

Benchmarking Clarifying Questions for Effective Collaboration in Grounded Instruction-Based Interactions

Anonymous ACL submission

Abstract

Motivated by the adaptability of human intelligence across various tasks and multi-modal environments, the research community is actively engaged in developing interactive agents capable of engaging in natural conversations with humans and assisting them in real-world tasks. These agents need the ability to request feedback in the form of situated clarifying questions when communication breaks down or instructions are unclear. This paper delves into an extensive investigation of the production of clarifying questions within the context of human-centered AI instruction-based interaction, using a Minecraft environment as a grounding framework. The unique challenges presented by this scenario include the agent's requirement to navigate and complete tasks in a complex, virtual environment, relying on natural language instructions and action states.

In this paper, we made the following contributions: (i) a crowd-sourcing tool for collecting grounded language instructions along with clarifying questions in times when instructions are not clear at scale with low costs; (ii) a substantial dataset of grounded language instructions accompanied by clarifying questions; and (iii) several state-of-the-art baselines for requesting feedback in case of unclear instructions. These contributions are suitable as a foundation for further research.

1 Introduction

One of the long-lasting goals of AI agents (Wingo, 1972) is the ability to seamlessly interact with humans in natural language to help humans learn new skills (Narayan-Chen et al., 2019a; Kiseleva et al., 2022a; Zhang et al., 2021; Wang et al., 2023a) or assist in solving tasks (Shridhar et al., 2019; Kiseleva et al., 2022b). To achieve the latter, the agent must understand and respond to human language to execute instructions in a given environment (Skrynnik et al., 2022; Kiseleva et al., 2022a,b). Over the years, researchers have proposed many tasks

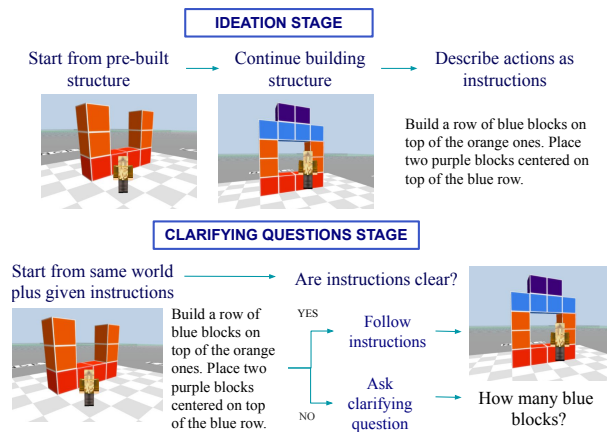


Figure 1: An example of human-agent interactive collaboration, where the goal is to build a given structure, and the agent needs to decide whether to follow the instruction or ask a clarifying question

to tackle this human-AI collaboration challenge, many centered around humans providing instructions to the agent to solve a goal (Gluck and Laird, 2018; Shridhar et al., 2020). An example is the blocks world task, where the agent must understand human instructions to move blocks on a grid (Wingo, 1972; Bisk et al., 2016). Other setups use Minecraft (Gray et al., 2019a), to move objects around (Abramson et al., 2020), to simulate human behavior (Park et al., 2023), or to simulate household tasks (Shridhar et al., 2019; Wang et al., 2023b). However, the instructions humans provide are often inherently ambiguous for most tasks. To complete these tasks successfully, agents need to engage in conversation by asking clarifying questions (Aliannejadi et al., 2021a; Shi et al., 2022; Press et al., 2022), which creates a naturally friendly interface for humans (Nass and Moon, 2000).

We aim to provide an in-depth investigation into the production of clarifying questions for grounded instruction-based interaction using a Minecraft environment, which has shown its effectiveness in studying human-AI collaboration (Fan et al.; Wang et al., 2023a; Kanervisto et al., 2022). This scenario presents a unique challenge, as the agent must

Table 1: Comparison between relevant platforms.

Dataset	Settings	Size of dataset	Collaborative instructional (AI/Human)	Availability of Data collection tool	Availability of Training environment
SHRDLURN(Wang et al., 2016)	Building game	100 games (10,223 utterances)	✓	✗	✓
Voxelurn(Wang et al., 2017)	building structures	230 structures (36,589 utterances)	✓	✗	✓
CEREAL-BAR(Suhr et al., 2019)	collaborative game	1202	✓	✗	✓
ALFRED(Shridhar et al., 2019)	Household tasks	25,743	✗	✗	✓
CVDN(Thomason et al., 2019)	Navigation	2050	✓	✓	✓
TEACH(Padmakumar et al., 2022)	Household tasks	3215	✓	✗	✓
MineDojo (Fan et al., 2022)	Minecraft	730K YouTube videos, 7K Wiki pages, 340K Reddit posts	✓	N/A	N/A
MineRL (Guss et al., 2019)	Minecraft	500 video hours	✗	✗	✓
HoloAssist (Wang et al., 2023b)	Physical tasks	166 video hours	✓	N/A	N/A
Ours	Collaborative building	9,111 utterances/1,1142 clarifying questions	✓	✓	✓

068 navigate and complete tasks in a complex, virtual
069 environment, relying solely on natural language
070 instructions. To ensure successful task completion,
071 the agent must identify gaps in the instructions
072 and pose relevant clarifying questions, as demon-
073 strated in Fig. 1. By tackling this problem head-on,
074 we intend to pave the way for more effective and
075 user-friendly human-AI agent interactions.

076 A significant challenge hindering the exploration
077 of building interactive agents (Narayan-Chen et al.,
078 2019b; Bara et al., 2021) is the scarcity of appro-
079 priate datasets, and scalable and easily extendable
080 data collection tools. These deficiencies have im-
081 peded progress in the field and pose a considerable
082 obstacle to developing effective solutions. Our
083 work addresses this challenge by proposing a novel
084 dataset and scalable data collection methodology,
085 thus contributing to the field’s progress. We believe
086 our work will enable researchers to explore new
087 avenues and enhance user experience in human-AI
088 interactions by addressing this important obstacle.
089 In summary, our main contributions are:

090 **C1 Crowdsourcing Tool for Collecting Interac-**
091 **tive Grounded Language Instructions** specif-
092 ically designed for efficiently gathering interac-
093 tive grounded language instructions within
094 a Minecraft-like environment. With low costs,
095 we can do so at a large scale because it does
096 not require multiple players to be online simul-
097 taneously (Sec. 3).

098 **C2 Extendable Dataset of Human-to-Human**
099 **Grounded Language Instructions** that is ac-
100 companied by clarifying questions (Sec. 4).
101 This dataset represents a valuable resource for
102 various research directions, including but not
103 limited to building structures based on given
104 instructions or predicting clarifying questions.

105 **C3 Baselines for Predicting Clarifying Ques-**
106 **tions** on the aforementioned dataset which
107 serves as a benchmark for evaluating the per-
108 formance of future models (Sec. 5).

2 Related Work

109 Natural Language Interfaces (NLIs) have been
110 a subject of study in various disciplines, includ-
111 ing human-computer interaction and information
112 search, for several decades. Early works such
113 as (Woods et al., 1972; Codd, 1974; Hendrix et al.,
114 1978) laid the foundation for understanding and
115 designing effective interfaces for human language
116 communication with computers. 117

118 **Evolution of NLIs and Applications:** Tab. 1
119 demonstrates a comprehensive set of related plat-
120 forms. In recent years, there has been a resurgence
121 of interest in NLIs due to advances in language un-
122 derstanding capabilities driven by large-scale deep
123 learning models (Devlin et al., 2018; Liu et al.,
124 2019; Clark et al., 2020; Adiwardana et al., 2020;
125 Roller et al., 2020; Brown et al., 2020; OpenAI,
126 2023; Chowdhery et al., 2022) and the increasing
127 demand for various applications such as virtual
128 assistants, dialog systems (Li et al., 2019, 2020c;
129 Burtsev et al., 2017; Li et al., 2020b, 2021), and
130 question answering systems (Liu and Lane, 2017,
131 2018; Dinan et al., 2020; Zhang et al., 2019). NLIs
132 now extend beyond traditional databases to encom-
133 pass knowledge bases (Copestake and Jones, 1990;
134 Berant et al., 2013) to robots (Tellex et al., 2011),
135 personal assistants (Kiseleva et al., 2016b,a), and
136 other forms of interaction (Fast et al., 2018; Desai
137 et al., 2016; Young et al., 2013; Su et al., 2017).

138 **Agent Interactivity and Learning:** The focus
139 has shifted towards interactivity and continuous
140 learning, enabling agents to interact with users,
141 learning new tasks from instructions (Li et al.,
142 2020a), assessing their uncertainty (Yao et al.,
143 2019), asking clarifying questions (Aliannejadi
144 et al., 2020a, 2021b; Arabzadeh et al., 2022), and
145 leveraging feedback from humans to correct mis-
146 takes (Elgohary et al., 2020; Nguyen et al., 2022;
147 Nguyen and au2, 2019). Currently, LLMs are also
148 being studied to asses uncertainty and their own
149 errors (Press et al., 2022; Ren et al., 2023).

Grounded Language Understanding: This paper focuses on grounded language understanding—connecting natural language instructions with real-world or simulated environment context and taking corresponding actions (Hermann et al., 2017; Mitsuda et al., 2022). This is crucial to enabling more effective communication between humans and intelligent agents. Our work focuses specifically on tackling grounded language understanding in the context of collaborative building tasks performed by agents, as highlighted in (Carta et al., 2023; Kiseleva et al., 2021, 2022b; Mehta et al., 2023; Mohanty et al., 2022; Skrynnik et al., 2022).

Leveraging Minecraft for Grounded Language Understanding: We select Minecraft for grounded language understanding due to its distinct advantages. Szlam et al. (2019) highlights the benefits of an open interactive assistant in Minecraft, offering a cost-effective alternative to real-world assistants. The game’s 3D voxel gridworld and adherence to simple physics rules provide ample research scenarios for reinforcement learning experimentation. Minecraft’s interactive nature, player interactions, and dialog exchanges offer diverse opportunities for grounded natural language understanding (Yao et al., 2020; Srinet et al., 2020; Narayan-Chen et al., 2019b). The game’s immense popularity ensures enthusiastic player interaction, facilitating rich human-in-the-loop studies. Minecraft’s advantage extends to the availability of the highly developed set of tools for logging agents interactions and deploying agents for evaluation with human-in-the-loop, including *Malmo* (Johnson et al., 2016), *Craftassist* (Gray et al., 2019b), *TaskWorldMod* (Ogawa et al., 2020), *MC-Saar-Instruct* (Köhn et al., 2020) and *IGLU GridWorld* (Zholus et al., 2022).

3 Data Collection Tool

Narayan-Chen et al., 2019b proposed a setup for a collaborative building task within the Minecraft environment where an Architect is provided with a target structure that needs to be built by the Builder. The Architect provides instructions through a chat on how to create the target structure, and the Builder can ask clarifying questions if an instruction is unclear (Zhang et al., 2021). This approach required installing Microsoft’s Project Malmo (Johnson et al., 2016) client, which provides an API for Minecraft agents to chat, build, and the ability to save and load game states, which makes it limited to lab-based

studies. The setup collects multi-turn interactions between the Architect and the Builder, collaboratively working towards building a given target structure. However, having multiple players online adds unnecessary complications, such as waiting while one of the players is typing, and costs.

We have developed and released an open-source data collection tool¹ that is specifically designed to facilitate the collection of data for multi-modal collaborative building tasks. Our tool eliminates the need for participants to install a local client and allows multiple participants simultaneously annotating data, consequently streamlining the data collection process. As such, it enables 1) *Integration with Crowdsourcing platforms:* Our work has the ability to merge and integrate seamlessly into any crowdsourcing platforms for efficient participant scaling and collecting more data. 2) *Bidirectional Dataset:* While most datasets are one-way, our dataset is bidirectional. It can be used to teach both architects and builders, facilitating more comprehensive language understanding in collaborative building tasks. and 3) *Game Environment for Testing:* We employ a game-type environment, which is more scalable and easier for testing compared to video-based approaches. This choice of environment enhances the practicality and efficiency of our approach. We have used the Amazon Mechanical Turk (MTurk) as the crowd-sourcing platform after obtaining approval from the Institutional Review Board (IRB). Each participant or annotator submits a HIT (Human Intelligence Task). A HIT is comprised the CraftAssist (Gray et al., 2019b) voxelworld and a form which is customizable for different tasks. The form includes rules for a given task and a segment where task instructions or clarifying questions for the building task. The CraftAssist voxelworld is a framework that provides tools and a platform for dialog-enabled interactive agents that learn from natural language interactions. The library provides a 3-d voxelworld grid where agents perform building actions that can be recorded as action states and retrieved for following sessions. Current actions supported by the integrated CraftAssist framework include picking, placing, and removing blocks of different colors within the voxelworld. Agents can also jump to place blocks. These actions enable agents to create structures of varying complexity. Fig. 5 in the appendix illustrates the MTurk views of the task with the embedded voxelworld.

¹<https://bit.ly/42ZUNf7>

Table 2: Statistics of Multi-Turn Dataset

Target Structures	31
Completed Games	127
Median Duration of Completed Games	16 mins
Utterances	811
Avg. Length of Instructions	19.32 words
Clarifying Questions	126

4 Datasets

We built corpora of multi-modal data, which could be used towards solving wide-ranging NLP and RL tasks, including training interactive agents by demonstrations given natural language instructions (Skrynnik et al., 2022). Our research initially concentrates on multi-turn interactions, following a similar approach as presented by (Narayan-Chen et al., 2019b) (Sec. 4.1). To enhance the size of our dataset, we subsequently expanded our data collection efforts to a Single-Turn dataset (Sec. 4.2) to gather a larger corpus of data more efficiently. The datasets and accompanying code for analysis and visualization is openly available ².

4.1 Multi-Turn Dataset

The Multi-Turn dataset comprises dialog-behavior sequences, which we called *game* (Appendix Fig. 2). The sequences either start from scratch for a given goal structure or build on intermediate results. In each *turn*, an annotator takes on the role of either the Architect or the Builder. Architects provide the next step instruction, while the Builder starts with an empty world and executes the instruction or poses a clarifying question. We have improved the data collection process by introducing asynchronous turn-taking. This means the tool no longer relies on having the same annotators online throughout the game. We have implemented checks to prevent a single annotator from taking on both architect and builder roles for the same structure. Importantly, this asynchronous approach allows for the simultaneous launch of multiple structures. Annotators can work on different structures concurrently without waiting for responses, saving time and making process scalable.

Tab. 2 shows the summary of the Multi-Turn dataset. There are 31 goal structures presented to annotators to build. We process and clean the data by filtering out missing and low-quality submissions such as very short instructions. Finally, we have 127 completed game sessions, with the median duration of a game being around 16 minutes. A game

²<https://bit.ly/43WhnGC>

session is considered complete when the Builder can complete building a given goal structure after interacting with and following instructions provided by the Architect. This is denoted by the Architect marking the structure as “*complete*”. Across all the games, we had 811 utterances or dialog interactions between the Architect and Builder annotators. The average length of instructions provided by the Architects was around 19 words, and the number of clarifying questions asked by the Builders – 126 (for all the filtered games).

To provide a deeper understanding of the covered structures in our multi-turn dataset, we performed manual labeling on the 31 structures. The labels and the corresponding number of structures in the dataset in brackets, are as follows: 1. *flat* [7]: all blocks on the ground 2. *flying* [27]: there are blocks that cannot be fully-added without removing some other blocks 3. *diagonal* [6]: some blocks are adjacent (in the vertical axis) diagonally 4. *tricky* [6]: some blocks are hidden or they should be placed in a specific order 5. *tall* [25]: a structure cannot be built without the agent being high enough (the placement radius is 3 blocks) We consider different categories of the structures to make sure the agent is using different skills and abilities and also to make sure the target structures are diverse. For instance, if all the structures are flat, the agent will never learn to use other actions, such as flying. This diversity is essential for training a robust and adaptable agent.

4.2 Single-Turn Dataset

From our extensive study on Multi-Turn data collection, we identified certain challenges that crowd-sourced annotators encountered when engaging in the collaborative building task and issuing instructions for specific target structures. To enhance the crowd-sourcing process, we decided to simplify the task. Our approach involved removing the added complexity of building a predefined target structure. Instead, participants were free to perform free-form building actions within the voxelworld while providing instructions that should allow another worker to rebuild the same structure. This modification led to creating Single-Turn task segments, where annotators collaborated asynchronously to construct the same structure. With this adjustment, we were able to collect data at a faster pace, resulting in a larger corpus comprising of natural language instructions, corresponding actions performed based on those

instructions, and a set of clarifying questions. We record and save actions performed by annotators in a key-value pair format that stores the movement of the agent and positional changes of blocks within the voxelworld.

To provide diverse starting canvases for annotators, we utilized the Multi-Turn dataset to load different world states, which served as varying initial conditions for the building process. The process of collecting single-turn instructions and associated clarifying questions is in (Fig. 1):

- An annotator is assigned a world state from the Multi-Turn dataset as the starting point for their building task (Fig. 1: Ideation Stage).
- The annotator is prompted to perform a sequence of actions for a duration of one minute.
- Then, the annotator is required to describe their set of actions in the form of an instruction.
- Another annotator is shown the instruction and asked to perform the steps mentioned. If the instruction is unclear, the annotator specifies it as thus and asks clarification questions (Fig. 1: Clarification Question Stage).

The instructors answered these clarifying questions, and the data related to these clarifying questions has also been released with this dataset. Tab. 3 presents comprehensive statistics on the Single-Turn dataset, currently the largest dataset available for interactive grounded language understanding. We processed and cleaned the collected Single-Turn dataset by following a heuristic approach, which included filtering out samples where the length of instruction was very short. We also checked whether the instruction was in English and evaluated jobs to remove submissions by annotators who provided low-quality instructions, such as providing the same instruction repeatedly. As shown in Tab. 3, the Single-Turn corpus comprises 8,136 pairs of actions and instructions. On average, an instruction has 18 words, which indicates the instructions are descriptive enough for a one-minute of building.

In addition to the processing steps for cleaning instructions, for the clarifying questions, if an annotator marked the instruction as ambiguous, they were supposed to issue a clarifying question else the submission would be filtered out with a warning provided to the annotator. This was to ensure that every instruction annotated as “not clear” is accompanied by at least one clarifying question. Out of 8,136 instructions, 1,056 (12.98%) were annotated as *Not Clear*, thus being ambiguous, and

Table 3: Statistics of Single-Turn Dataset.

Instructions (train/test)		Avg. Length (in words)	
Total	8136 (6843/1293)	Instructions	18.29
Clear	7080 (5951/1129)	Clarifying Questions	12.05
Ambiguous	1056 (892/164)		

7,080 (87.02%) as *Clear* instructions. The average length of clarifying questions is around 12 words. Tab. 6 in the appendix exemplifies a few instructions marked as being unclear, along with clarifying questions issued by annotators. Majority of clarifying questions fall into the following categories:

- *Color*: Questions clarifying the color of the blocks to be used.
- *Direction/Orientation*: Questions clarifying the direction and orientation in the world.
- *Number of blocks*: Questions that clarify the number of blocks to be placed.
- *Identifying blocks to be changed*: Questions clarifying which blocks need to be changed.

It is important to note that we reassessed the annotations for 100 randomly selected instructions to gauge the level of agreement among the annotators. The agreement rate among the three annotators for these 100 instructions falls within the range interpreted as “fair” according to the Krippendorff agreement measure. This suggests that the interpretation of ambiguous instructions can be highly subjective, and moreover, emphasizes the complexity of such a task. While one annotator may perceive an instruction as clear, another may find it ambiguous. Furthermore, different annotators may ask different clarifying questions about the same instruction, as they may identify unclear aspects from various perspectives.

The Single-Turn approach offers several advantages over the sequential nature of the Multi-Turn process, one of which is the independence of each sample, allowing for easier utilization in different tasks. Each turn can be interpreted as a complete set of information, enabling flexibility and versatility in its application as they do not rely on the context of previous turns. This independence allows researchers to extract valuable insights and information from individual turns without considering the entire dialogue sequence. Furthermore, the Single-Turn approach allows for collecting multiple clarifying questions for each instruction augmenting the richness and diversity of the dataset, enabling a deeper understanding of the nuances and challenges in generating clarifying questions.

5 Baselines Models and Evaluation

We have developed baselines for the prediction of clarifying questions in the Architect-Builder task mentioned in (Sec. 4.2) using the Single-Turn dataset. As such, we focus on the following key research questions:

- **When to ask clarifying questions?:** Predicting whether an instruction provided by the Architect is ambiguous or insufficient for the Builder to complete a task successfully indicating further clarification is required.
- **What clarifying question to ask?** When faced with an instruction that is considered ambiguous, this research question focuses on determining the appropriate question to ask for clarification.

It is worth noting that issues related to determining *When* and *What* clarifying questions to ask have gained significant attention in the domains of NLP and information retrieval (IR) (Aliannejadi et al., 2019, 2021b, 2020b; Arabzadeh et al., 2022). However, as far as we are aware, this aspect has not been explored to a great extent in the context of interacting with agents. The following sections present end-to-end pipelines that show promising performance in addressing each research question. All the baselines are made publicly available at ³

In addition to the baselines discussed in the following sections, we ran initial experiments using Large Language Models that highlight their application in solving this task is not a straightforward endeavor. The grounded nature of the task poses challenges when directly employing LLMs. Our experiments have shown that the transformation of voxel world information into textual format and the subsequent prompt engineering required to address these tasks using LLMs can be a complex and resource-intensive process. We recognize the potential benefits of exploring the use of larger language models for this task, which aligns with our future research direction. Further details on employing LLMs for this task can be found in Appendix A.2.

5.1 When: Clarification Need Prediction

We report the performance of baselines in Tab. 4 and utilize the F-1 Score as the evaluation metric as it provides a balanced measure of precision and recall for this classification task of predicting ambiguity in instructions.

³<https://bit.ly/3qZ7QQD>

Table 4: Results of the baselines on *When* to ask clarifying questions.

Baseline	F-1 score
Fine-tuned BERT (Sec. 5.1.1)	0.732
Text-Grid Cross Modularity (Sec. 5.1.2)	0.757
Textual Grid world State (Sec. 5.1.3)	0.761

5.1.1 BERT fine-tuning

Due to the substantial amount of training data in our collected dataset, one straightforward baseline (Aliannejadi et al., 2021b) to determine whether an instruction requires a clarifying question would be fine-tuning LMs such as BERT (Devlin et al., 2018) followed by a classification layer. This approach has shown promising performance on similar classification tasks (Arabzadeh et al., 2022) demonstrated in Tab. 4.

5.1.2 Text-Grid Cross Modularity

This baseline (Shi et al., 2023) has shown improvement over the BERT fine-tuning approach (Sec. 5.1.1) and consists of the following major components: 1) *Utterance Encoder*, where Architect and Builder annotations would be added before each architect utterance A_t and each builder utterance B_t , respectively. Then, the dialogue utterances are represented as $D_t = architectA_t \oplus builderB_t$ at the turn t , where \oplus is the operation of sequence concatenation. The dialogue is encoded through pre-trained language models such as BERT. 2) *World state encoder* aims to represent the pre-built structure using a voxel-based grid. Each grid state is encoded as a 7-dimensional one-hot vector, representing either an empty space or a block of one of six colors. This encoding results in a $7 \times 11 \times 9 \times 11$ representation of the world state. The structure of the World State Encoder is similar to the approach presented in (Jayannavar et al., 2020). It comprises 3D-convolutional layers followed by a Rectified Linear Unit (ReLU) activation function. This configuration allows the encoder to extract meaningful features from the voxel-based grid representation of the world state. By applying convolutional layers and non-linear activation, the World State Encoder captures spatial dependencies and abstract representations of the pre-built structure. 3) *Fusion module* consists of three sub-modules: one Single-Modality and two Cross-Modality. The former modules are based on self-attention layers, and the latter on cross-attention layers. These take as input the world state representation and dialogue history representation. Between every successive pair of grid

single-modality modules or text single-modality modules, there is a cross-modality module. 4) *Linear projection layer*, this component contains one linear projection to obtain a scalar value for the final binary classification through the Sigmoid function. Finally, the combination of the four aforementioned components obtained F-1 score of 0.757 on the task of *When*. While this approach might seem like the model is deciding whether the follower needs to speak, it aligns with the setup where the agent must decide whether to ask clarifying questions or to act on the most likely action that might lead to a successful task completion.

5.1.3 Text-Grid World State

This baseline focuses on mapping the GridWorld state to a textual context, which is then added as a prefix to the verbalizations of the Architect-Agent. This approach utilizes an automated description of the number of blocks per color in the pre-built structures. For instance, a voxel world can be automatically converted into a textual description like ‘There are 4 levels. There are 15 different blocks. At level 0, there are 3 green blocks. Above the 1st level, there are 2 purple, 2 yellow, and 1 green block. Above level 2, there are 3 green blocks. Above the 3rd level, there are 2 yellow and 2 green blocks.’ This description provides important contextual information about the voxel world and contributes to the improved performance of the simple LLM fine-tuning baseline. We note that the proposed approach could be applied to fine-tune any widely used Language Model such as BERT. However, the reported performance was achieved using DeBERTa-v3-base⁴. Overall, including a textual description of the voxelworld has enhanced the simple LM fine-tuning baseline by 4% in terms of performance (Tab. 4). This approach showcases the importance of incorporating relevant contextual information to enhance the understanding and classification of language-guided collaborative tasks.

5.2 What: Clarifying Question Retrieval

What to ask as clarifying questions has shown to be quite a challenging task (Aliannejadi et al., 2019). As such, similar to (Aliannejadi et al., 2020a, 2021b), we simplify the task by ranking a pool of clarifying questions based on their relevance to ambiguous instructions to place the most pertinent clarifying questions at the top of the ranked list. At inference time, the pool was designed to include

⁴<https://bit.ly/3TldRTY>

Table 5: Performance of the baselines on *What* to ask as clarifying question.

Baseline	MRR@20
BM25	0.3410
Text-World Fusion Ranker (5.2.1)	0.5360
State-Instruction Concatenation Ranker (5.2.2)	0.5960

all clarifying questions in the test set. Given that the relevance judgments for this task are sparse. Namely, only one clarifying question per ambiguous instruction is annotated. We evaluate the task using the Mean Reciprocal Rank (MRR) at cut-off 20. This evaluation approach is consistent with well-known benchmarks like MS MARCO (Nguyen et al., 2016). Tab. 5 presents the performance of BM25 (Robertson et al., 2009, 1995), which is a widely used and well-known ranking function used for Information Retrieval, followed by the two introduced baselines, measured using the MRR@20.

5.2.1 Text-World Fusion Ranker

In this baseline, the instruction and the state of the voxel world are represented individually as text representation and world representation, respectively. Further, the encoded text representation and the world representation are concatenated, and the vector is passed through a two-layer MLP to obtain the final representation. The model is trained using a CrossEntropy loss function over 10-fold cross-validation. At inference time, the ensemble predictions of the 10 models are used for the final predictions. In the following, we elaborate on each of the text and world representations:

Text Representation (TR): A frozen DeBERTa-v3-base model has demonstrated promising performance for ranking. This baseline encodes the instructions, followed by a separator and a question. The last 4 layers of DeBERTa are concatenated and passed through a two-layer BiLSTM to acquire TR.

World Representation (WR): WR is utilized to create a 3D grid. The 3D grid represents a three-dimensional matrix representing the voxel world environment. Each block within this grid is represented by a specific number corresponding to different colors in the matrix. This is subsequently passed through a 1D convolutional network to simplify the height dimension (y), and then the resulting vector is passed through a 2D convolutional network to reduce the width/length (x, z) dimensions. The underlying assumption is that height occupies a different semantic space from the interchangeable x, z dimensions. For example, an instruction might include references to a *tower* or *column*, which would

623 be a stack of objects in the y direction, while a *wall*
624 could extend in the x or z direction. Ultimately, the
625 size of the 3D grid is reduced by an AvgPooling
626 layer to a 1D vector. This assumption is essentially
627 made to make the 3D structure simplified into a 2D
628 and then into a 1D representation to reduce the com-
629 plexity of the representation. This simplification
630 is akin to dimensionality reduction techniques and
631 helps make the problem more manageable (Huang
632 et al., 2022; Sainburg et al., 2020; Cao et al., 2018).

633 In addition, it has been revealed that certain
634 straightforward post-processing tricks relying on
635 certain assumptions about the content of questions
636 given a world and instruction could be helpful. For
637 example, the size of the ranking pool could be
638 reduced by excluding questions that don't overlap
639 with the given instructions. If the instruction doesn't
640 mention a color like *blue*, and *blue* is also absent in
641 the world, it can be assumed that the question will
642 not reference the word *blue*. While these heuristic
643 rules may seem somewhat aggressive, they have
644 proven useful in excluding additional questions
645 irrelevant to the instruction, as we see that Text-
646 World Fusion Ranker utilized these approaches.

647 5.2.2 State-Instruction Concatenation Ranker

648 To comprehend the concept of relevance, the ap-
649 proach of aligning queries and relevant items closely
650 in embedding space while distancing queries from ir-
651 relevant items in the same space has proven to be ef-
652 fective (Izacard et al., 2021; Reimers and Gurevych,
653 2019; Karpukhin et al., 2020; Zhan et al., 2021).
654 Similarly, in this baseline, each positive question is
655 paired with sampled irrelevant negative questions
656 drawn from the candidate questions. The similarity
657 between the instruction and the question is then
658 measured using a BERT-like pre-trained LM.

659 To include information from the world state and
660 pre-built structure, state information, such as the
661 colors and numbers of initialized blocks, is encoded
662 in natural language and then concatenated with the
663 instruction. It has been shown that clarifying ques-
664 tions about the same instruction can differ based
665 on the world states (Shi et al., 2022; Aliannejadi
666 et al., 2019; Deng et al., 2023). To avoid redundant
667 state information and improve the model's gener-
668 alization, randomly selecting only one color type
669 of block as the state information has proven help-
670 ful. The state information and raw instruction are
671 then concatenated and labeled with the keywords
672 *state* and *instruction*, respectively. For example,
673 the input could be: *state*: There are 9 green blocks;

674 *instruction*: put a green block on top of the yellow
675 and the two blue ones. To balance the data distri-
676 bution Easy Data Augmentation (EDA) has been
677 adapted (Wei and Zou, 2019), which could expand
678 the dataset by synonym replacement, random inser-
679 tion, random swap, and random deletion, according
680 to a pre-defined ratio. Moreover, taking inspiration
681 from DAPT (Gururangan et al., 2020), datasets
682 such as (Kiseleva et al., 2022b; Narayan-Chen et al.,
683 2019b; Shi et al., 2022; Zholus et al., 2022) are used
684 for performing domain-adaptive fine-tuning. Fur-
685 ther, we propose to use the Fast Gradient Method
686 (FGM), inspired by adversarial training, to mitigate
687 the overfitting problem (Goodfellow et al., 2014).
688 Finally, taking cues from (Gao et al., 2021), the
689 list-wise loss is used to train the model.

690 6 Conclusions and Future Work

691 In conclusion, our paper addresses the crucial issue
692 of enabling natural interaction in grounded human-
693 AI agent collaboration. We achieve this by allowing
694 agents to clarify instructions through a familiar and
695 friendly interface, such as the use of clarifying ques-
696 tions. A significant obstacle hindering progress
697 in this field has been the scarcity of appropriate
698 datasets and scalable, extensible data collection
699 tools. To address this challenge, we developed a
700 crowdsourcing tool specifically designed for col-
701 lecting interactive grounded language instructions
702 within a Minecraft-like environment at a large scale.
703 We created a dataset of human-to-human grounded
704 language instructions, accompanied by clarifying
705 questions which could be useful for a wide range
706 of natural language understanding and reinforce-
707 ment learning tasks. Furthermore, we established
708 baselines for predicting clarifying questions, pro-
709 viding a benchmark for evaluating the performance
710 of future models and algorithms in this domain. As
711 future work, we plan to investigate how LLMs can
712 be applied for our task. We also plan to develop an
713 evaluation framework that incorporates human judg-
714 ments or task-specific metrics to provide a better
715 understanding of the performance and limitations
716 of the proposed methods. Additionally, we plan
717 to conduct comprehensive user studies to evaluate
718 the usability of the generated clarifying questions
719 in real-world scenarios. We anticipate that ex-
720 ploring these future directions will contribute to
721 even a greater understanding of the challenges and
722 potential solutions involved generating clarifying
723 questions for instruction based interactions.

7 Limitations

Our paper centers on the utilization of a Minecraft-like environment to examine human-AI interaction. While this emphasis may not comprehensively encapsulate the intricacies of real-world scenarios, it affords the opportunity to scrutinize specific facets of the problem in an isolated and safe environment. Nevertheless, there are constraints within the task scenarios, including the consideration of potential variations in task complexity. This may constrain the understanding of how the generation of clarifying questions may vary in different contexts. Consequently, the generalizability and applicability of our findings to real-world settings may be influenced by these factors. However, we believe the suggested environment and the data collection tool allow exploration for further scenarios.

Moreover, the paper relies on a crowdsourcing tool for data collection, which introduces the possibility of biases in the dataset. The demographic composition, skill levels, and motivations of the crowd workers may impact the quality and representativeness of the collected data. To mitigate these biases, we introduced sophisticated training and tests for crowd-source workers to enable them to complete tasks. To address any potential ethical issues, all the crowd-source workers signed an IRB.

References

- Josh Abramson, Arun Ahuja, Iain Barr, Arthur Brussee, Federico Carnevale, Mary Cassin, Rachita Chhaparia, Stephen Clark, Bogdan Damoc, Andrew Dudzik, et al. 2020. Imitating interactive intelligence. *arXiv preprint arXiv:2012.05672*.
- Daniel Adiwardana, Minh-Thang Luong, David R So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, et al. 2020. Towards a human-like open-domain chatbot. *arXiv preprint arXiv:2001.09977*.
- Mohammad Aliannejadi, Julia Kiseleva, Aleksandr Chuklin, Jeff Dalton, and Mikhail Burtsev. 2020a. *Convai3: Generating clarifying questions for open-domain dialogue systems (clariq)*.
- Mohammad Aliannejadi, Julia Kiseleva, Aleksandr Chuklin, Jeff Dalton, and Mikhail Burtsev. 2020b. *Convai3: Generating clarifying questions for open-domain dialogue systems (clariq)*. *arXiv preprint arXiv:2009.11352*.
- Mohammad Aliannejadi, Julia Kiseleva, Aleksandr Chuklin, Jeff Dalton, and Mikhail Burtsev. 2021a. *Building and evaluating open-domain dialogue corpora with clarifying questions*. In *Proceedings of the*

2021 Conference on Empirical Methods in Natural Language Processing, pages 4473–4484, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.

- Mohammad Aliannejadi, Julia Kiseleva, Aleksandr Chuklin, Jeffrey Dalton, and Mikhail Burtsev. 2021b. *Building and evaluating open-domain dialogue corpora with clarifying questions*. *arXiv preprint arXiv:2109.05794*.

- Mohammad Aliannejadi, Hamed Zamani, Fabio Crestani, and W Bruce Croft. 2019. *Asking clarifying questions in open-domain information-seeking conversations*. In *Proceedings of the 42nd international acm sigir conference on research and development in information retrieval*, pages 475–484.

- Negar Arabzadeh, Mahsa Seifkar, and Charles LA Clarke. 2022. *Unsupervised question clarity prediction through retrieved item coherency*. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, pages 3811–3816.

- Cristian-Paul Bara, Sky CH-Wang, and Joyce Chai. 2021. *Minecraft: Theory of mind modeling for situated dialogue in collaborative tasks*. *arXiv preprint arXiv:2109.06275*.

- Jonathan Berant, Andrew Chou, Roy Frostig, and Percy Liang. 2013. *Semantic parsing on freebase from question-answer pairs*. In *Proceedings of the 2013 conference on empirical methods in natural language processing*, pages 1533–1544.

- Yonatan Bisk, Deniz Yuret, and Daniel Marcu. 2016. *Natural language communication with robots*. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 751–761.

- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. *Language models are few-shot learners*.

- Mikhail Burtsev, Aleksandr Chuklin, Julia Kiseleva, and Alexey Borisov. 2017. *Search-oriented conversational ai (scai)*. In *Proceedings of the ACM SIGIR International Conference on Theory of Information Retrieval*, pages 333–334.

- Zhihao Cao, MU Shaomin, XU Yongyu, and Mengping Dong. 2018. *Image retrieval method based on cnn and dimension reduction*. In *2018 International Conference on Security, Pattern Analysis, and Cybernetics (SPAC)*, pages 441–445. IEEE.

831	Thomas Carta, Clément Romac, Thomas Wolf, Sylvain Lamprier, Olivier Sigaud, and Pierre-Yves Oudeyer. 2023. Grounding large language models in interactive environments with online reinforcement learning. <i>arXiv preprint arXiv:2302.02662</i> .	886
832		887
833		888
834		889
835		890
836	Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. 2022. Palm: Scaling language modeling with pathways. <i>arXiv preprint arXiv:2204.02311</i> .	891
837		892
838		893
839		894
840		895
841		
842	Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. 2020. Electra: Pre-training text encoders as discriminators rather than generators. <i>arXiv preprint arXiv:2003.10555</i> .	896
843		897
844		898
845		
846	Edgar F Codd. 1974. <i>Seven steps to rendezvous with the casual user</i> . IBM Corporation.	899
847		900
848	Ann Copestake and Karen Sparck Jones. 1990. Natural language interfaces to databases.	901
849		
850	Yang Deng, Shuaiyi Li, and Wai Lam. 2023. Learning to ask clarification questions with spatial reasoning. In <i>Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 2113–2117.	902
851		903
852		904
853		905
854		906
855	Aditya Desai, Sumit Gulwani, Vineet Hingorani, Nidhi Jain, Amey Karkare, Mark Marron, Subhajt Roy, et al. 2016. Program synthesis using natural language. In <i>Proceedings of the 38th International Conference on Software Engineering</i> , pages 345–356. ACM.	907
856		908
857		909
858		910
859		911
860	Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. In <i>Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL)</i> .	912
861		913
862		914
863		915
864		
865		
866	Emily Dinan, Varvara Logacheva, Valentin Malykh, Alexander Miller, Kurt Shuster, Jack Urbanek, Douwe Kiela, Arthur Szlam, Iulian Serban, Ryan Lowe, et al. 2020. The second conversational intelligence challenge (convai2). In <i>The NeurIPS’18 Competition</i> , pages 187–208. Springer, Cham.	916
867		917
868		918
869		919
870		920
871		
872	Ahmed Elgohary, Saghar Hosseini, and Ahmed Hassan Awadallah. 2020. <i>Speak to your parser: Interactive text-to-SQL with natural language feedback</i> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 2065–2077, Online. Association for Computational Linguistics.	921
873		922
874		923
875		924
876		925
877		
878		
879	Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. 2022. Minedojo: Building open-ended embodied agents with internet-scale knowledge. In <i>Thirty-sixth Conference on Neural Information Processing Systems Datasets and Benchmarks Track</i> .	926
880		927
881		928
882		929
883		930
884		931
885		932
		933
		934
		935
		936
		937
		938
		939
		940

941	Gautier Izacard, Mathilde Caron, Lucas Hosseini, Sebastian Riedel, Piotr Bojanowski, Armand Joulin, and Edouard Grave. 2021. Towards unsupervised dense information retrieval with contrastive learning. <i>arXiv preprint arXiv:2112.09118</i> .	intelligent assistants. In <i>Proceedings of the 39th International ACM SIGIR conference on Research and Development in Information Retrieval</i> , pages 45–54.	998 999 1000
942			
943			
944			
945			
946	Prashant Jayannavar, Anjali Narayan-Chen, and Julia Hockenmaier. 2020. Learning to execute instructions in a minecraft dialogue. In <i>Proceedings of the 58th annual meeting of the association for computational linguistics</i> , pages 2589–2602.	Julia Kiseleva, Kyle Williams, Jiepu Jiang, Ahmed Hassan Awadallah, Aidan C Crook, Imed Zitouni, and Tasos Anastasakos. 2016b. Understanding user satisfaction with intelligent assistants. In <i>Proceedings of the 2016 ACM on Conference on Human Information Interaction and Retrieval</i> , pages 121–130.	1001 1002 1003 1004 1005 1006
947			
948			
949			
950			
951	Matthew Johnson, Katja Hofmann, Tim Hutton, and David Bignell. 2016. The malmo platform for artificial intelligence experimentation. In <i>IJCAI</i> , pages 4246–4247. Citeseer.	Arne Köhn, Julia Wichlacz, Christine Schäfer, Alvaro Torralba, Jörg Hoffmann, and Alexander Koller. 2020. Mc-saar-instruct: a platform for minecraft instruction giving agents. In <i>Proceedings of the 21th Annual Meeting of the Special Interest Group on Discourse and Dialogue</i> , pages 53–56.	1007 1008 1009 1010 1011 1012
952			
953			
954			
955	Anssi Kanervisto, Stephanie Milani, Karolis Ramanauskas, Nicholay Topin, Zichuan Lin, Junyou Li, Jianing Shi, Deheng Ye, Qiang Fu, Wei Yang, Weijun Hong, Zhongyue Huang, Haicheng Chen, Guangjun Zeng, Yue Lin, Vincent Micheli, Eloi Alonso, François Fleuret, Alexander Nikulin, Yury Belousov, Oleg Svidchenko, and Aleksei Shpilman. 2022. <i>Minerl diamond 2021 competition: Overview, results, and lessons learned</i> .	Toby Jia-Jun Li, Tom Mitchell, and Brad Myers. 2020a. Interactive task learning from GUI-grounded natural language instructions and demonstrations. In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics: System Demonstrations</i> .	1013 1014 1015 1016 1017 1018
956			
957			
958			
959			
960			
961			
962			
963			
964	Vladimir Karpukhin, Barlas Öğuz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. Dense passage retrieval for open-domain question answering. <i>arXiv preprint arXiv:2004.04906</i> .	Ziming Li, Julia Kiseleva, and Maarten De Rijke. 2019. Dialogue generation: From imitation learning to inverse reinforcement learning. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 33, pages 6722–6729.	1019 1020 1021 1022 1023
965			
966			
967			
968			
969	Julia Kiseleva, Ziming Li, Mohammad Aliannejadi, Shrestha Mohanty, Maartje ter Hoeve, Mikhail Burtsev, Alexey Skrynnik, Artem Zholus, Aleksandr Panov, Kavya Srinet, Arthur Szlam, Yuxuan Sun, Katja Hofmann, Marc-Alexandre Côté, Ahmed Awadallah, Linar Abdrazakov, Igor Churin, Putra Manggala, Kata Naszadi, Michiel van der Meer, and Taewoon Kim. 2022a. Interactive grounded language understanding in a collaborative environment: Iglu 2021. In <i>NeurIPS 2021 Competitions and Demonstrations Track</i> , pages 146–161. PMLR.	Ziming Li, Julia Kiseleva, and Maarten de Rijke. 2020b. Rethinking supervised learning and reinforcement learning in task-oriented dialogue systems. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 3537–3546, Online. Association for Computational Linguistics.	1024 1025 1026 1027 1028 1029
970			
971			
972			
973			
974			
975			
976			
977			
978			
979			
980	Julia Kiseleva, Ziming Li, Mohammad Aliannejadi, Shrestha Mohanty, Maartje ter Hoeve, Mikhail Burtsev, Alexey Skrynnik, Artem Zholus, Aleksandr Panov, Kavya Srinet, et al. 2021. Neurips 2021 competition iglu: Interactive grounded language understanding in a collaborative environment. <i>arXiv preprint arXiv:2110.06536</i> .	Ziming Li, Julia Kiseleva, and Maarten de Rijke. 2021. Improving response quality with backward reasoning in open-domain dialogue systems. In <i>Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , pages 1940–1944.	1030 1031 1032 1033 1034 1035
981			
982			
983			
984			
985			
986			
987	Julia Kiseleva, Alexey Skrynnik, Artem Zholus, Shrestha Mohanty, Negar Arabzadeh, Marc-Alexandre Côté, Mohammad Aliannejadi, Milagro Teruel, Ziming Li, Mikhail Burtsev, et al. 2022b. Interactive grounded language understanding in a collaborative environment: Retrospective on iglu 2022 competition. In <i>NeurIPS 2022 Competition Track</i> , pages 204–216. PMLR.	Ziming Li, Sungjin Lee, Baolin Peng, Jinchao Li, Julia Kiseleva, Maarten de Rijke, Shahin Shayandeh, and Jianfeng Gao. 2020c. Guided dialogue policy learning without adversarial learning in the loop. In <i>Findings of the Association for Computational Linguistics: EMNLP 2020</i> , pages 2308–2317, Online. Association for Computational Linguistics.	1036 1037 1038 1039 1040 1041 1042
988			
989			
990			
991			
992			
993			
994			
995	Julia Kiseleva, Kyle Williams, Ahmed Hassan Awadallah, Aidan C Crook, Imed Zitouni, and Tasos Anastasakos. 2016a. Predicting user satisfaction with	Bing Liu and Ian Lane. 2017. Iterative policy learning in end-to-end trainable task-oriented neural dialog models. In <i>2017 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)</i> , pages 482–489. IEEE.	1043 1044 1045 1046 1047
996			
997			
		Bing Liu and Ian Lane. 2018. Adversarial learning of task-oriented neural dialog models. In <i>Proceedings of the SIGDIAL 2018 Conference</i> , pages 350–359.	1048 1049 1050
		Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. <i>Roberta</i> :	1051 1052 1053

1054	A robustly optimized BERT pretraining approach.	Aishwarya Padmakumar, Jesse Thomason, Ayush Shrivastava, Patrick Lange, Anjali Narayan-Chen, Span-	1108
1055	CoRR, abs/1907.11692.	dana Gella, Robinson PIRAMUTHU, Gökhan Tür, and	1109
1056	Nikhil Mehta, Milagro Teruel, Patricio Figueroa Sanz,	Dilek Hakkani-Tür. 2022. Teach: Task-driven embod-	1110
1057	Xin Deng, Ahmed Hassan Awadallah, and Julia Kisel-	ied agents that chat. In <i>Thirty-Sixth AAAI Conference</i>	1111
1058	eva. 2023. Improving grounded language understand-	<i>on Artificial Intelligence, AAAI 2022, Thirty-Fourth</i>	1112
1059	ing in a collaborative environment by interacting	<i>Conference on Innovative Applications of Artificial</i>	1113
1060	with agents through help feedback. <i>arXiv preprint</i>	<i>Intelligence, IAAI 2022, The Twelveth Symposium on</i>	1114
1061	arXiv:2304.10750.	<i>Educational Advances in Artificial Intelligence, EAAI</i>	1115
1062	Koh Mitsuda, Ryuichiro Higashinaka, Yuhei Oga, and	<i>2022 Virtual Event, February 22 - March 1, 2022,</i>	1116
1063	Sen Yoshida. 2022. Dialogue collection for recording	pages 2017–2025. AAAI Press.	1117
1064	the process of building common ground in a collab-	Joon Sung Park, Joseph C O’Brien, Carrie J Cai, Mered-	1119
1065	orative task. In <i>Proceedings of the Thirteenth Lan-</i>	ith Ringel Morris, Percy Liang, and Michael S Bern-	1120
1066	<i>guage Resources and Evaluation Conference,</i> pages	stein. 2023. Generative agents: Interactive simulacra	1121
1067	5749–5758, Marseille, France. European Language	of human behavior. <i>arXiv preprint arXiv:2304.03442.</i>	1122
1068	Resources Association.	Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt,	1123
1069	Shrestha Mohanty, Negar Arabzadeh, Milagro Teruel,	Noah A Smith, and Mike Lewis. 2022. Measuring	1124
1070	Yuxuan Sun, Artem Zhohus, Alexey Skrynnik,	and narrowing the compositionality gap in language	1125
1071	Mikhail Burtsev, Kavya Srinet, Aleksandr Panov,	models. <i>arXiv preprint arXiv:2210.03350.</i>	1126
1072	Arthur Szlam, Marc-Alexandre Côté, and Julia	Nils Reimers and Iryna Gurevych. 2019. Sentence-bert:	1127
1073	Kiseleva. 2022. Collecting interactive multi-modal	Sentence embeddings using siamese bert-networks.	1128
1074	datasets for grounded language understanding. <i>arXiv</i>	<i>arXiv preprint arXiv:1908.10084.</i>	1129
1075	preprint arXiv:2211.06552.	Allen Z. Ren, Anushri Dixit, Alexandra Bodrova, Sumeet	1130
1076	Anjali Narayan-Chen, Prashant Jayannavar, and Ju-	Singh, Stephen Tu, Noah Brown, Peng Xu, Leila	1131
1077	lia Hockenmaier. 2019a. Collaborative dialogue in	Takayama, Fei Xia, Jake Varley, Zhenjia Xu, Dorsa	1132
1078	minecraft. In <i>Proceedings of the 57th Annual Meet-</i>	Sadigh, Andy Zeng, and Anirudha Majumdar. 2023.	1133
1079	<i>ing of the Association for Computational Linguistics,</i>	Robots that ask for help: Uncertainty alignment for	1134
1080	pages 5405–5415.	large language model planners.	1135
1081	Anjali Narayan-Chen, Prashant Jayannavar, and Ju-	Stephen Robertson, Hugo Zaragoza, et al. 2009. The	1136
1082	lia Hockenmaier. 2019b. Collaborative dialogue in	probabilistic relevance framework: Bm25 and beyond.	1137
1083	Minecraft. In <i>Proceedings of the 57th Annual Meet-</i>	<i>Foundations and Trends® in Information Retrieval,</i>	1138
1084	<i>ing of the Association for Computational Linguistics,</i>	3(4):333–389.	1139
1085	pages 5405–5415, Florence, Italy. Association for	Stephen E Robertson, Steve Walker, Susan Jones, Miche-	1140
1086	Computational Linguistics.	line M Hancock-Beaulieu, Mike Gatford, et al. 1995.	1141
1087	Clifford Nass and Youngme Moon. 2000. Machines	Okapi at trec-3. <i>Nist Special Publication Sp,</i> 109:109.	1142
1088	and mindlessness: Social responses to computers.	Stephen Roller, Emily Dinan, Naman Goyal, Da Ju,	1143
1089	<i>Journal of social issues,</i> 56(1):81–103.	Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott,	1144
1090	Khanh Nguyen and Hal Daumé III au2. 2019. Help,	Kurt Shuster, Eric M Smith, et al. 2020. Recipes	1145
1091	anna! visual navigation with natural multimodal	for building an open-domain chatbot. <i>arXiv preprint</i>	1146
1092	assistance via retrospective curiosity-encouraging	arXiv:2004.13637.	1147
1093	imitation learning.	Tim Sainburg, Marvin Thielk, and Timothy Q Gen-	1148
1094	Khanh Nguyen, Yonatan Bisk, and Hal Daumé III au2.	ter. 2020. Finding, visualizing, and quantifying la-	1149
1095	2022. A framework for learning to request rich and	tent structure across diverse animal vocal repertoires.	1150
1096	contextually useful information from humans.	<i>PLoS computational biology,</i> 16(10):e1008228.	1151
1097	Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao,	Zhengxiang Shi, Yue Feng, and Aldo Lipani. 2022.	1152
1098	Saurabh Tiwary, Rangan Majumder, and Li Deng.	Learning to execute actions or ask clarification ques-	1153
1099	2016. Ms marco: A human generated machine read-	tions. <i>arXiv preprint arXiv:2204.08373.</i>	1154
1100	ing comprehension dataset. <i>choice,</i> 2640:660.	Zhengxiang Shi, Jerome Ramos, To Eun Kim, Xi Wang,	1155
1101	Haruna Ogawa, Hitoshi Nishikawa, Takenobu Tokunaga,	Hossein A. Rahmani, and Aldo Lipani. 2023. When	1156
1102	and Hikaru Yokono. 2020. Gamification platform for	and what to ask through world states and text instruc-	1157
1103	collecting task-oriented dialogue data. In <i>Proceed-</i>	tions: Iglu nlp challenge solution.	1158
1104	<i>ings of the 12th Language Resources and Evaluation</i>	Mohit Shridhar, Jesse Thomason, Daniel Gordon,	1159
1105	<i>Conference,</i> pages 7084–7093, Marseille, France.	Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke	1160
1106	European Language Resources Association.	Zettlemoyer, and Dieter Fox. 2019. ALFRED: A	1161
1107	OpenAI. 2023. Gpt-4 technical report.	benchmark for interpreting grounded instructions for	1162
		everyday tasks. <i>CoRR, abs/1912.01734.</i>	1163

1164	Mohit Shridhar, Jesse Thomason, Daniel Gordon,	Sida I. Wang, Percy Liang, and Christopher D. Manning.	1219
1165	Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke	2016. Learning language games through interaction.	1220
1166	Zettlemoyer, and Dieter Fox. 2020. Alfred: A bench-	<i>CoRR</i> , abs/1606.02447.	1221
1167	mark for interpreting grounded instructions for ev-		
1168	eryday tasks. In <i>Proceedings of the IEEE/CVF con-</i>	Xin Wang, Taein Kwon, Mahdi Rad, Bowen Pan, Ishani	1222
1169	<i>ference on computer vision and pattern recognition</i> ,	Chakraborty, Sean Andrist, Dan Bohus, Ashley Fe-	1223
1170	pages 10740–10749.	niello, Felipe Vieira Frujeri, Neel Joshi, and Marc	1224
		Pollefeys. 2023b. Holoassist: an egocentric human	1225
		interaction dataset for interactive ai assistants in the	1226
1171	Alexey Skrynnik, Zoya Volovikova, Marc-Alexandre	real world. In <i>ICCV 2023.</i>	1227
1172	Côté, Anton Voronov, Artem Zholus, Negar		
1173	Arabzadeh, Shrestha Mohanty, Milagro Teruel,	Jason W. Wei and Kai Zou. 2019. EDA: easy data	1228
1174	Ahmed Awadallah, Aleksandr Panov, Mikhail Burt-	augmentation techniques for boosting performance	1229
1175	sev, and Julia Kiseleva. 2022. Learning to solve voxel	on text classification tasks. <i>CoRR</i> , abs/1901.11196.	1230
1176	building embodied tasks from pixels and natural lan-		
1177	guage instructions. <i>arXiv preprint arXiv:2211.00688.</i>	Terry Winograd. 1972. Understanding natural language.	1231
		<i>Cognitive psychology</i> , 3(1):1–191.	1232
1178	Kavya Srinet, Yacine Jernite, Jonathan Gray, and Arthur		
1179	Szlam. 2020. CraftAssist instruction parsing: Semantic	W. A. Woods, Ronald M Kaplan, and Bonnie L. Webber.	1233
1180	parsing for a voxel-world assistant. In <i>Proceedings</i>	1972. The lunar sciences natural language informa-	1234
1181	<i>of the 58th Annual Meeting of the Association for</i>	tion system: Final report. <i>BBN Report 2378.</i>	1235
1182	<i>Computational Linguistics</i> , pages 4693–4714, Online.		
1183	Association for Computational Linguistics.	Ziyu Yao, Yu Su, Huan Sun, and Wen-tau Yih. 2019.	1236
		Model-based interactive semantic parsing: A unified	1237
1184	Yu Su, Ahmed Hassan Awadallah, Madian Khabsa,	framework and a text-to-SQL case study. In <i>Proceed-</i>	1238
1185	Patrick Pantel, Michael Gamon, and Mark Encarna-	<i>ings of the 2019 Conference on Empirical Methods</i>	1239
1186	cacion. 2017. Building natural language interfaces to	<i>in Natural Language Processing and the 9th Interna-</i>	1240
1187	web apis. In <i>Proceedings of the 2017 ACM on Con-</i>	<i>national Joint Conference on Natural Language Pro-</i>	1241
1188	<i>ference on Information and Knowledge Management</i> ,	<i>cessing (EMNLP-IJCNLP)</i> , pages 5447–5458, Hong	1242
1189	pages 177–186. ACM.	Kong, China. Association for Computational Linguis-	1243
		tics.	1244
1190	Alane Suhr, Claudia Yan, Jacob Schluger, Stanley Yu,	Ziyu Yao, Yiqi Tang, Wen-tau Yih, Huan Sun, and	1245
1191	Hadi Khader, Marwa Mouallem, Iris Zhang, and	Yu Su. 2020. An imitation game for learning semantic	1246
1192	Yoav Artzi. 2019. Executing instructions in situated	parsers from user interaction.	1247
1193	collaborative interactions. <i>CoRR</i> , abs/1910.03655.		
		Steve Young, Milica Gašić, Blaise Thomson, and Jason D	1248
1194	Arthur Szlam, Jonathan Gray, Kavya Srinet, Yacine	Williams. 2013. Pomdp-based statistical spoken di-	1249
1195	Jernite, Armand Joulin, Gabriel Synnaeve, Douwe	alog systems: A review. <i>Proceedings of the IEEE</i> ,	1250
1196	Kiela, Haonan Yu, Zhuoyuan Chen, Siddharth Goyal,	101(5):1160–1179.	1251
1197	Demi Guo, Danielle Rothermel, C. Lawrence Zit-		
1198	nick, and Jason Weston. 2019. Why Build an Assis-	Jingtao Zhan, Jiaxin Mao, Yiqun Liu, Jiafeng Guo,	1252
1199	tant in Minecraft? <i>arXiv:1907.09273 [cs]</i> . ArXiv:	Min Zhang, and Shaoping Ma. 2021. Optimizing	1253
1200	1907.09273.	dense retrieval model training with hard negatives.	1254
		In <i>Proceedings of the 44th International ACM SI-</i>	1255
1201	Stefanie Tellex, Thomas Kollar, Steven Dickerson,	<i>GIR Conference on Research and Development in</i>	1256
1202	Matthew R Walter, Ashis Gopal Banerjee, Seth Teller,	<i>Information Retrieval</i> , pages 1503–1512.	1257
1203	and Nicholas Roy. 2011. Understanding natural lan-		
1204	guage commands for robotic navigation and mobile	Yi Zhang, Sujay Kumar Jauhar, Julia Kiseleva, Ryen	1258
1205	manipulation. In <i>Twenty-Fifth AAAI Conference on</i>	White, and Dan Roth. 2021. Learning to decompose	1259
1206	<i>Artificial Intelligence.</i>	and organize complex tasks. In <i>Proceedings of the</i>	1260
		<i>2021 Conference of the North American Chapter of the</i>	1261
1207	Jesse Thomason, Michael Murray, Maya Cakmak, and	<i>Association for Computational Linguistics: Human</i>	1262
1208	Luke Zettlemoyer. 2019. Vision-and-dialog naviga-	<i>Language Technologies</i> , pages 2726–2735, Online.	1263
1209	tion. <i>CoRR</i> , abs/1907.04957.	Association for Computational Linguistics.	1264
1210	Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Man-	Yizhe Zhang, Siqu Sun, Michel Galley, Yen-Chun Chen,	1265
1211	dlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and	Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing	1266
1212	Anima Anandkumar. 2023a. Voyager: An open-	Liu, and Bill Dolan. 2019. Dialogpt: Large-scale	1267
1213	ended embodied agent with large language models.	generative pre-training for conversational response	1268
1214	<i>arXiv preprint arXiv:2305.16291.</i>	generation. <i>arXiv preprint arXiv:1911.00536.</i>	1269
1215	Sida I. Wang, Samuel Ginn, Percy Liang, and Christo-	Artem Zholus, Alexey Skrynnik, Shrestha Mohanty,	1270
1216	pher D. Manning. 2017. Naturalizing a program-	Zoya Volovikova, Julia Kiseleva, Artur Szlam, Marc-	1271
1217	ming language via interactive learning. <i>CoRR</i> ,	Alexandre Coté, and Aleksandr I Panov. 2022. Iglu	1272
1218	abs/1704.06956.	gridworld: Simple and fast environment for embodied	1273
		dialog agents. <i>arXiv preprint arXiv:2206.00142.</i>	1274

A Appendix

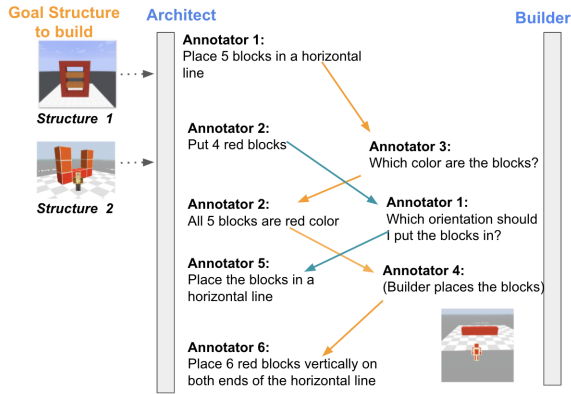


Figure 2: Example of multi-turn data collection, where the Architect can see the goal structure and provides instructions for the Builder. The blue arrows indicate turns for the first goal structure, the orange arrows indicate turns for the second goal structure. Annotators can switch roles between architect and builder for different structures.

A.1 Data Collection Details

Multi-turn Data Collection: In Figure 2, we illustrate an example of multi-turn data collection. In this scenario, the Architect can observe the goal structure and offer instructions to the Builder. The blue arrows represent the turns associated with the first goal structure, while the orange arrows correspond to the turns related to the second goal structure. Annotators can switch roles between architect and builder for different structures. Fig. 2 illustrates this concept of our data collection methodology with different annotators (1, 3, 2, 4, and 6) collaborating to construct Structure 1. Annotators can switch roles between architect and builder for different structures.

Fig. 3 illustrates the overall design of the tool. Our tool can be integrated with crowd-sourcing platforms to provide an interface for participants to complete tasks. Fig. 5 demonstrates MTurk views of the Data Collection Tool (Sec.3) for the Multi-Turn Dataset (Sec.4.1). We have the Architect Task, where the Architect provides instructions to the Builder based on the provided target structure. Next, we have the Builder Task, where instructions and the current structure built so far are shown. The Builder can mark the instructions as unclear or will follow the instructions by adjusting blocks in the voxelworld. Finally, we have the Intermediate Architect Task, where the Architect is shown the

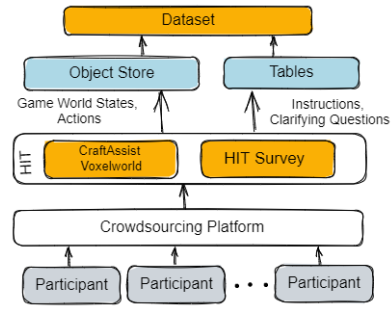


Figure 3: The architecture of the data crowdsourcing collection tool

progress of the structure built so far and provides the next instruction.

Examples of Single-Turn Dataset: Tab. 6 provides examples of instructions marked as unclear in the Single-Turn Dataset along with different kinds of clarifying questions posed by annotators (Sec.4.2). Clarifying questions consist of topics such as color, direction, and identification of blocks.

A.2 Large Language Models as baselines

In our earlier discussion in Section 5, we highlighted that applying Large Language Models (LLMs) to the task of determining when and what to ask as clarifying questions in our designed environment is not a straightforward process. This complexity arises primarily due to the multimodal nature of the task and the significant engineering efforts required to create effective prompts. While we do not suggest that leveraging LLMs for this task is impossible, it is important to clarify that our paper’s primary focus lies in benchmarking and dataset creation. Integrating LLMs into our study falls beyond the scope of this research.

Additionally, We address the challenge of determining when and what clarifying questions to ask by employing a combination of classification and retrieval methods instead of relying on text generation. Our decision was influenced by several factors, including the absence of well-established evaluation procedures for LLM text generation and the need to handle complex structures like action states of the world, which serve as inputs to our current pre-trained model baselines. Nonetheless, we conducted preliminary experiments using GPT-3.5-Turbo to explore their potential applicability to this task.

In these experiments, we randomly selected 50 instructions and utilized their previous utterances as information to reconstruct the pre-built structure. We then prompted the LLM to determine whether,

Table 6: Examples of Unclear Instructions and Clarifying Questions

Instruction	Clarifying Question
Place four blocks to the east of the highest block, horizontally.	Which color blocks?
Destroy 2 purple blocks and then build 3 green blocks diagonally.	Which two purple blocks need to be destroyed?
Destroy the 3 stacked red blocks on the east side. Replace them with 3 stacked blue boxes	Which three of the four stacked red blocks on the east side need to be destroyed?
Make a rectangle that is the width and height of the blue shape and fill it in with purple blocks.	Which side I need to make the rectangle is not clear
Facing South remove the rightmost purple block. Place a row of three orange blocks to the left of the upper leftmost purple block. Place two orange blocks above and below the leftmost orange block.	Which one of the rightmost blocks should be removed?
Facing north and purple-green blocks will be arranged one by one.	Where would you like to place the purple and green blocks exactly?

You are participating in a game set in a Minecraft-like world. In this game, there are two roles: the Architect and the Builder.

1. The Architect: This player provides instructions for building structures in the game environment.
2. The Builder: This player’s role is to follow the instructions given by the Architect and execute them within the game environment.

During the game, the Builder has two response options:

- If the instruction provided by the Architect is clear and can be executed without any need for further clarification or questions, the Builder responds with “yes. The instruction is clear.”
- If the instruction is unclear or requires clarification from the Architect before it can be executed, the Builder responds with “no” and generates a clarification question.

reply only a “Yes. The instruction is clear” or a “No” followed by a relevant clarification question.

Previous Dialogue: <Architect> Facing North, Build a blue block in the left most corner.
 <Architect> Destroy the blue block and build a purple block there.
 <Architect> Facing East, place one green block on the very top right corner of the map.
 Starting Grid world of 3D blocks in the format $(x, y, z, color) : (-5, -1, -5, purple)(-5, 9, 5, 'green')$

Current Instruction: <Architect> In the northeast corner place one blue block. In the southwest corner place one purple block then a red block on top of that.

Figure 4: Example of using LLMs for solving the clarifying question need task.

1344 given the pre-existing instructions from the Archi-
 1345 tect, the new instruction was clear or if the builder
 1346 needed to pose clarifying questions. An example of
 1347 such a prompt can be found in Figure 4, focusing on
 1348 the “when to ask clarifying questions” task. How-
 1349 ever, the results were far from satisfactory when
 1350 compared to the baseline models, yielding an F1
 1351 score of 0.45, which was significantly lower than
 1352 the F1 scores achieved by our baseline models, as
 1353 reported in Table 4. All baseline models achieved
 1354 F1 scores above 0.732.

1355 We believe that the performance of LLMs can
 1356 be enhanced through improved prompt engineering
 1357 and a better representation of the voxel world. How-
 1358 ever, we decided not to include these findings in the

1359 paper to avoid potential misinterpretations. Our pri-
 1360 mary aim was to establish a benchmark with clear
 1361 and reproducible baselines. While we acknowledge
 1362 the potential of LLMs for this task, we consider this
 1363 aspect as part of our future work, which extends
 1364 beyond the scope of the current study.

Instructions:

The goal of this HIT is to give a short instruction to a builder explaining what blocks to place to progress from the current structure depicted in the lower section of images towards the target structure depicted in the upper section of images. The images in the upper and lower sections depict different views of target and currently built structures, respectively.

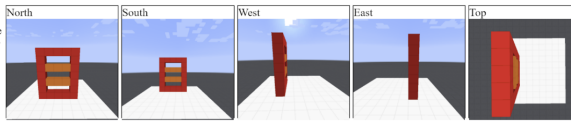
You can assume that the builder that will follow your instructions has the following capabilities:

- **Move** (the builder can be told to move around game area)
- **Build** (the builder can be told to place blocks of different colors)
- **Destroy** (the builder can be told to destroy things)

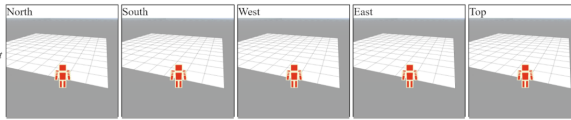
Further notes:

- Give the instructions with the understanding that the builder *cannot* see the target structure. Only a *single* view of the currently built structure will be available to them. This view will be selected by you.
- Make sure to choose (and keep in mind) the view of the **current** structure on which your instruction is based. This is how the builder will be oriented. If it is the first instruction in a game, views contain only a grid. In this case, pick 'North,' View 'Top' is there for your orientation and cannot be picked along with instruction.
- The instruction should be about a minutes' worth of work for the builder to complete.
- Please use grammatical English in your instructions.
- If the current structure is already the same as the target, mark it so.

Target structure from different perspectives:



Current structure built so far (base your instruction on one of these views):



Is the current structure exactly like the target structure? Pick 'yes' only if the images of current structure in the bottom row match that of the structure to the target one in the upper row.

Choose the viewpoint on which you based your instruction on:

Provide an instruction for what the builder should do next to progress to the target:

Facing North, place 5 red blocks in a horizontal line

(a) Architect Task

Please wait! It may take some time for voxel world to load

Important: Click on the 'Submit' at the bottom of the page when you are done

Execute:

Facing North, place 5 red blocks in a horizontal line

For the north viewpoint of the intermediate structure (note if it's the first step, the game area below is empty):



The previous builder found this instruction unclear and has asked the following clarifying question:

do you mean I need to put the five blocks parallel to the board while looking towards north side?

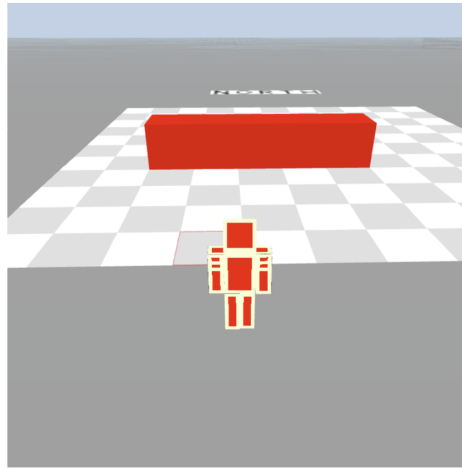
Here's the answer the architect gave to that question:

yes, that's correct

Is this instruction clear now?

Given the information above, please execute the clarified instruction using the following keyboard command to operate in the voxel world:

- **mouse click** on the game area: to activate builder and see the cursor
- **w/arrow**: move forward/left/backward/right
- **space**: jump
- **double click space**: enable flying mode
- **shift**: Move downwards. When you are in flying mode, keep pressing shift until the agent hit the ground. Once the agent hit the ground, flying mode will be turned off.
- **mouse click**: break block
- **1/2/3/4/5/6**: place a blue/yellow/green/orange/purple/red block
- **esc**: leave the Voxel world area



(b) Builder Task

Instructions:

The goal of this HIT is to give a short instruction to a builder explaining what blocks to place to progress from the current structure depicted in the lower section of images towards the target structure depicted in the upper section of images. The images in the upper and lower sections depict different views of target and currently built structures, respectively.

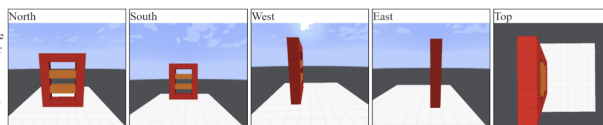
You can assume that the builder that will follow your instructions has the following capabilities:

- **Move** (the builder can be told to move around game area)
- **Build** (the builder can be told to place blocks of different colors)
- **Destroy** (the builder can be told to destroy things)

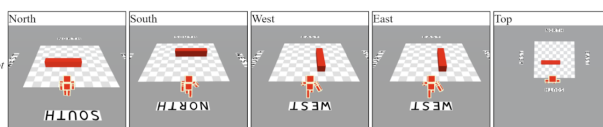
Further notes:

- Give the instructions with the understanding that the builder *cannot* see the target structure. Only a *single* view of the currently built structure will be available to them. This view will be selected by you.
- Make sure to choose (and keep in mind) the view of the **current** structure on which your instruction is based. This is how the builder will be oriented. If it is the first instruction in a game, views contain only a grid. In this case, pick 'North,' View 'Top' is there for your orientation and cannot be picked along with instruction.
- The instruction should be about a minutes' worth of work for the builder to complete.
- Please use grammatical English in your instructions.
- If the current structure is already the same as the target, mark it so.

Target structure from different perspectives:



Current structure built so far (base your instruction on one of these views):



Is the current structure exactly like the target structure? Pick 'yes' only if the images of current structure in the bottom row match that of the structure to the target one in the upper row.

Choose the viewpoint on which you based your instruction on:

Provide an instruction for what the builder should do next to progress to the target:

Place 6 red blocks on both sides of the horizontal line

(c) Intermediate Structure Architect Task

Figure 5: View of MTurk Data Collection Tool