

Function	Type	Example Utterances
Information Providing	Inform (Init.)	I have very dark small dot in the center
	Inform (Cont.)	It also has a small light grey one further down from the group
	Agreement	Yes I have one like that. / same here.
	Agreement (Strong)	Exactly! / perfect. mine too.
	Agreement (Partial)	not sure its the one / more of a line.
Information Seeking	Disagreement	Yes, but the small is medium dark, not completely black
		I don't have that one. / mine are not in those locations.
	Question (Prop.)	the middle one is the darkest of the 3?
	Question (Set)	where is it in relation to the large med grey?
	Question (Choice)	Which should we choose? / the black or the grey?
Commissives	Question (Check)	It's the darkest dot in the circle, right?
	Offer	lets click the upper left one that's bigger and darker gray
Directives	Request	tell me about your tiniest dot? / pick one at the bottom
	Suggestion	Please describe it in relation to other dots in the circle

Table 3: Illustrative utterances in the dataset, grouped by the *task dimension* of communicative functions (?)

*mension* of communicative functions (?), including information transfer functions (Information Providing/Seeking) and action discussion functions (Commissives/Directives). With additional annotations, our dataset can be extended for other dialogue tasks, such as dialogue act recognition.

**Discourse Level Phenomena** Naturally, we found many coreference and anaphoric expressions in our dataset. *Coreference resolution* is the task of mapping mentions of entities to their referents. In our dataset, we found two characteristics that complicate this task. First, due to continuous and partially-observable context, mentions are usually ambiguous and referents may be missing. Thus players must keep track of various possibilities and investigate them through interaction. Secondly, players often use *groupings* (such as *three in a line*, *a cluster of 4 dots*) where mentions refer to *sets* of entities. This strategy could be effective but adds complexity to coreference resolution.

On the other hand, *anaphoric relation* is the relation between a mention and following mentions which refer to the previous mention. This can occur both within utterances (“a medium size black one, with a very light slightly smaller one to *it's* left”) and across utterances (“Does *the lighter dot* appear to be slightly larger?”). Similar to coreference resolution, this is a challenging subtask of common grounding at the discourse level which can be studied on our dataset.

## 5 Experiments

### 5.1 Experiment Overview

In this experiment, we formulate a natural language understanding task based on *target selection*: specifically, we try to predict which target a player selected, given the player's observation and the corresponding dialogue. This is an essential subtask of collaborative referring, where players choose their final selection based on the created common ground. Since the number of entities in view is fixed at 7, we can formulate this as a simple classification problem. Our

baseline models are kept as simple as possible, with minimal preprocessing and hyperparameter tuning.

### 5.2 Methods

Two main components of the models are as follows:

**Context Embedding** The structured form of the context is represented as a 28 dimensional real-valued vector, where each of the 7 observable entities is represented as a 4 dimensional vector (x-value, y-value, size, color). Each dimension is further normalized in the range of (-1,1).

The simplest way to embed context is to directly apply a multi-layered perceptron (MLP) over the context vector. However, without feature engineering this simple approach may have difficulty in capturing relevant information, such as relations between entities. Therefore, in the second approach we use the Relation Network module (?) to create additional features about relations between entities. Specifically, we embed each combination of the entities (total of 21 pairs) with a shared MLP and append the sum of these vectors as additional input.

**Dialogue Embedding** Utterances are all tokenized and lowercased, and tokens which occur less than 10 times are treated as a unique *unknown* token. We insert tokens which represent *speaker id* to each utterance at the beginning, and another token to indicate the end of the dialogue. Then, we embed these tokens with a shared MLP and run a bidirectional GRU (?) over the embedded tokens. Finally, we take the last output of the bi-GRU as the final representation of the dialogue.

For prediction, we simply concatenate the context and dialogue embeddings and run another MLP. However, as we've seen in Section 4, there are nonlinguistic selection bias in our dataset, so it is possible to make predictions without dialogue embeddings. Therefore, we also train models to predict only from the context embeddings using MLP.

	Full	Uncorrelated	Success Only
Random	14.28	14.28	14.28
Context Only (MLP)	$27.90 \pm 0.6$	28.74	29.59
Context Only (RN)	$31.94 \pm 0.9$	30.22	32.40
Context + Dialogue (MLP)	$40.27 \pm 1.3$	40.89	43.82
Context + Dialogue (RN)	$43.09 \pm 0.8$	44.00	49.44
Humans	-	82.50	90.79

Table 4: Results of the target selection experiment. Models are trained 10 times initialized with different seeds for the Full testset, and the models with best validation loss are used for the additional testset results (Uncorrelated and Success Only).

Following common practice, we split the dataset into training, validation and test set with a proportion of 8:1:1, and all models are tuned on the validation set. The loss function is calculated using cross entropy. All components of the neural networks consist of single layer with 128 hidden units, and dropout rate of 0.5 is applied at each layer to avoid overfitting. All parameters are initialized uniformly within the range of  $(-0.01, 0.01)$ . Models are trained with the Adam optimizer (?) with initial learning rate of 0.001, and we clip gradients whose  $L^2$  norm is greater than 0.1. The experiment is run 10 times initialized with different seeds, and we report the mean and standard deviation of the selection accuracies on the full testset.

For further analyses, models with the best validation loss in the previous experiment are also tested on two variants of the testset. First, we create an uncorrelated testset by randomly removing one from each correlated pair in the current testset (same dialogue but different context). Secondly, we further removed dialogues where players failed to coordinate on the same entity from the uncorrelated testset, since this may affect target selection performance. The statistical significance of the results for each pair of methods are tested on the uncorrelated testset using paired student's t-test. Finally, we take 100 random samples from the uncorrelated testset (including 76 successful) to report human performance based on average accuracy of two annotators.

### 5.3 Results

We show the results of our experiment in Table 4. As we can see, models trained only with the context embeddings perform significantly better than random ( $p$ -value  $< 10^{-7}$ ). This verifies that we can indeed take advantage of selection bias to make better predictions.

In addition, we found that embedding context with Relation Network consistently outperforms MLP, but not at a statistically significant level ( $p$ -value  $> 0.1$ ). Therefore, the simplest strategy of using MLP works decently, but a better architecture may improve the overall performance.

Finally, models trained with both context and dialogue embeddings significantly outperform models trained only with the context embeddings ( $p$ -value  $< 10^{-9}$ ). This indicates that even our simplest models can learn to ground linguistic meanings based on the context to make better predictions. When the testset only includes successful cases, models perform better but human performance improves even more achieving over 90% accuracy. Overall, our target selection task is challenging due to the complexity of common

grounding, and we still have a huge room for improvement.

## 6 Conclusion and Future Work

The main contributions can be summarized as follows:

- We proposed a simple and general idea of incorporating continuous and partially-observable context to the dialogue tasks, which makes common grounding difficult in a natural way.
- Following this idea, we formulated a novel dialogue task based on collaborative referring which enables clear evaluation and analysis of complex models.
- We collected a largescale dataset of 6,760 dialogues, which fulfills essential requirements of natural language corpora and will be publicly available online.
- Our analysis of the dataset verified the difficulty of common grounding and revealed various phenomena that need to be considered.
- We evaluated and analyzed simple baseline models on an important subtask of collaborative referring and showed that there is still room for further improvement.

In future work, we will evaluate and analyze dialogue models based on our task, especially to identify the current limitations of end-to-end approaches in terms of common grounding. Models can be trained in a variety of ways, including supervised learning, reinforcement learning with humans, and reinforcement learning based on *self-play* (?). Overall, we expect our task to be a fundamental testbed for developing dialogue systems with sophisticated common grounding abilities.

## Acknowledgements

This work was supported by JSPS KAKENHI Grant Numbers 16K12546, 18H03297.

Nihil libero saepe molestias voluptatum consequatur blanditiis, molestiae nostrum sed quidem accusantium, soluta dolores tenetur sunt sequi similique reiciendis, dolorum repudiandae non consequatur labore porro laudantium, ratione ea cupiditate alias. Nulla voluptas ipsa placeat necessitatibus, sint fuga facilis a recusandae rerum officiis, voluptatem quas culpa a cum facilis quos omnis placeat atque? Animi laboriosam libero suscipit impedit laborum ipsam nostrum blanditiis facilis, expedita eligendi voluptatem tenetur cumque tempora voluptates quo reprehenderit, eligendi reiciendis neque ad magni veritatis explicabo optio inventore corrupti. Qui fugit sequi laudantium maiores sint aspernatur voluptate natus porro in ea, eaque illum molestiae hic totam laboriosam adipisci alias, architecto