# MDC-R: The Minecraft Dialogue Corpus with Reference

**Chris Madge[1], Maris Camilleri[1], Paloma Carretero Garcia[1], Vanja Karan[2],**
**Juexi Shao[1], Prashant Jayannavar[3], Julian Hough[4], Benjamin Roth[2], Massimo Poesio[1]**
[1]Queen Mary University of London, [2]Universität Wien,
[3]University of Illinois, [4]Swansea University
**Correspondence:** m.poesio@qmul.ac.uk

## Abstract

We introduce the Minecraft Dialogue Corpus with Reference (MDC-R). MDC-R is a new language resource that supplements the original Minecraft Dialogue Corpus (MDC) with expert annotations of anaphoric and deictic reference. MDC's task-orientated, multi-turn, situated dialogue in a dynamic environment has motivated multiple annotation efforts, owing to the interesting linguistic phenomena that this setting gives rise to. We believe it can serve as a valuable resource when annotated with reference, too. Here, we discuss our method of annotation and the resulting corpus, and provide both a quantitative and a qualitative analysis of the data. Furthermore, we carry out a short experiment demonstrating the usefulness of our corpus for referring expression comprehension.

## 1 Introduction

Dialogue games are games in which conversational agents impersonating characters can learn to perform tasks or improve their communicative ability by interacting with human players or other artificial agents. Such games offer an exciting opportunity to study how conversational agents can carry out interaction (Johnson et al., 2016; Urbanek et al., 2019; Narayan-Chen et al., 2019; Szlam et al., 2019; Bara et al., 2021; Kiseleva et al., 2022). Virtual world dialogue games in particular may approach the complexity of the real world (Shridhar et al., 2020; Puig et al., 2018; Kolve et al., 2017) and virtual agents operating in such virtual worlds need to be able to develop a variety of interactional skills to be perceived as "real" . Dialogue games have been argued to be the best benchmark for (grounded) spoken language understanding (Schlangen, 2023). Games are also engaging, so human participants are more likely to be motivated (Von Ahn, 2006; Chamberlain et al., 2008; Poesio et al., 2013, 2019; Szlam et al., 2019).

Virtual world games are a particularly promising domain to study reference to entities in a visual situation (Johnson et al., 2017; Qi et al., 2020; Islam et al., 2022) and anaphoric reference in dialogue (Yu et al., 2022a). Most research on referring expression comprehension has focused on static images and most research on anaphoric reference has focused on written text (Poesio et al., 2023), but dialogue presents unique challenges for both anaphoric and deictic reference, as reference in dialogue is more fluid. Participants propose, negotiate and if necessary repair their common understanding of referents (Clark and Wilkes-Gibbs, 1986). Some recent efforts have looked at reference in situated dialogue (Loáiciga et al., 2022; Gigliobianco et al., 2024), but this research focuses on static environments, rather than a dynamic one in which one of the interlocutors may manipulate the environment over the course of the dialogue.

In this work, we introduce a corpus of conversations between human agents carrying out tasks in the virtual world Minecraft Collaborative Dialogue Task (Narayan-Chen et al., 2019). This task sees two parties collaborate to build a structure in a virtual world, annotated for deictic and anaphoric reference. This yields a dataset that poses unique annotation challenges and opportunities, given its dynamic environment that continuously changes throughout the conversation as it is manipulated by the human participants, introducing complex grounding, misunderstandings and ambiguity. The voxel world task-orientated dialogue requires addressing abstract shapes composed of blocks, and compositions of those abstract shapes. This results in another interesting phenomena known as conceptual pacts (Brennan and Clark, 1996), in which collaborating participants in a conversation come to an agreement over terms referring to entities.

The contributions of this paper are as follows: 1. a new approach to annotating reference in dynamically changing visual worlds; 2. a new annotation

for deictic and anaphoric reference of the popular Minecraft Dialogue Corpus [1]; 3. an analysis of the issues with reference arising; and 4. a baseline model for reference resolution in this corpus.

The rest of the paper is structured as follows. We review related work in Section 2, followed by a discussion of the base MDC dataset in Section 3 and 4. Next, Section 5 outlines our annotation methodology. We then describe our prediction experiments and provide corresponding discussion in Section 6 and Section 7, respectively.

## 2 Related Work

### 2.1 Virtual world dialogue games

Aside from the Minecraft Collaborative Building Task (Narayan-Chen et al., 2019), discussed in Section 3, there have been a few other works that use virtual world games to collect task orientated dialogue between human interlocutors.

Ogawa et al. (2020) created a platform for collection of situated dialogue as an extension to *Minecraft*. As part of their work, they propose a scenario inspired by the *Map Task* (Anderson et al., 1991), in a *Minecraft* setting. *Mindcraft* is a *Minecraft*-like game in which two players, with differing skillsets collaborate to complete a shared goal, communicating to form a common understanding of their respective skills in relation to the given mission (Bara et al., 2021). The aim of the work is to understand how players establish a mental model in a situated task, and how this is affected by communication and the shared environment.

In the *CerealBar* game (Suhr et al., 2019) two separate uniquely constrained participants, must cooperate to collect cards in a 3D environment. The leader, possessing greater observability but comparatively reduced movement, instructs the follower (who is unable to respond).

Virtual world games for dialogue are not limited to 3D interfaces. *LIGHT* (Urbanek et al., 2019) is a text-based multiplayer adventure-fantasy world, designed for studying grounded dialogue in a situated environment. Participants are able to both speak and act in text form. A corpus of 11k episodes was collected, some of which have been used in reference-based tasks (Yu et al., 2022a).

---

[1]Freely available from `https://github.com/arciduca-project/MDC-R`

### 2.2 Referring Expression Comprehension

In a conversation in which the participants share a visual scene two possible types of reference are possible: *anaphoric* reference to entities introduced in the language, and *deictic* reference to objects in the visual scene that may or may not have been mentioned before. Conversational agents interacting in a virtual world need both the ability to comprehend referring expressions and to generate them (van Deemter, 2016).

**The REC task** In recent years, Referring Expression Comprehension (REC) (Kazemzadeh et al., 2014; Mao et al., 2016; Nagaraja et al., 2016) has become a fundamental task in vision-and-language research. In this field, it is typically defined as identifying the correct bounding box for a target object described by a free-form text. Over the years, the scope of REC has expanded to include various types of textual inputs — ranging from conversational queries (De Vries et al., 2017) and scene-aware contexts (Chen et al., 2023) to knowledge-based cues (Wang et al., 2020; You et al., 2022) — which collectively challenge multimodal models to harness more advanced reasoning skills. In parallel, researchers have leveraged simulated environments (Johnson et al., 2017; Qi et al., 2020; Islam et al., 2022) to study compositional and embodied reasoning in settings where objects are meticulously customised. Despite the flexibility of these virtual worlds, existing REC studies have yet to explore simulated settings enriched with extended, dialogue-centric text, highlighting a significant gap that our work aims to address.

**Generalized REC** Generalized Referring Expression Comprehension (GREC) (Liu et al., 2023a; He et al., 2023a) extends traditional REC by allowing expressions to refer to multiple or no targets, thereby enhancing its applicability in real-world scenarios. Recently, RECANTFormer (Hemanthage et al., 2024) has employed a one-stage method to improve GREC recognition. Our corpus aligns with the GREC framework by providing bounding boxes for multiple targets within a single image and unique annotations for each target, as detailed in Section 5.1.

**Visual Coreference** In parallel with the REC work, recent research in multimodal coreference resolution have underscored the significance of integrating visual cues into language understanding. For instance, (Goel et al., 2023b) and (Goel

et al., 2023a) explore coreference resolution in image narrations, demonstrating that grounding textual entities in specific visual regions can substantially enhance reference disambiguation. In addition, research in visual dialogue has made notable progress with works such as (Yu et al., 2019) and (Yu et al., 2022b), which leverage visual context to resolve pronoun references and thereby improve dialogue coherence and interpretability. Loáiciga et al. (2021) delve into the annotation of anaphoric phenomena within situated dialogue, emphasizing how contextual information is critical for resolving referential ambiguities.

## 2.3 Datasets for studying reference in dialogue

REX-J (Spanger et al., 2012) task sees two participants collaborating in separate roles on a tangram puzzle. The solver role communicates with the operator role, proposing a solution, while the operator manipulates the tangram pieces. The PentoRef task (Zarrieß et al., 2016) also makes use of two human participants, an instruction giver (who can see the target structure), and an instruction follower (can manipulate pieces), collaborating to reproduce the target structure. The followers actions are used to infer the effectiveness of the instruction giver's utterances. Reference annotation on the resultant dataset is carried out by experts (Zarrieß et al., 2016). The corpus of Loáiciga et al. (2021) features the *Cups* task, in which two participants collaborate to identify discrepancies in an almost identical 3D environment, given images taken from two different perspectives. The dataset is expert annotated for reference with MMAX2 (Müller and Strube, 2006), according to ARRAU (Poesio et al., 2024a) (similar to the approach used in this work).

All of the aforementioned datasets and methods of gathering dialogue promote interesting referential phenomena, including deictic reference and spatial relations with the two party instruction giver/follower paradigm being a popular format for soliciting this type of dialogue. In some respects, our offering could be seen as a combination and natural extension of the properties of prior corpora, in that we provide annotations over a dataset that incorporates a collaborative task-orientated dialogue, with action aligned utterances, and changing perspectives in a dynamic 3D situated environment.



Figure 1: Example of Minecraft Collaborative Builder Task from (Narayan-Chen et al., 2019)

## 3 The Minecraft Dialogue Corpus

### 3.1 The Data

The Minecraft Dialogue Corpus (Narayan-Chen et al., 2019) is a collection of conversations among human participants performing the Minecraft Collaborative Building Task, illustrated in Figure 1. In these conversations, two humans, *the Architect* – giving instructions and *the Builder* – executing them, cooperate to replicate a 3D structure in a 3D voxel based $11 \times 9 \times 11$ Cartesian coordinate based Minecraft world, with blocks of 6 different colours, accessed through Malmo (Johnson et al., 2016). The Architect has full observability over both the target environment and the builders actions, but may not directly effect change to the environment, only converse with the builder. The Builder is not able to see the target structure, but has an avatar that they can use to navigate and manipulate the environment with the goal of constructing the target structure. This motivates a multi turn situated dialogue exhibiting complex phenomena, including references to actions, objects and abstract shapes. But also sees these develop and sometimes change as the target structure is created. Figure 1 shows an example of this from the Builder's perspective. No conversational constraints are imposed; bidirectional communication, clarification etc. are all permissible. 509 dialogues were collected, with an average length of 30 utterances, for 15,926 utterances total, and 113,116 tokens.

### 3.2 Annotations

The richness of linguistic phenomena observed in the MDC has motivated several previous efforts to extend MDC with various forms of annotation.

Bonn et al. (2020) annotated 185 of the dialogues with an extended form of Abstract Meaning Representation (AMR) (Banarescu et al., 2013) that
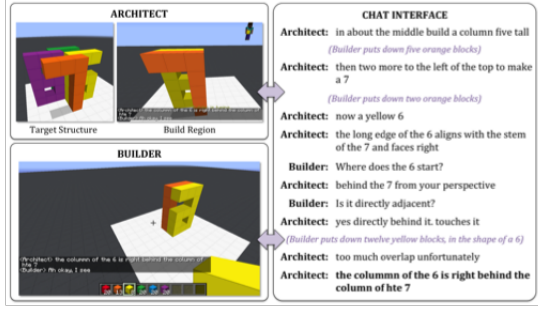
Figure 1: In the Minecraft Collaborative Building Task, the Architect (**A**) has to instruct a Builder (**B**) to build a target structure. A can observe **B**, but remains invisible to **B**. Both players communicate via a chat interface. (NB: We show **B**'s actions in the dialogue as a visual aid to the reader.)

Figure 2: References in a MDC dialogue (Narayan-Chen et al., 2019)

incorporates spatial relations. Bonial et al. (2021) apply a separate extension of AMR on MDC, with the Dialogue-AMR (Bonial et al., 2020) representation, a form of AMR that captures the illocutionary force of dialogue acts.

Thompson et al. (2024) produced the Minecraft Structured Dialogue Corpus (MSDC), that, through Segmented Discourse Representation Theory (SDRT) (Asher and Lascarides, 2003; Lascarides and Asher, 2007), gives a representation of the logical forms of the discourse in MDC, connecting dialogue and agent actions to deliver a macrostructure that links narrative arcs and discourse relations (e.g. corrections, confirmations, acknowledgements etc.).

## 4 Reference in the MDC

The MDC dialogues contain a rich variety of examples of reference in dialogue, as illustrated by the dialogue in Figure 2.

The dialogues contain plenty of deictic references to the visual situation (*the 6*, *the 7*), many of which are bridging references (e.g., *the middle* in the first utterance, *the top*). But perhaps the most distinctive aspect of these conversations is that the visual situation dynamically changes throughout the dialogues as new objects are built. For instance, the object that can be described as *the 7* is the result of the Builder's actions of putting down five orange blocks, then adding two orange blocks, as per Architect instructions.

Another distinctive feature of the corpus is that the objects that emerge as a result of the Builder actions can be referred to using their distinctive shapes. In Figure 2, the Builder first builds an object that looks like a 7, and can therefore be referred to as *the 7*, then an object that can be referred to as *the 6*. In Figure 3, the overall object can be referred as *the arch*, and the set of blocks in the middle can be referred to as *the bell*, even if those terms have not been previously introduced.

## 5 Our Annotation

### 5.1 Annotation Scheme

The annotation scheme for MDC-R is based on the scheme developed for the CODI/CRAC 2021 and 2022 shared tasks on anaphoric reference in dialogue (Yu et al., 2022a), which in turn is an extension of the ARRAU annotation scheme (Poesio et al., 2024a,b) to cover more dialogue phenomena. For MDC-R, we extended the CODI/CRAC scheme to cover more complex reference phenomena, building on the proposals in (Loáiciga et al., 2022). Reference to objects in a visual situation could be done with strong reliability in the TRAINS portion of the ARRAU corpus (Uryupina et al., 2020), in which the map was shown in an image separate from the annotation tool. (Loáiciga et al., 2022) obtained a Krippendorff $\alpha$ value of $0.55$ for the approach to annotation of references to the visual situation in which the image is shown on demand.

The key challenges we had to tackle were how to specify references in the visual scene, and how to handle the constantly changing dynamic environment. To address this challenge, we made some changes to the annotation tool MMAX2 (Müller and Strube, 2006). Similar to (Loáiciga et al., 2022), we support viewing an object labelled image of the current state of the world before each utterance, as seen by the builder. In our adaptation of the proposal by Loaiciga et al. to the Minecraft world situation, each block is given a unique alphanumerical tuple label, as shown in Figure 3. The annotators then enter the labels of the blocks referred to by the noun phrase in the 'Object' slot. Figures 4 and 5 show how references to such states are annotated. In Figure 4 we can see the state of the world before the utterance *on the block to the left of the green block*, which refers to the tower of 8 green blocks a0,a1,a2,a3,a4,a5,a6,a7. We choose this labelling format to create a concise sufficiently sized address space that could be easily identified and entered by annotators, without dependency on colour (which can sometimes be misjudged e.g. brown/orange in `B12-A26-C11-1522940682033` (Narayan-Chen et al., 2019)), or introducing further ambiguity. Figure 5 shows how annotators entered the labels for the blocks in the Object slot
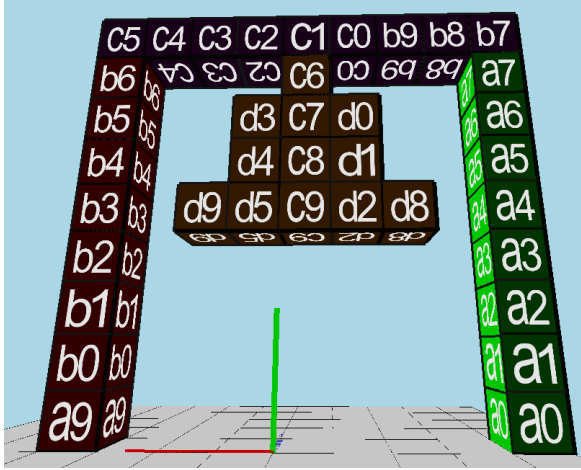
Figure 3: Labelling MDC with alphanumerical tuple indexes
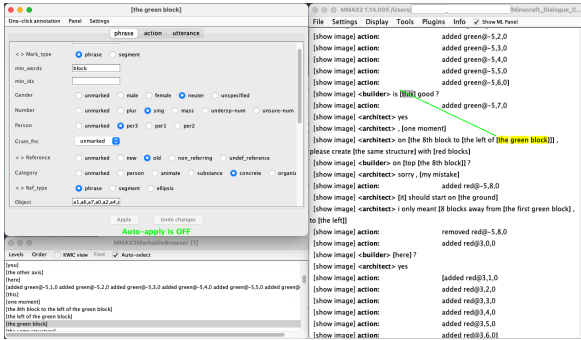


Figure 4: Reference annotation and grounding in MMAX2

of the referring mention. A reliability test with this revised scheme and guidelines (2 experienced annotators double-coded 5 dialogues containing a total of about 400 markables), resulting in a preliminary $\kappa = 0.43$. We expect this result to be a lower bound as many of the disagreements are due to the coders identifying slightly different sets of blocks.

## 5.2 Bounding Boxes

In the process of generating these labels, each label is tied to the 3D Cartesian coordinates of the block it is assigned to. This allows us to include a 2D bounding box that identifies the bounds of the block in the perspective based image. This makes MDC-R a versatile corpus, suitable both for versions of the REC task in which systems have to output labels, and for the more traditional version of the REC task in which systems have to output a bounding box (see Section 7 for an example of how MDC-R can be used in this way).

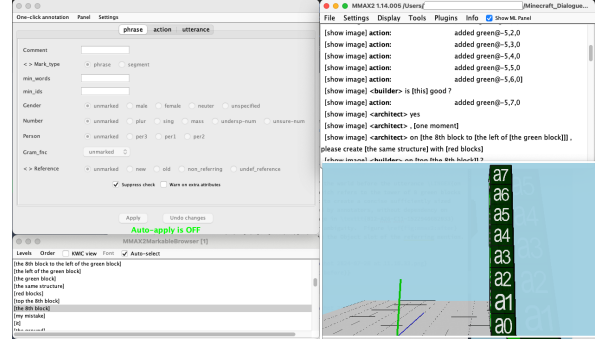These images highlight one of the fundamen-



Figure 5: Annotating references in MMAX2

tal properties of this corpus: that reference to objects in the scene translates to reference to *sets* of blocks that constitute these objects. This property means that a generalized notion of reference is required–one in which it is possible to refer not only to objects, but to sets of objects (He et al., 2023b; Hemanthage et al., 2024).

## 5.3 Annotation and Statistics

We selected at random a subset of 100 dialogues from MDC and converted them into a format suitable for labelling in MMAX2. The data was annotated by two professional linguists. Table 1 describes the corpus in terms of the frequency of specific reference properties that occur.

| Statistic | Count | Statistic | Count |
|---|---|---|---|
| Documents | 101 | Tokens | 29,174 |
| Utterances | 3,343 | Actions | 5,793 |
| Markables | 7,600 | Discourse old | 1960 |
| Bridging | 1,053 | Discourse Deixis | 500 |
| Plural | 24 | Ambiguous | 149 |

Table 1: Corpus Statistics

## 6 Reference in MDC-R: Observations

The MDC-R was created as a resource both to study linguistics reference in dynamically changing 3D settings and to develop models able to carry out this type of interpretation. We discuss in this Section how a selection of interesting linguistic phenomena in this Section end up being captured in the annotation corpus, and its suitability for modelling development in the next.

### 6.1 Dynamically changing states of the world

The state of the world described in the image in Figure 3 is the result of several rounds of interaction between Architect and Builder. Starting from a

completely empty 3D world, the Architect typically instructs the Builder to build the separate components one at a time–e.g., in the dialogue we have been discussing, the Architect first asks the Builder to create the (right) pillar shown in Figure 4, before going on asking to build the left pillar, then to put a beam on top of the two pillars, before moving on building the bell. All of these parts can be referred using terms such as *the pillar* or *the beam*, which in our scheme end up referring to sets of blocks.

## 6.2 The effect of perspective

The Architect and the Builder do not necessarily have the same perspective. Thus, the type of left/right confusions observed in the CUPS corpus (Dobnik et al., 2020) can be found in the MDC-R as well. One example is the following exchange in B36-A38-C66-1524262976017: *"if you walk to the right side of the yellow triangle"*, *"my right?"*, *"Now your left, sorry"*.

## 6.3 Grounding

In addition to misunderstandings due to perspective, a range of other types of misunderstandings can be observed, which means grounding these terms often requires rounds of clarifications. For instance, in Figure 4, the Architect initially uses the term *the green block* to refer to the pillar built by the Builder (and shown in Figure 5). The uncertain Builder intially interprets *the 8th block on top of the green block* as a reference to the block at the top of the pillar (block a7 in the picture), but then asks a clarification request.

Notice that a fully explicit annotation of this example would require keeping track of Architect and Builder's separate views of the common ground, as done e.g., in (Poesio et al., 2004), but this is not currently possible in MDC-R, so the annotators were instructed in these cases to mark the reference intended by the speaker of the utterance.

## 6.4 Intensional objects

Often, the Architect starts a round of building by telling the Builder what the objective is. For instance, in B29-A8-C1-152286385634, for which the target is Figure 3, the Architect starts by telling the Builder they are going to create, *"a gate with a bell all on one vertical plane"*, as illustrated in Figure 6. This NP is not a deictic reference. However, after the 'bell' is constructed, the Architect's reference to *the bell* is treated as deictic, and as bridging to the first mention (Figure 7).
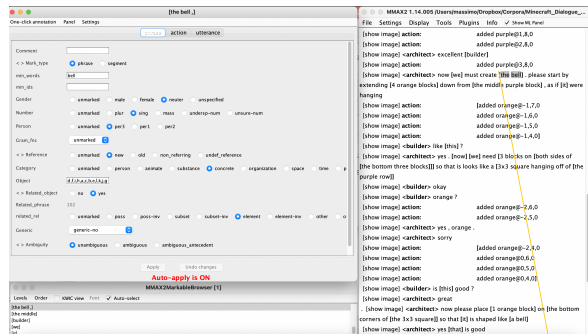


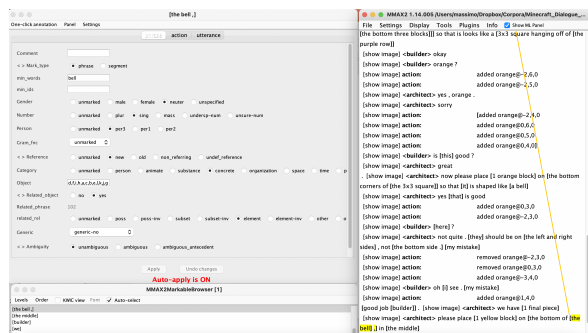Figure 6: Introducing the part of the construction to be built next, thee bell



Figure 7: Referring to the bell after it has been built.

# 7 Using the MDC-R for REC

To test the usefulness of the corpus for REC, we ran experiments predicting the bounding boxes which an expression refers to in the image of the state of the world associated with that utterance.

The primary challenge in this task is that the model must comprehend the dialogues provided so far and infer object attributes (e.g., shape, colour, position) based on the mention, and then detect the bounding box of the object in the image, resulting in cross-modal reasoning. A further challenge arises from the dynamic and embodied environment in which the data is collected. Specifically, the perspective of the screenshots varies depending on the architect's position during gameplay. This places a higher demand on the model to understand the voxelised world.

## 7.1 Data

To perform the REC task, we construct a dataset comprising 1,150 bounding boxes for blocks, 423 referring expressions, and 101 distinct scenarios. For each scenario, we extract the mention and dialogue history up to the point of the mention, providing a textual input of the scene. The corresponding screenshot at that specific round of dialogue serves

as the visual input. The ground truth label is defined as the merged bounding box encompassing all referenced blocks.

## 7.2 Models

We selected two models, Qwen2-VL (Wang et al., 2024) and MDETR (Kamath et al., 2021; He et al., 2023a; Liu et al., 2023b).

**Qwen2-VL**: An advanced multimodal large language model (MLLM), widely used in visual grounding, achieves state-of-the-art (SOTA) performance on RefCOCO, RefCOCO+, and RefCOCOg (Kazemzadeh et al., 2014; Mao et al., 2016; Nagaraja et al., 2016), which are widely regarded as the most classic REC datasets.

**MDETR**: This model follows a one-stage design for REC, predicting target regions directly without relying on object proposals. We adopt the implementation from He et al. (2023a), trained on the gRefCOCO (He et al., 2023a; Liu et al., 2023a) dataset, which supports multi-box prediction. To handle long referring expressions, we replace the original text encoder with Longformer (Beltagy et al., 2020) and fine-tune the modified model on gRefCOCO.

## 7.3 Implementation Details

We adopted two settings of the REC experiment, based on the distinct advantages of two baselines. **Classic REC**: We used Qwen2-VL to predict a single bounding box. The input and output format template are shown on Table 6. The input consists of both the system and the user perspectives, and the output is expected to contain a bounding box with a special token which could be parsed and extracted. **GREC**: We used the MDETR model to predict a set of bounding boxes.

## 7.4 Evaluation Metrics

To evaluate the result of Qwen2-VL, we employed the **Merge-box protocol** (Hemanthage et al., 2024): all ground-truth bounding boxes are merged into a single *minimal enclosing bounding box*. This merged bounding box serves as a unified ground truth label for all blocks in one image. To assess model performance, we compute the **Intersection over Union (IoU)** between the predicted bounding box and the merged ground-truth bounding box. Additionally, we report **Accuracy@0.25**, **Accuracy@0.5**, and **Accuracy@0.75**, where a prediction is considered correct if its IoU exceeds thresholds of 0.25, 0.5, or 0.75, respectively.

To evaluate the performance of MDETR (He et al., 2023a), we report the mean **F1 score** 1 at IoU thresholds of 0.25, 0.5, and 0.75. Specifically, given a set of predicted bounding boxes and ground-truth bounding boxes, a prediction is considered a true positive (TP) if it matches (IoU $\geq$ Threshold). Predicted bounding boxes that do not match any ground-truth box are counted as false positives (FP), while ground-truth boxes without any matching prediction are treated as false negatives (FN).

$$\text{F1} = \frac{2TP}{2TP+FP+FN} \quad (1)$$

In addition, for both GREC and classic REC, we calculate the **Mean IoU (mIoU)**, which represents the average IoU across all predictions in the dataset, providing a more comprehensive evaluation of overall performance.

## 7.5 Results

The results are given in Table 2. While performing much better than random,[2] we can see that the baseline model exhibits much lower results on this dataset than on other classical REC data. These results are noticeably lower than those typically reported in standard REC datasets (Kazemzadeh et al., 2014; Mao et al., 2016; Nagaraja et al., 2016) reported by Wang et al. (2024), underscoring the challenge of accurately predicting the target bounding box in Minecraft world screenshots using dialogue context and referring expressions.

| Metric | | | |
|---|---|---|---|
| mIoU (%) | Acc@0.25 (%) | Acc@0.5 (%) | Acc@0.75 (%) |
| 30.4 | 38.7 | 28.8 | 21.2 |

Table 2: Evaluation of Qwen2-VL performance, where Acc denotes Accuracy under Classic REC.

| Metric | | | |
|---|---|---|---|
| mIoU (%) | F1@0.25 (%) | F1@0.5 (%) | F1@0.75 (%) |
| 19.6 | 19.8 | 9 | 2.1 |

Table 3: F1 performance on MDETR under GREC.

Table 3 shows that MDETR performs poorly under the GREC setting. We hypothesized two main causes: (1) the domain gap in long-form dialogue and virtual Minecraft scenes poses a significant challenge due to the lack of aligned data;

---

[2]As a sanity check we tried random bounding box limits (mIoU = 5.1) and limits over the entire image (mIoU = 10.8).

(2) the model struggles with small object localization, as indicated by the sharp drop in F1 scores at higher IoU thresholds. A detailed example comparing the two evaluation settings is provided in Appendix 9.2.

Furthermore, Table 4 provides a detailed analysis of Qwen2-VL's performance across noun phrase (NP) categories, where each NP(n) denotes a referring expression for n block(s) referenced in one Minecraft world screenshot. The table shows that instances with smaller NP categories (i.e., NP1 to NP5) constitute the majority of the corpus.

| NP Category | mIoU (%) | Accuracy@0.5 (%) | Quantity |
|---|---|---|---|
| NP1 | 21.8 | 20.3 | 246 |
| NP2 | 27.6 | 22.7 | 66 |
| NP3 | 39.5 | 36.4 | 33 |
| NP4 | 36.3 | 31.8 | 22 |
| NP5 | 19.8 | 22.2 | 9 |
| NP6 | 76.7 | 85.7 | 7 |
| NP7 | 20.6 | 0.0 | 1 |
| NP8 | 49.2 | 53.8 | 13 |
| NP9 | 41.5 | 42.9 | 7 |
| NP10 | 71.2 | 100.0 | 1 |
| NP12 | 95.0 | 100.0 | 2 |
| NP13 | 4.8 | 0.0 | 2 |
| NP14 | 0.0 | 0.0 | 1 |
| NP15 | 27.4 | 0.0 | 1 |
| NP16 | 98.5 | 100.0 | 4 |
| NP18 | 65.8 | 100.0 | 2 |
| NP19 | 96.8 | 100.0 | 1 |
| NP20 | 98.4 | 100.0 | 3 |
| NP28 | 95.8 | 100.0 | 2 |

Table 4: mIoU scores and quantities for noun phrase (NP) analysis, where NP(n) represents one mention referring to n block(s) in the Minecraft world.

A counter-intuitive observation is that the model performs poorly on cases involving fewer blocks (NP<6) while achieving higher mIoU scores for cases with more blocks. One possible explanation is that bounding boxes covering a larger area (associated with higher NP values) may inherently yield a higher overlap, thus inflating the mIoU metric. This discrepancy highlights a limitation in the current evaluation framework.

In our corpus, each block is annotated with a unique ID and precise coordinates. Although these detailed annotations provide the potential to enhance evaluation accuracy, they simultaneously introduce additional challenges for model design.

In summary, the results emphasize the complexity of the REC task within the embodied visual context of Minecraft and the textual context of multi-round dialogue. Furthermore, our corpus motivates further exploration into refined evaluation methodologies and model architectures.

## 7.6 Error analysis

Manual inspection of the top 50 best- and worst-scoring examples from the results of Qwen2-VL revealed interesting patterns. Please see the Appendix for image examples.

**Good performance** occurs when: (1) the number of blocks in the scene is low, (2) the object consists of many blocks (covering a large part of the structure), (3) object is in foreground, (4) object is clearly separated from the rest of the scene in terms of colour/position (shadows or interleaving the object with other objects is detrimental).

**Bad performance** occurs when (1) there are many of distractor blocks that are not part of the object (both the dialog and scene are complex) (2) object is obscured fully or partially by either the edges of the scene or by other blocks, (3) the perspective of the image makes judging distances (depth perception) difficult. These observations are in line with expectations and illustrate the challenging nature of this dataset.

**Quantifying difficulty** The two sets of factors above correspond to situations where the ambiguity and complexity in the image is low or high, respectively. We aimed to expand this qualitative analysis by more explicitly quantifying aspects of the data that might influence performance. To this end, we calculate pearson's $\rho$ between the model performance and different data aspects (on N1 - N10, to avoid outliers). Results, provided in Table 5, confirm the findings the qualitative analysis, indicating that complexity and ambiguity, especially in terms of many distractor blocks and smaller referenced objects pose significant challenges. Furthermore, we reveal that dialog-related aspects also play a role, with longer dialogs leading to more challenging cases by producing more complex reference chains and repair structures.

| | | | |
|---|---|---|---|
| Num. object blocks (OB) | 0.21 | Num. utterances | -0.25 |
| Num. scene blocks (SB) | -0.29 | Num. architect utterances | -0.23 |
| Object/Scene ratio (OB/SB) | 0.46 | Num builder uterrances | -0.23 |
| Num. distractor blocks (SB - OB) | -0.34 | Num block remove actions | -0.17 |
| Gold bounding-box area | 0.32 | Num. block add actions | -0.26 |

Table 5: Pearson's $\rho$ between IoU and data aspects.

## 8 Conclusion

We introduced the Minecraft Dialogue Corpus with Reference, a multimodal dialogue dataset that

adds deictic and anaphoric reference annotations to it. The new version of the MDC provides both anaphoric and deictic reference labels for objects in the scenes, ids for the blocks, bounding boxes for the objects referred to, and a dynamic setting in which the actions of the speakers modify the world state (as well as speaker perspectives on the scene) as the dialog progresses. Moreover, MDC-R is the first dataset which provides information about the constituent parts of the referenced objects. As such, the dataset can be used for a wide variety of diverse tasks and will constitute both a valuable tool for studying collaborative dialogue as well as a very challenging benchmark for the reasoning abilities of next-generation multi-modal LLMs.

## Limitations

The key limitation of this work is that we were only able to annotate part of the original MDC in the time available. We hope to carry out more annotations in the future. A second limitation is that only one perspective is annotated, that of the speaker–as discussed in the paper, in some cases it would be good to explicitly annotate the different points of view of Architect and Builder. We also hope to address this limitation in future work.

## Acknowledgements

## References

Anne H Anderson, Miles Bader, Ellen Gurman Bard, Elizabeth Boyle, Gwyneth Doherty, Simon Garrod, Stephen Isard, Jacqueline Kowtko, Jan McAllister, Jim Miller, et al. 1991. The hcrc map task corpus. *Language and speech*, 34(4):351–366.

Nicholas Asher and Alex Lascarides. 2003. *Logics of conversation*. Cambridge University Press.

Laura Banarescu, Claire Bonial, Shu Cai, Madalina Georgescu, Kira Griffitt, Ulf Hermjakob, Kevin Knight, Philipp Koehn, Martha Palmer, and Nathan Schneider. 2013. Abstract meaning representation for sembanking. In *Proceedings of the 7th linguistic annotation workshop and interoperability with discourse*, pages 178–186.

Cristian-Paul Bara, Sky CH-Wang, and Joyce Chai. 2021. Mindcraft: Theory of mind modeling for situated dialogue in collaborative tasks. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, page 1112–1125. Association for Computational Linguistics.

Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.

Claire Bonial, Mitchell Abrams, David Traum, and Clare Voss. 2021. Builder, we have done it: evaluating & extending dialogue-amr nlu pipeline for two collaborative domains. In *Proceedings of the 14th International Conference on Computational Semantics (IWCS)*, pages 173–183.

Claire Bonial, Lucia Donatelli, Mitchell Abrams, Stephanie Lukin, Stephen Tratz, Matthew Marge, Ron Artstein, David Traum, and Clare Voss. 2020. Dialogue-amr: abstract meaning representation for dialogue. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 684–695.

Julia Bonn, Martha Palmer, Jon Cai, and Kristin Wright-Bettner. 2020. Spatial AMR: Expanded spatial annotation in the context of a grounded minecraft corpus. In *Proceedings of the 12th Conference on Language Resources and Evaluation (LREC 2020),*.

Susan E Brennan and Herbert H Clark. 1996. Conceptual pacts and lexical choice in conversation. *Journal of experimental psychology: Learning, memory, and cognition*, 22(6):1482.

Jon Chamberlain, Massimo Poesio, Udo Kruschwitz, et al. 2008. Phrase detectives: A web-based collaborative annotation game. In *Proceedings of the international conference on semantic systems (I-Semantics' 08)*, pages 42–49.

Zhihong Chen, Ruifei Zhang, Yibing Song, Xiang Wan, and Guanbin Li. 2023. Advancing visual grounding with scene knowledge: Benchmark and method. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15039–15049.

Herbert H Clark and Deanna Wilkes-Gibbs. 1986. Referring as a collaborative process. *Cognition*, 22(1):1–39.

Harm De Vries, Florian Strub, Sarath Chandar, Olivier Pietquin, Hugo Larochelle, and Aaron Courville. 2017. Guesswhat?! visual object discovery through multi-modal dialogue. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 5503–5512.

Simon Dobnik, John D Kelleher, and Christine Howes. 2020. Local alignment of frame of reference assignment in english and swedish dialogue. In *German Conference on Spatial Cognition*, pages 251–267. Springer.

Sebastiano Gigliobianco, Dimosthenis Kontogiorgos, and David Schlangen. 2024. Learning task-oriented dialogues through various degrees of interactivity. In *Proceedings of the 28th Workshop on the Semantics and Pragmatics of Dialogue*.

Arushi Goel, Basura Fernando, Frank Keller, and Hakan Bilen. 2023a. Semi-supervised multimodal coreference resolution in image narrations. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 4567–4578.

Arushi Goel, Basura Fernando, Frank Keller, and Hakan Bilen. 2023b. Who are you referring to? coreference resolution in image narrations. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5247–5257.

Shuting He, Henghui Ding, Chang Liu, and Xudong Jiang. 2023a. GREC: Generalized referring expression comprehension. *arXiv preprint arXiv:2308.16182*.

Shuting He, Henghui Ding, Chang Liu, and Xudong Jiang. 2023b. Grec: Generalized referring expression comprehension. *Preprint*, arXiv:2308.16182.

Bhathiya Hemanthage, Hakan Bilen, Phil Bartie, Christian Dondrup, and Oliver Lemon. 2024. Recantformer: Referring expression comprehension with varying numbers of targets. In *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*, pages 21784–21798.

Md Mofijul Islam, Reza Mirzaiee, Alexi Gladstone, Haley Green, and Tariq Iqbal. 2022. Caesar: An embodied simulator for generating multimodal referring expression datasets. *Advances in Neural Information Processing Systems*, 35:21001–21015.

Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Li Fei-Fei, C Lawrence Zitnick, and Ross Girshick. 2017. Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2901–2910.

Matthew Johnson, Katja Hofmann, Tim Hutton, and David Bignell. 2016. The malmo platform for artificial intelligence experimentation. In *IJCAI*, volume 16, pages 4246–4247.

Aishwarya Kamath, Mannat Singh, Yann LeCun, Gabriel Synnaeve, Ishan Misra, and Nicolas Carion. 2021. Mdetr-modulated detection for end-to-end multi-modal understanding. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 1780–1790.

Sahar Kazemzadeh, Vicente Ordonez, Mark Matten, and Tamara Berg. 2014. ReferItGame: Referring to objects in photographs of natural scenes. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 787–798, Doha, Qatar. Association for Computational Linguistics.

Julia Kiseleva, Ziming Li, Mohammad Aliannejadi, Shrestha Mohanty, Maartje ter Hoeve, Mikhail Burtsev, Alexey Skrynnik, Artem Zholus, Aleksandr Panov, Kavya Srinet, et al. 2022. Interactive grounded language understanding in a collaborative environment: IGLU 2021. In *NeurIPS 2021 Competitions and Demonstrations Track*, pages 146–161. PMLR.

Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Matt Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, et al. 2017. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474*.

Alex Lascarides and Nicholas Asher. 2007. Segmented discourse representation theory: Dynamic semantics with discourse structure. In *Computing meaning*, pages 87–124. Springer.

Chang Liu, Henghui Ding, and Xudong Jiang. 2023a. GRES: Generalized referring expression segmentation. In *CVPR*.

Chang Liu, Henghui Ding, and Xudong Jiang. 2023b. Gres: Generalized referring expression segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 23592–23601.

Sharid Loáiciga, Simon Dobnik, and David Schlangen. 2021. Annotating anaphoric phenomena in situated dialogue. In *Proceedings of the 1st Workshop on Multimodal Semantic Representations (MMSR)*, pages 78–88.

Sharid Loáiciga, Simon Dobnik, and David Schlangen. 2022. Anaphoric phenomena in situated dialog: A first round of annotations. In *Proceedings of the Fifth Workshop on Computational Models of Reference, Anaphora and Coreference*, pages 31–37.

Junhua Mao, Jonathan Huang, Alexander Toshev, Oana Camburu, Alan L Yuille, and Kevin Murphy. 2016. Generation and comprehension of unambiguous object descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 11–20.

Christoph Müller and Michael Strube. 2006. Multi-level annotation of linguistic data with mmax2. *Corpus technology and language pedagogy: New resources, new tools, new methods*, 3:197–214.

Varun K Nagaraja, Vlad I Morariu, and Larry S Davis. 2016. Modeling context between objects for referring expression understanding. In *Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11–14, 2016, Proceedings, Part IV 14*, pages 792–807. Springer.

Anjali Narayan-Chen, Prashant Jayannavar, and Julia Hockenmaier. 2019. Collaborative dialogue in minecraft. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 5405–5415.

Haruna Ogawa, Hitoshi Nishikawa, Takenobu Tokunaga, and Hikaru Yokono. 2020. Gamification platform for collecting task-oriented dialogue data. In

*Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 7084–7093.

M. Poesio, R. Delmonte, A. Bristot, L. Chiran, and S. Tonelli. 2004. The VENEX corpus of anaphoric information in spoken and written Italian. Unpublished. Available online at http://cswww.essex.ac.uk/staff/poesio/publications/VENEX04.pdf.

Massimo Poesio, Maris Camilleri, Paloma Carretero Garcia, and Ron Artstein. 2024a. *The ARRAU 3 Annotation Manual*, v. 1.1 edition. Queen Mary University of London.

Massimo Poesio, Maris Camilleri, Paloma Carretero Garcia, Juntao Yu, and Mark-Christoph Müller. 2024b. The arrau 3.0 corpus. In *Proceedings of the 5th Workshop on Computational Approaches to Discourse (CODI )*, pages 127–138. Association for Computational Linguistics.

Massimo Poesio, Jon Chamberlain, Udo Kruschwitz, Silviu Paun, Alexandra Uma, and Juntao Yu. 2019. A crowdsourced corpus of multiple judgments and disagreement on anaphoric interpretation. In *Proc. of NAACL*, page 1778–1789, Minneapolis. Association for Computational Linguistics (ACL).

Massimo Poesio, Jon Chamberlain, Udo Kruschwitz, Livio Robaldo, and Luca Ducceschi. 2013. Phrase detectives: Utilizing collective intelligence for internet-scale language resource creation. *ACM Transactions on Intelligent Interactive Systems*, 3(1).

Massimo Poesio, Juntao Yu, Silviu Paun, Abdulrahman Aloraini, Pengcheng Lu, Janosch Haber, and Derya Cokal. 2023. Computational models of anaphora. *Annual Review of Linguistics*, 9:561–587.

Xavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Torralba. 2018. Virtualhome: Simulating household activities via programs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8494–8502.

Yuankai Qi, Qi Wu, Peter Anderson, Xin Wang, William Yang Wang, Chunhua Shen, and Anton van den Hengel. 2020. Reverie: Remote embodied visual referring expression in real indoor environments. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9982–9991.

David Schlangen. 2023. Dialogue games for benchmarking language understanding: Motivation, taxonomy, strategy. *arXiv preprint arXiv:2304.07007*.

Mohit Shridhar, Jesse Thomason, Daniel Gordon, Yonatan Bisk, Winson Han, Roozbeh Mottaghi, Luke Zettlemoyer, and Dieter Fox. 2020. Alfred: A benchmark for interpreting grounded instructions for everyday tasks. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 10740–10749.

Philipp Spanger, Masaaki Yasuhara, Ryu Iida, Takenobu Tokunaga, Asuka Terai, and Naoko Kuriyama. 2012. Rex-j: Japanese referring expression corpus of situated dialogs. *Language Resources and Evaluation*, 46:461–491.

Alane Suhr, Claudia Yan, Charlotte Schluger, Stanley Yu, Hadi Khader, Marwa Mouallem, Iris Zhang, and Yoav Artzi. 2019. Executing instructions in situated collaborative interactions. *arXiv preprint arXiv:1910.03655*.

Arthur Szlam, Jonathan Gray, Kavya Srinet, Yacine Jernite, Armand Joulin, Gabriel Synnaeve, Douwe Kiela, Haonan Yu, Zhuoyuan Chen, Siddharth Goyal, et al. 2019. Why build an assistant in minecraft? *arXiv preprint arXiv:1907.09273*.

Kate Thompson, Julie Hunter, and Nicholas Asher. 2024. Discourse structure for the minecraft corpus. In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 4957–4967.

Jack Urbanek, Angela Fan, Siddharth Karamcheti, Saachi Jain, Samuel Humeau, Emily Dinan, Tim Rocktäschel, Douwe Kiela, Arthur Szlam, and Jason Weston. 2019. Learning to speak and act in a fantasy text adventure game. *arXiv preprint arXiv:1903.03094*.

Olga Uryupina, Ron Artstein, Antonella Bristot, Federica Cavicchio, Francesca Delogu, Kepa J Rodriguez, and Massimo Poesio. 2020. Annotating a broad range of anaphoric phenomena, in a variety of genres: the arrau corpus. *Natural Language Engineering*, 26(1):95–128.

Kees van Deemter. 2016. *Computational Models of Referring*. The MIT Press.

Luis Von Ahn. 2006. Games with a purpose. *Computer*, 39(6):92–94.

Peng Wang, Shuai Bai, Sinan Tan, Shijie Wang, Zhihao Fan, Jinze Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Yang Fan, Kai Dang, Mengfei Du, Xuancheng Ren, Rui Men, Dayiheng Liu, Chang Zhou, Jingren Zhou, and Junyang Lin. 2024. Qwen2-vl: Enhancing vision-language model's perception of the world at any resolution. *arXiv preprint arXiv:2409.12191*.

Peng Wang, Dongyang Liu, Hui Li, and Qi Wu. 2020. Give me something to eat: Referring expression comprehension with commonsense knowledge. In *Proceedings of the 28th ACM International Conference on Multimedia*, pages 28–36.

Haoxuan You, Rui Sun, Zhecan Wang, Kai-Wei Chang, and Shih-Fu Chang. 2022. Find someone who: Visual commonsense understanding in human-centric grounding. In *Findings of the Association for Computational Linguistics: EMNLP 2022*, pages 5444–5454, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.

Juntao Yu, Sopan Khosla, Ramesh Manuvinakurike, Lori Levin, Vincent Ng, Massimo Poesio, Michael Strube, and Carolyn Rosé. 2022a. The CODI-CRAC 2022 shared task on anaphora, bridging, and discourse deixis in dialogue. In *Proceedings of the CODI-CRAC 2022 Shared Task on Anaphora, Bridging, and Discourse Deixis in Dialogue*, pages 1–14.

Xintong Yu, Hongming Zhang, Ruixin Hong, Yangqiu Song, and Changshui Zhang. 2022b. Vd-pcr: Improving visual dialog with pronoun coreference resolution. *arXiv preprint arXiv:2205.14693*.

Xintong Yu, Hongming Zhang, Yangqiu Song, Yan Song, and Changshui Zhang. 2019. What you see is what you get: Visual pronoun coreference resolution in dialogues. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 5123–5132.

Sina Zarrieß, Julian Hough, Casey Kennington, Ramesh Manuvinakurike, David DeVault, Raquel Fernández, and David Schlangen. 2016. PentoRef: A corpus of spoken references in task-oriented dialogues. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC'16)*, pages 125–131, Portorož, Slovenia. European Language Resources Association (ELRA).

# 9 Appendix

## 9.1 Prompt

| Input format | |
|---|---|
| System | You are a helpful agent to understand a mention in the last sentence of a dialogue in a Minecraft scenario, and detect the bounding boxes of the target block(s). <\|dialogue_start\|><\|dialogue_end\|> is the dialogue. <\|mention_start\|><\|mention_end\|> is the mention. You need to output the bounding box in <\|box_start\|><\|box_end\|>. |
| User | <\|vision_start\|>picture<\|vision_end\|> <\|dialogue_start\|>dialogue<\|dialogue_end\|> detect the bounding box of <\|mention_start\|>mention<\|mention_end\|>. |
| **Output format** | |
| Agent | ...<\|box_start\|>bounding_box<\|box_end\|> |

Table 6: The input and output formats. In the input, an additional instruction from the system perspective is added to guide the model on how to understand the task.

## 9.2 Example for two distinct evaluation settings

We present model predictions under two different evaluation settings. The gold bounding box is shown in green, and the predicted one in red. As shown in Figure 8, under the REC setting, the model easily recognizes the large 3×3 blue square, as the target occupies a prominent area in the image. In contrast, as illustrated in Figure 9, the GREC setting requires the model to precisely locate each individual block that forms the square. This significantly increases the difficulty, demanding a deeper understanding of the expression and a precise reference to the dialogue to consider all mentioned blocks.
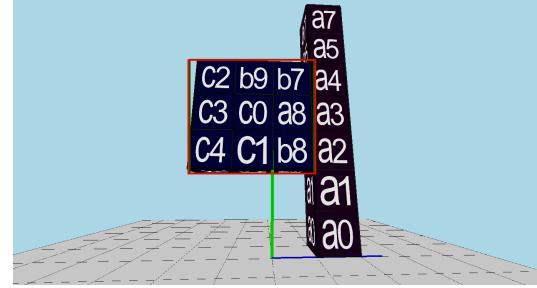


Figure 8: Example case of recognizing merged blocks by Qwen2-VL under the REC setting.
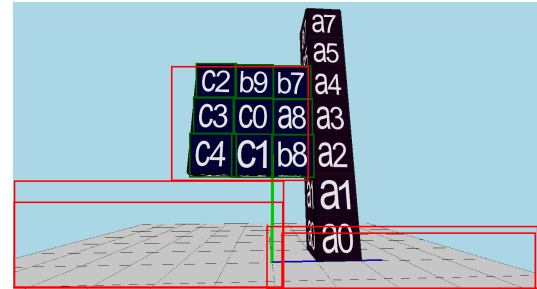


Figure 9: Example case of recognizing independent blocks by MDETR under the GREC setting.

## 9.3 Example model outputs

This section provides images of the model outputs for particularly easy (Figure 10, Figure 11, Figure 12, Figure 13, and Figure 14) or difficult (Figure 15, Figure 16, Figure 17) cases. The gold bounding box is green, while the predicted one is red.
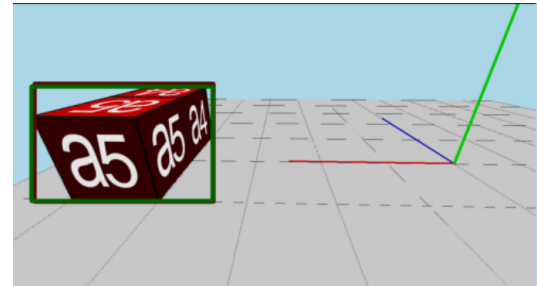


Figure 10: Example case where there are only a few blocks in the scene, making the task less challenging.
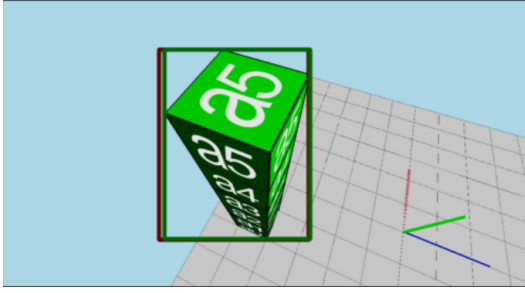
Figure 11: Example of a case where the object covers a significant part of the structure (often the whole structure). This makes the example less challenging even if there are many blocks involved.
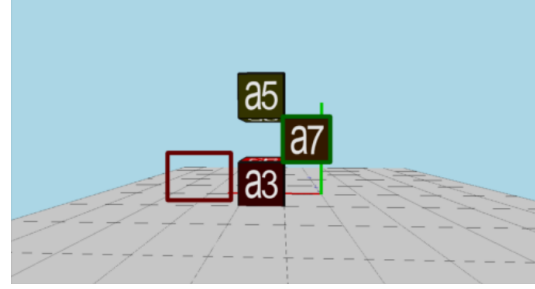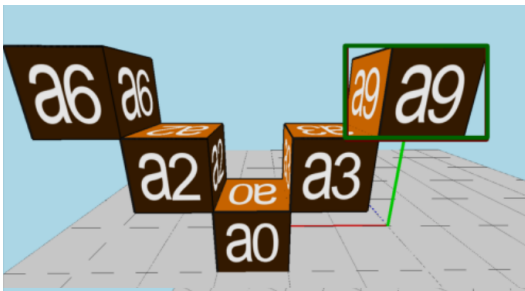


Figure 12: Example case where the object block is in the foreground and easy to see.



Figure 13: Example case where the object is clearly separated by colour.



Figure 14: Example case where the object is clearly separated in terms of space.



Figure 17: Example case where the perspective makes depth perception very difficult, making the task harder.



Figure 15: Example case where there is a large number of distractor blocks which do not belong to the referenced object. These cases are extremely challenging.
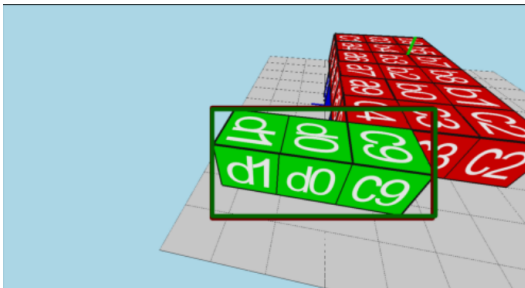


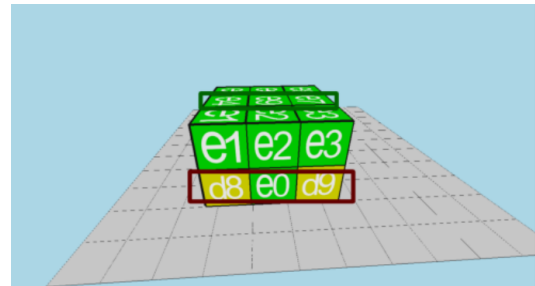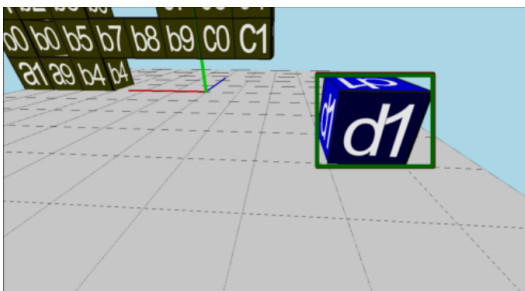Figure 16: Example case where the object is partially obscured by. These are often very challenging.