, Jonathan Berant, Quoc Le, Kenneth D Forbus, and Ni Lao. 2016. Neural symbolic machines: Learning semantic parsers on free-base with weak supervision. *arXiv preprint arXiv:1611.00020*.
- Cynthia Matuszek, Evan Herbst, Luke Zettlemoyer, and Dieter Fox. 2013. Learning to parse natural language commands to a robot control system. In *Experimental Robotics*, pages 403–415. Springer.
- Arvind Neelakantan, Quoc V Le, Martin Abadi, Andrew McCallum, and Dario Amodei. 2016. Learning a natural language interface with neural programmer. *arXiv preprint arXiv:1611.08945*.
- Junhyuk Oh, Valliappa Chockalingam, Satinder Singh, and Honglak Lee. 2016. Control of memory, active perception, and action in minecraft. *arXiv preprint arXiv:1605.09128*.
- Junhyuk Oh, Satinder Singh, Honglak Lee, and Pushmeet Kohli. 2017. Zero-shot task generalization with multi-task deep reinforcement learning. *arXiv preprint arXiv:1706.05064*.
- Patti J Price. 1990. Evaluation of spoken language systems: The atis domain. In *Speech and Natural Language: Proceedings of a Workshop Held at Hidden Valley, Pennsylvania, June 24-27, 1990*.
- Sebastian Riedel, Matko Bosnjak, and Tim Rocktäschel. 2016. Programming with a differentiable forth interpreter. *CoRR, abs/1605.06640*.
- Tianmin Shu, Caiming Xiong, and Richard Socher. 2017. Hierarchical and interpretable skill acquisition in multi-task reinforcement learning. *arXiv preprint arXiv:1712.07294*.
- Lappoon R Tang and Raymond J Mooney. 2001. Using multiple clause constructors in inductive logic programming for semantic parsing. In *European*

Conference on Machine Learning, pages 466–477. Springer.

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

Chen Tessler, Shahar Givony, Tom Zahavy, Daniel J Mankowitz, and Shie Mannor. 2017. A deep hierarchical approach to lifelong learning in minecraft. In *AAAI*, volume 3, page 6.

Jesse Thomason, Shiqi Zhang, Raymond J Mooney, and Peter Stone. 2015. Learning to interpret natural language commands through human-robot dialog. In *Twenty-Fourth International Joint Conference on Artificial Intelligence*.

Hiroto Udagawa, Tarun Narasimhan, and Shim-Young Lee. 2016. Fighting zombies in minecraft with deep reinforcement learning. Technical report, Technical report, Stanford University.

Sida I Wang, Samuel Ginn, Percy Liang, and Christopher D Manning. 2017. Naturalizing a programming language via interactive learning. *arXiv preprint arXiv:1704.06956*.

Sida I Wang, Percy Liang, and Christopher D Manning. 2016. Learning language games through interaction. *arXiv preprint arXiv:1606.02447*.

Yushi Wang, Jonathan Berant, and Percy Liang. 2015. Building a semantic parser overnight. In *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, volume 1, pages 1332–1342.

Kexin Yi, Jiajun Wu, Chuang Gan, Antonio Torralba, Pushmeet Kohli, and Josh Tenenbaum. 2018. Neural-symbolic vqa: Disentangling reasoning from vision and language understanding. In *Advances in Neural Information Processing Systems*, pages 1039–1050.

Victor Zhong, Caiming Xiong, and Richard Socher. 2017. Seq2sql: Generating structured queries from natural language using reinforcement learning. *arXiv preprint arXiv:1709.00103*.

A Action Tree structure

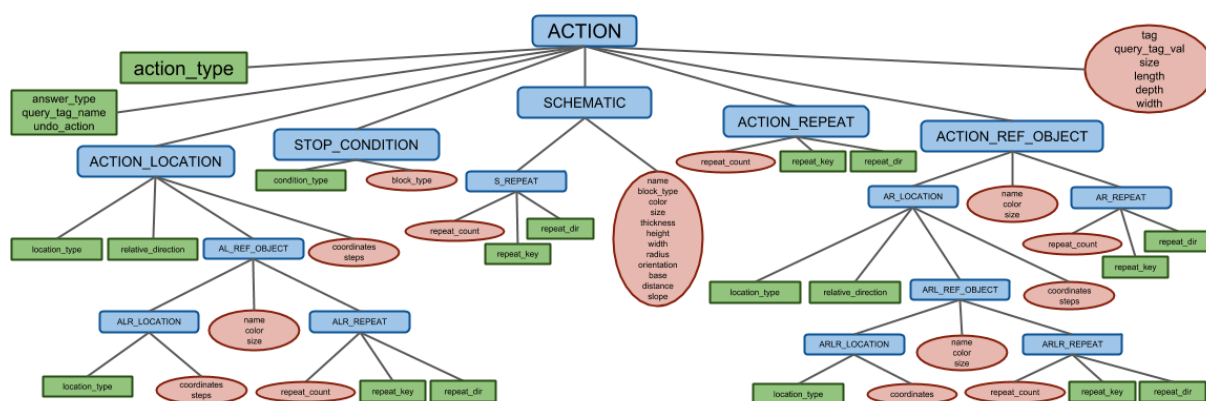


Figure 3: Action space grammar

This section describes the details of the action tree. We support the following actions in our dataset : Build, Copy, Noop, Spawn, Resume, Fill, Destroy, Move, Undo, Stop, Dig, Tag, FreeBuild and Answer. The detailed action tree for each action has been presented in the following subsections. Figure 4 shows an example for a BUILD action.

```

0      1      2      3      4      5      6
"Make three oak wood houses to the
7      8      9      10     11     12
left of the dark grey church."

```

```

{"Build": {
  "schematic": {
    "has_block_type_": [2, 3],
    "has_name_": [4, 4],
    "repeat": {
      "repeat_key": "FOR",
      "repeat_count": [1, 1]
    }
  },
  "location": {
    "relative_direction": "LEFT",
    "location_type": "BlockObject",
    "location_reference_object": {
      "has_colour_": [10, 11],
      "has_name_": [12, 12]
    }
  }
}}

```

Figure 4: An example action tree. The word indices are numbered here for clarity.

A.1 Build Action

This is the action to Build a schematic at an optional location. The Build action tree is shown in 5 and the action can have one of the following as its child:

- location only
- schematic only
- location and schematic both
- neither

A.2 Copy Action

This is the action to copy a block object to an optional location. The copy action is represented as a "Build" with an optional "action reference object" in the tree. The tree is shown in 6.

Copy action can have one the following as its child:

```

{ "Build" : {
  "location" : {
    "location_type" : Coordinates / BlockObject / Mob / AgentPos / SpeakerPos / SpeakerLook,
    "steps" : span,
    "relative_direction" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK' / 'AWAY',
    "coordinates" : span,
    "location_reference_object" : {
      "repeat" : {
        "repeat_key" : 'FOR' / 'ALL'
        "repeat_count" : span,
        "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
      }
      "has_name_" : span,
      "has_size_" : span,
      "has_colour_" : span,
      "location" : {
        "location_type" : Coordinates / AgentPos / SpeakerPos / SpeakerLook,
        "coordinates" : span
      } } },
    "schematic" : {
      "repeat" : {
        "repeat_key" : 'FOR'
        "repeat_count" : span,
        "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
      }
      "has_block_type_" : span,
      "has_name_" : span,
      "has_size_" : span,
      "has_orientation_" : span,
      "has_thickness_" : span,
      "has_colour_" : span,
      "has_height_" : span,
      "has_length_" : span,
      "has_radius_" : span,
      "has_slope_" : span,
      "has_width_" : span,
      "has_base_" : span,
      "has_distance_" : span
    },
    "repeat" : {
      "repeat_key" : 'FOR'
      "repeat_count" : span,
      "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
    } } }
}

```

Figure 5: Details of Build action tree

- action reference object
- action reference object and location
- neither

A.3 Spawn Action

This action indicates that the specified object should be spawned in the environment. The tree is shown in: 7

Spawn action has the child: action reference object.

A.4 Fill Action

This action states that a hole / negative shape at an optional location needs to be filled up. The tree is explained in : 8

Fill action can have one of the following as its child:

- location
- nothing

```

{ "Build" : {
  "location" : {
    "location_type" : Coordinates / BlockObject / Mob / AgentPos / SpeakerPos / SpeakerLook,
    "steps" : span,
    "relative_direction" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK' / 'AWAY',
    "coordinates" : span,
    "location_reference_object" : {
      "repeat" : {
        "repeat_key" : 'FOR' / 'ALL'
        "repeat_count" : span,
        "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
      }
      "has_name_" : span,
      "has_size_" : span,
      "has_colour_" : span,
      "location" : {
        "location_type" : Coordinates / AgentPos / SpeakerPos / SpeakerLook,
        "coordinates" : span
      } } },
    "action_reference_object" : {
      "has_size_" : span,
      "has_colour_" : span,
      "has_name_" : span,
      "repeat" : {
        "repeat_key" : 'FOR' / 'ALL'
        "repeat_count" : span,
        "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
      }
      "location" : {
        "location_type" : Coordinates / AgentPos / SpeakerPos / SpeakerLook,
        "coordinates" : span
      } } },
    "repeat" : {
      "repeat_key" : 'FOR'
      "repeat_count" : span,
      "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
    } } }
}

```

Figure 6: Details of Copy action tree

A.5 Destroy Action

This action indicates the intent to destroy a block object at an optional location. The tree is shown in: [9](#)
 Destroy action can have one of the following as the child:

- action reference object
- nothing

A.6 Move Action

This action states that the agent should move to the specified location, the corresponding tree is in: [10](#)
 Move action can have one of the following as its child:

- location
- stop condition (stop moving when a condition is met)
- location and stop condition
- neither

A.7 Dig Action

This action represents the intent to dig a hole / negative shape of optional dimensions at an optional location. The tree is in [11](#)

Dig action can have one of the following as its child:

```

{ "Spawn" : {
  "action_reference_object" : {
    "repeat" : {
      "repeat_key" : 'FOR'
      "repeat_count" : span,
      "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
    }
    "has_name_" : span,
  },
  "repeat" : {
    "repeat_key" : 'FOR'
    "repeat_count" : span,
    "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
  } } }

```

Figure 7: Details of Spawn action tree

```

{ "Fill" : {
  "location" : {
    "relative_direction" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK' / 'AWAY',
    "location_type" : Coordinates / BlockObject / Mob / AgentPos / SpeakerPos / SpeakerLook,
    "coordinates" : span,
    "steps" : span,
    "location_reference_object" : {
      "has_name_" : span,
      "repeat" : {
        "repeat_key" : 'FOR' / 'ALL'
        "repeat_count" : span,
        "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
      },
      "has_colour_" : span,
      "has_size_" : span,
      "location" : {
        "location_type" : Coordinates / AgentPos / SpeakerPos / SpeakerLook,
        "coordinates" : span
      } } },
    "repeat" : {
      "repeat_key" : 'FOR' / 'ALL'
      "repeat_count" : span,
      "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
    } } }

```

Figure 8: Details of Fill action tree

- nothing
- location
- stop condition
- location and stop condition and / or size, length, depth, width

A.8 Tag Action

This action represents that an action reference object should be tagged with the given tag and the tree is shown in: [12](#)

Tag action can have the following as its children:

- tag
- action reference object

A.9 FreeBuild Action

This action represents that the agent should complete an already existing half-finished block object, using its mental model. The action tree is explained in: [13](#)

FreeBuild action can have one of the following as its child:

```

{ "Destroy" : {
  "action_reference_object" : {
    "location" : {
      "relative_direction" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK' / 'AWAY',
      "location_type" : Coordinates / BlockObject / Mob / AgentPos / SpeakerPos / SpeakerLook
      "coordinates" : span,
      "steps" : span,
      "location_reference_object" : {
        "has_name_" : span,
        "repeat" : {
          "repeat_key" : 'FOR' / 'ALL'
          "repeat_count" : span,
          "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
        },
        "has_colour_" : span,
        "has_size_" : span,
        "location" : {
          "location_type" : Coordinates / AgentPos / SpeakerPos / SpeakerLook,
          "coordinates" : span
        }
      },
      "has_colour_" : span,
      "has_name_" : span,
      "repeat" : {
        "repeat_key" : 'FOR' / 'ALL',
        "repeat_count" : span,
        "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
      }
    }
  }
}
} } }

```

Figure 9: Details of Destroy action tree

- action reference object only
- action reference object and location

A.10 Answer Action

This action represents the agent answering a question about the environment. This is similar to the setup in Visual Question Answering. The tree is represented in: [14](#)

Answer action has the following as its children: action reference object, answer type, query tag name and query tag val. Answer type represents the type of expected answer : counting, querying a specific attribute or querying everything ("what is the size of X" vs "what is X")

A.11 Noop Action

This action indicates no operation should be performed, the tree is shown in : [15](#)

A.12 Resume Action

This action indicates that the previous action should be resumed, the tree is shown in: [16](#)

A.13 Undo Action

This action states the intent to revert the specified action, if any. The tree is in [17](#). Undo action can have on of the following as its child:

- undo action
- nothing (meaning : undo the last action)

A.14 Stop Action

This action indicates stop and the tree is shown in [18](#)

```

{ "Move" : {
  "location" : {
    "relative_direction" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK' / 'AWAY',
    "location_type" : Coordinates / BlockObject / Mob / AgentPos / SpeakerPos / SpeakerLook,
    "coordinates" : span,
    "steps" : span,
    "location_reference_object" : {
      "has_name_" : span,
      "repeat" : {
        "repeat_key" : 'FOR' / 'ALL',
        "repeat_count" : span,
        "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
      }
      "has_colour_" : span,
      "has_size_" : span,
      "location" : {
        "location_type" : Coordinates / AgentPos / SpeakerPos / SpeakerLook,
        "coordinates" : span
      }
    }
  },
  "stop_condition" : {
    "condition_type" : AdjacentToBlockType / Never,
    "block_type" : span
  },
  "repeat" : {
    "repeat_key" : 'FOR',
    "repeat_count" : span,
    "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
  }
} } }

```

Figure 10: Details of Move action tree

```

{ "Dig" : {
  "location" : {
    "relative_direction" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK' / 'AWAY',
    "location_type" : Coordinates / BlockObject / Mob / AgentPos / SpeakerPos / SpeakerLook,
    "coordinates" : span,
    "steps" : span,
    "location_reference_object" : {
      "has_name_" : span,
      "repeat" : {
        "repeat_key" : 'FOR' / 'ALL',
        "repeat_count" : span,
        "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
      },
      "has_colour_" : span,
      "has_size_" : span,
      "location" : {
        "location_type" : Coordinates / AgentPos / SpeakerPos / SpeakerLook,
        "coordinates" : span
      }
    }
  },
  "stop_condition" : {
    "condition_type" : AdjacentToBlockType / Never,
    "block_type" : span
  },
  "repeat" : {
    "repeat_key" : 'FOR',
    "repeat_count" : span,
    "repeat_dir": 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
  },
  "has_size_" : span,
  "has_length_" : span,
  "has_depth_" : span,
  "has_width_" : span } }

```

Figure 11: Details of Dig action tree


```

{ "Tag" : {
  "action_reference_object" : {
    "location" : {
      "location_type" : Coordinates / BlockObject / Mob / AgentPos / SpeakerPos / SpeakerLook
      "relative_direction" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK' / 'AWAY',
      "coordinates" : span,
      "location_reference_object" : {
        "repeat" : {
          "repeat_key" : 'FOR' / 'ALL',
          "repeat_count" : span,
          "repeat_dir" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
        },
        "has_name_" : span,
        "has_size_" : span,
        "has_colour_" : span,
        "location" : {
          "location_type" : Coordinates / AgentPos / SpeakerPos / SpeakerLook,
          "coordinates" : span
        }
      }
    },
    "has_size_" : span,
    "has_colour_" : span,
    "has_name_" : span,
    "repeat" : {
      "repeat_key" : 'FOR' / 'ALL',
      "repeat_count" : span,
      "repeat_dir" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
    }
  }
}
"tag" : span } }

```

Figure 12: Details of Tag action tree

```

{ "FreeBuild" : {
  "action_reference_object" : {
    "location" : {
      "location_type" : Coordinates / BlockObject / Mob / AgentPos / SpeakerPos / SpeakerLook,
      "relative_direction" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK' / 'AWAY',
      "coordinates" : span,
      "location_reference_object" : {
        "repeat" : {
          "repeat_key" : 'FOR' / 'ALL',
          "repeat_count" : span,
          "repeat_dir" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
        },
        "has_name_" : span,
        "has_size_" : span,
        "has_colour_" : span,
        "location" : {
          "location_type" : Coordinates / AgentPos / SpeakerPos / SpeakerLook,
          "coordinates" : span
        }
      }
    },
    "has_size_" : span,
    "has_colour_" : span,
    "has_name_" : span,
    "repeat" : {
      "repeat_key" : 'FOR' / 'ALL',
      "repeat_count" : span,
      "repeat_dir" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
    }
  },
  "location" : {
    "relative_direction" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK' / 'AWAY',
    "location_type" : Coordinates / BlockObject / Mob / AgentPos / SpeakerPos / SpeakerLook,
    "coordinates" : span,
    "steps" : span,
    "location_reference_object" : {
      "has_name_" : span,
      "repeat" : {
        "repeat_key" : 'FOR' / 'ALL',
        "repeat_count" : span,
        "repeat_dir" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK'
      },
      "has_colour_" : span,
      "has_size_" : span,
      "location" : {
        "location_type" : Coordinates / AgentPos / SpeakerPos / SpeakerLook,
        "coordinates" : span
      }
    }
  }
}
} } }

```

Figure 13: Details of FreeBuild action tree

```

{ "Answer" : {
  "action_reference_object" : {
    "location" : {
      "location_type" : BlockObject / Mob / SpeakerLook,
      "relative_direction" : 'LEFT' / 'RIGHT' / 'UP' / 'DOWN' / 'FRONT' / 'BACK' / 'AWAY',
      "coordinates" : span,
      "location_reference_object" : {
        "has_name_" : span,
        "has_size_" : span,
        "has_colour_" : span,
        "location" : {
          "location_type" : Coordinates / AgentPos / SpeakerPos / SpeakerLook,
          "coordinates" : span
        } } },
      "has_size_" : span,
      "has_colour_" : span,
      "has_name_" : span,
    }
  }
  "answer_type" : "query_tag" / "query_all_tags" / "count"
  "query_tag_name" : "size" / "colour" / "name"
  "query_tag_val" : span }

```

Figure 14: Details of Answer action tree

```

{ "Noop" : { } }

```

Figure 15: Details of Noop action tree

```

{ "Resume" : { } }

```

Figure 16: Details of Resume action tree

```

{ "Undo" : { "undo_action" : ActionName } }

```

Figure 17: Details of Undo action tree

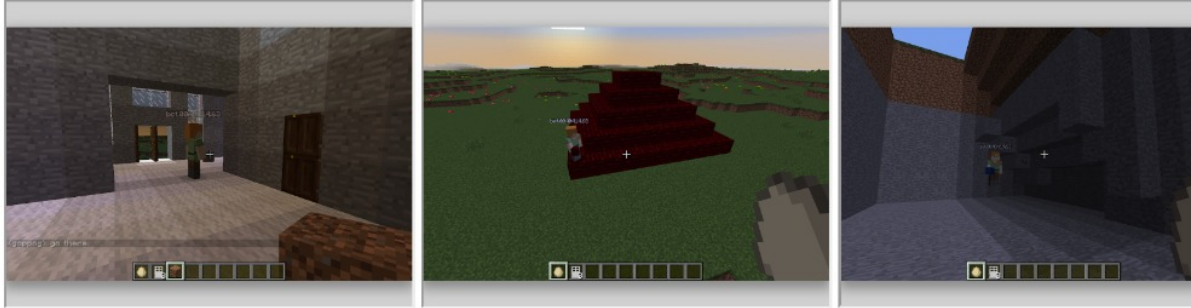
```

{ "Stop" : { } }

```

Figure 18: Details of Stop action tree

Give instructions to a robot assistant in a creative game.



Imagine you are playing a creative game and have access to a robot assistant (as shown in the pictures above) that things done.

Use your imagination and creativeness to come up with three different instructions you'd like to give to the bot. Keep in mind that the bot is non-violent. Each command should only correspond to a single instruction or query (**no that**).

You can assume that the bot has the following capabilities:

- **Move** (the bot can be told to move around to places, objects, things etc)
- **Build** (the bot can be told to build things)
- **Destroy** (the bot can be told to destroy things)
- **Tag** (the bot can be told to tag things)
- **Dig** (the bot can be told to dig up things)
- **Copy** (the bot can be told to make a copy of things)
- **Undo** (the bot can be told to revert things)
- **Fill** (the bot can be told to fill up structures)
- **Spawn** (the bot can be told to spawn objects)
- **Complete** (the bot can be told to complete unfinished things)
- **Answer** (the bot can be told to answer your questions)
- **Stop** (the bot can be told to stop)
- **Resume** (the bot can be told to resume an earlier action)

You can talk to the bot as you would to another human player to instruct it to perform actions me give it information.

Please use your imagination to come up with things you'd like to say to the bot, given it's capabilities. The instructions should be related to the bot's capabilities, please give us as much variety as you can. Give us three unique instructions in the boxes below. Remember that the bot is non-violent.

Figure 19: The task instructions shown to crowd-sourced workers for the Image and text prompts task

B Crowd-sourced task instructions

We have listed the instructions for each task mentioned in section 4.2 in the following subsections.

B.1 Image and Text Prompts

The instructions shown to workers are shown in 19.

B.2 Interactive Gameplay

The instructions shown to workers are shown in 20.

B.3 Annotation tool

The instructions shown to workers are shown in 21.

Welcome to the Minecraft Assistant project!

In this project, you will play Minecraft with a bot that is there to help you. Use the in-game chat to talk to the bot. Tell it where to go. Tell it what to do. Try to build something with its help.

For this project, you are asked to interact with the bot for **10 minutes**. You may ask it to do anything you'd like. Explore what the bot can and cannot do.

The bot may ask you questions so that it can learn over time. Please try to answer the questions it asks.

When you are done, you will be asked to fill out a quick survey to help the bot improve.

Figure 20: The task instructions shown to crowd-sourced workers for the interactive game play

Each of these sentences is spoken to an assistant who is tasked with helping the speaker. We are looking to determine the meaning of the commands given to the assistant.

For each command, answer a series of questions. Each question is either multiple-choice, or requires you to select which words in the sentence correspond to which part of the command.

For example, given the command **"Build a house next to the river"**

- For "What action is being instructed?", the answer is "Build"
- For "What should be built?", select the button for "house"
- For "Where should it be built?", the answer is "Relative to other object(s)"
- For "What other object(s)?", select the buttons for "the river"
- etc.

There may not be a suitable answer that captures the meaning of a sentence. Don't be afraid to select "Other" if there is no good answer.

Command: \${command}

What action is being instructed?

- Build something
- Move
- Destroy, remove, or kill something
- Dig a hole
- Fill a hole
- Assign a description, name, or tag to an object
- Stop current action
- Resume previous action
- Undo previous action
- Answer a question
- Another action not listed here
- This sentence is not a command or request to do something

How many times should this action be performed?

- Just once, or not specified
- Repeatedly, a specific number of times
- Repeatedly, once for each object
- Repeated forever
- Other

Figure 21: The task instructions shown to crowd-sourced workers for the annotation tool

C Confusion Matrices

To compute confusion matrices for the internal node predictions, we look at the gold labels, and add:

- 1 to the gold label count when it is present in the predictions
- $\frac{1}{\#predictions}$ for each predicted internal node that does not match a gold label node when the gold label is not present.

To compute confusion matrices for the categorical node predictions, we look at the gold labels, and add:

- 1 to the predicted class count when the node is present in the predictions (whether it is the gold class or not)
- 1 to the NO-PARENT result when the node's parent is absent in the predicted set
- 1 to the ABSENT result when the node's parent is present in the predicted set but the categorical node is absent

To compute confusion matrices for the categorical node predictions, we look at the gold labels, and add:

- 1 to the MATCH-SPAN result count when the span node is present in the predictions with the right span
- 1 to the MIS-SPAN result count when the span node is present in the predictions with the wrong span
- 1 to the NO-PARENT result when the node's parent is absent in the predicted set
- 1 to the ABSENT result when the node's parent is present in the predicted set but the categorical node is absent

```

{'action': [('total', 817.0), ('action', 1.0)],
 'action_location': [('total', 73.0),
                     ('action_location', 0.8219),
                     ('ar_location', 0.1438),
                     ('schematic', 0.0137),
                     ('action_ref_object', 0.0137),
                     ('arl_ref_object', 0.0068)],
 'action_ref_object': [('total', 125.0),
                      ('action_ref_object', 0.64),
                      ('schematic', 0.16),
                      ('action_location', 0.0787),
                      ('al_ref_object', 0.0707),
                      ('stop_condition', 0.0347),
                      ('s_repeat', 0.012),
                      ('alr_repeat', 0.004)],
 'al_ref_object': [('total', 16.0),
                  ('al_ref_object', 0.875),
                  ('ar_location', 0.0938),
                  ('arl_ref_object', 0.0312)],
 'ar_repeat': [('total', 5.0),
              ('ar_repeat', 0.8),
              ('schematic', 0.1),
              ('s_repeat', 0.1)],
 's_repeat': [('total', 7.0),
             ('s_repeat', 0.8571),
             ('action_repeat', 0.0714),
             ('action_ref_object', 0.0714)],
 'schematic': [('total', 279.0),
              ('schematic', 0.9534),
              ('action_ref_object', 0.0376),
              ('ar_repeat', 0.0018),
              ('action_location', 0.0018),
              ('al_ref_object', 0.0018),
              ('ar_location', 0.0018),
              ('action_repeat', 0.0018)]}

```

Figure 22: Confusion matrix for the internal node predictions by the SentenceRec model.

```

{'action:action_type': {'Answer': [('total', 52.0),
                                   ('Noop', 0.6923),
                                   ('Answer', 0.2308),
                                   ('Build', 0.0192),
                                   ('Dig', 0.0192),
                                   ('Destroy', 0.0192),
                                   ('Move', 0.0192)],
                        'Build': [('total', 358.0),
                                   ('Build', 0.8492),
                                   ('Noop', 0.0782),
                                   ('Spawn', 0.0559),
                                   ('Dig', 0.0168)],
                        'Destroy': [('total', 57.0),
                                    ('Destroy', 0.8246),
                                    ('Noop', 0.0702),
                                    ('Spawn', 0.0526),
                                    ('Move', 0.0175),
                                    ('Dig', 0.0175),
                                    ('Build', 0.0175)],
                        'Dig': [('total', 37.0),
                               ('Dig', 0.973),
                               ('Resume', 0.027)],
                        'Fill': [('total', 7.0),
                                  ('Fill', 0.5714),
                                  ('Build', 0.2857),
                                  ('Dig', 0.1429)],
                        'FreeBuild': [('total', 7.0),
                                       ('Build', 0.7143),
                                       ('Resume', 0.1429),
                                       ('Noop', 0.1429)],
                        'Move': [('total', 36.0),
                                ('Move', 0.8611),
                                ('Noop', 0.1389)],
                        'Noop': [('total', 113.0),
                                ('Noop', 0.8053),
                                ('Build', 0.0973),
                                ('Move', 0.0265),
                                ('Dig', 0.0265),
                                ('Spawn', 0.0177),
                                ('Answer', 0.0177),
                                ('Fill', 0.0088)],
                        'OtherAction': [('total', 112.0),
                                       ('Noop', 0.5446),
                                       ('Move', 0.1786),
                                       ('Spawn', 0.0804),
                                       ('Dig', 0.0625),
                                       ('Destroy', 0.0536),
                                       ('Build', 0.0268),
                                       ('Fill', 0.0179),
                                       ('Undo', 0.0089),
                                       ('FreeBuild', 0.0089),
                                       ('Answer', 0.0089),
                                       ('Stop', 0.0089)]},

```

Figure 23: Confusion matrix for the categorical node predictions by the SentenceRec model part 1.

```

'Resume': [('total', 3.0), ('Resume', 1.0)],
'Spawn': [('total', 16.0),
          ('Spawn', 0.875),
          ('Build', 0.125)],
'Stop': [('total', 10.0),
         ('Stop', 0.8),
         ('Move', 0.1),
         ('Destroy', 0.1)],
'Tag': [('total', 7.0),
        ('Answer', 0.4286),
        ('Build', 0.2857),
        ('Noop', 0.1429),
        ('Tag', 0.1429)],
'Undo': [('total', 2.0), ('Undo', 1.0)]},
'action_location:location_type': {'AgentPos': [('NO-PARENT', 1.0),
                                               ('total', 1.0)],
                                  'BlockObject': [('total', 31.0),
                                                  ('NO-PARENT', 0.2903),
                                                  ('Mob', 0.2581),
                                                  ('AgentPos', 0.2581),
                                                  ('BlockObject', 0.1935)],
                                  'Other': [('total', 4.0),
                                           ('AgentPos', 0.5),
                                           ('SpeakerPos', 0.25),
                                           ('NO-PARENT', 0.25)],
                                  'SpeakerLook': [('total', 19.0),
                                                  ('NO-PARENT', 0.5789),
                                                  ('SpeakerLook', 0.4211)],
                                  'SpeakerPos': [('total', 37.0),
                                                ('SpeakerPos', 0.7297),
                                                ('NO-PARENT', 0.2703)]},
'action_location:relative_direction': {'AWAY': [('total', 2.0), ('AWAY', 1.0)],
                                       'FRONT': [('total', 3.0),
                                                ('NO-PARENT', 0.6667),
                                                ('FRONT', 0.3333)],
                                       'LEFT': [('total', 2.0), ('LEFT', 1.0)],
                                       'UP': [('total', 2.0), ('UP', 1.0)]},
'ar_repeat:repeat_key': {'FOR': [('total', 5.0),
                                  ('FOR', 0.8),
                                  ('NO-PARENT', 0.2)]},
's_repeat:repeat_key': {'FOR': [('total', 7.0),
                                 ('FOR', 0.8571),
                                 ('NO-PARENT', 0.1429)]}}

```

Figure 24: Confusion matrix for the categorical node predictions by the SentenceRec model part 2.


```

{'action:has_depth_': [('total', 6.0),
                       ('ABSENT', 0.6667),
                       ('MIS-SPAN', 0.3333)],
 'action:has_size_': [('total', 9.0),
                      ('MATCH-SPAN', 0.6667),
                      ('ABSENT', 0.3333)],
 'action:has_width_': [('total', 2.0), ('MIS-SPAN', 1.0)],
 'action:tag': [('total', 7.0), ('ABSENT', 0.8571), ('MATCH-SPAN', 0.1429)],
 'action_ref_object:has_block_type_': [('total', 9.0), ('NO-PARENT', 1.0)],
 'action_ref_object:has_colour_': [('total', 2.0),
                                    ('MATCH-SPAN', 0.5),
                                    ('NO-PARENT', 0.5)],
 'action_ref_object:has_name_': [('total', 192.0),
                                  ('NO-PARENT', 0.5833),
                                  ('MATCH-SPAN', 0.3385),
                                  ('ABSENT', 0.0781)],
 'action_ref_object:has_size_': [('ABSENT', 1.0), ('total', 1.0)],
 'al_ref_object:has_name_': [('total', 29.0),
                             ('NO-PARENT', 0.5172),
                             ('MATCH-SPAN', 0.3448),
                             ('MIS-SPAN', 0.1379)],
 'ar_repeat:repeat_count': [('total', 5.0),
                             ('MATCH-SPAN', 0.8),
                             ('NO-PARENT', 0.2)],
 's_repeat:repeat_count': [('total', 7.0),
                            ('MATCH-SPAN', 0.8571),
                            ('NO-PARENT', 0.1429)],
 'schematic:has_block_type_': [('total', 37.0),
                                ('MATCH-SPAN', 0.7838),
                                ('ABSENT', 0.1351),
                                ('MIS-SPAN', 0.0541),
                                ('NO-PARENT', 0.027)],
 'schematic:has_colour_': [('total', 2.0), ('MATCH-SPAN', 1.0)],
 'schematic:has_height_': [('total', 7.0),
                            ('ABSENT', 0.5714),
                            ('MIS-SPAN', 0.2857),
                            ('NO-PARENT', 0.1429)],
 'schematic:has_name_': [('total', 336.0),
                          ('MATCH-SPAN', 0.7262),
                          ('NO-PARENT', 0.2083),
                          ('MIS-SPAN', 0.0625),
                          ('ABSENT', 0.003)],
 'schematic:has_size_': [('total', 21.0),
                          ('MATCH-SPAN', 0.5714),
                          ('NO-PARENT', 0.1905),
                          ('ABSENT', 0.1905),
                          ('MIS-SPAN', 0.0476)],
 'schematic:has_width_': [('total', 7.0),
                           ('ABSENT', 0.8571),
                           ('NO-PARENT', 0.1429)]}

```

Figure 25: Confusion matrix for the span node predictions by the SentenceRec model.