# CoDraw: Collaborative Drawing as a Testbed for Grounded Goal-driven Communication

**Jin-Hwa Kim**[*]
SK T-Brain[†]
jnhwkim@sktbrain.com

**Nikita Kitaev**[*]
University of California, Berkeley[†]
kitaev@cs.berkeley.edu

**Xinlei Chen, Marcus Rohrbach**
Facebook AI Research
{xinleic,mrf}@fb.com

**Byoung-Tak Zhang**
Seoul National University
btzhang@bi.snu.ac.kr

**Yuandong Tian**
Facebook AI Research
yuandong@fb.com

**Dhruv Batra & Devi Parikh**
Georgia Institute of Technology, Facebook AI Research
{parikh,dbatra}@gatech.edu

## Abstract

In this work, we propose a goal-driven collaborative task that combines language, perception, and action. Specifically, we develop a *Co*llaborative image-*Draw*ing game between two agents, called *CoDraw*. Our game is grounded in a virtual world that contains movable clip art objects. The game involves two players: a *Teller* and a *Drawer*. The Teller sees an abstract scene containing multiple clip art pieces in a semantically meaningful configuration, while the Drawer tries to reconstruct the scene on an empty canvas using available clip art pieces. The two players communicate with each other using natural language. We collect the CoDraw dataset of ~10K dialogs consisting of ~138K messages exchanged between human players. We define protocols and metrics to evaluate learned agents in this testbed, highlighting the need for a novel *crosstalk* evaluation condition which pairs agents trained independently on disjoint subsets of the training data. We present models for our task and benchmark them using both fully automated evaluation and by having them play the game live with humans.

## 1 Introduction

Building agents that can interact with humans in natural language while perceiving and taking actions in their environments is one of the fundamental goals in artificial intelligence. To this end, it will be necessary to *ground* language into perception and action (Harnad, 1990; Barsalou, 1999),

---

[*] The first two authors contributed equally to this work.
[†] Work performed while the authors were interns at Facebook AI Research.
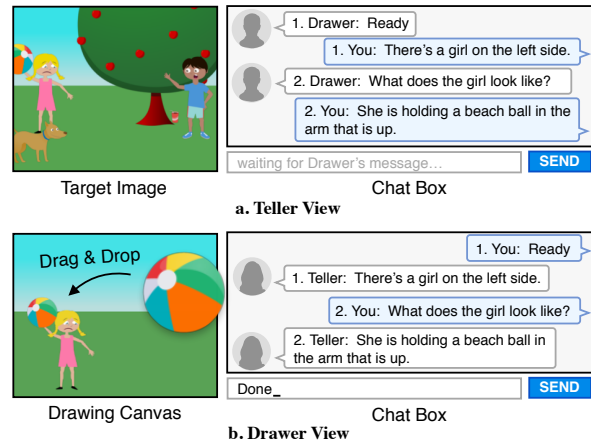


Figure 1: Overview of the proposed Collaborative Drawing (CoDraw) task. The game consists of two players – Teller and Drawer. The Teller sees an abstract scene, while the Drawer sees an initially empty canvas. Both players need to collaborate and communicate so that the Drawer can drag and drop clip art objects to reconstruct the target scene that is only visible to the Teller.

where, *e.g.*, nouns are connected to percepts and verbs relate to actions taken in an environment. Some approaches judge machine understanding of language purely based on the ability to mimic particular human utterances, but this has limitations: there are many ways to express roughly the same meaning, and conveying the correct information is often more important than the particular choice of words. An alternative approach, which has recently gained increased prominence, is to train and evaluate language capabilities in an *interactive* setting, where the focus is on *successfully communicating information* that an agent must share in order to achieve its goals.

In this paper, we propose the Collaborative Drawing (CoDraw) task, which combines grounded language understanding and learning effective goal-driven communication into a single, unified testbed. This task involves perception, communication, and actions in a partially observable environment. As shown in Figure 1, our game is grounded in a virtual world constructed from clip art objects (Zitnick et al., 2013; Zitnick and Parikh, 2013). Two players, Teller and Drawer, play the game. The Teller sees an abstract scene consisting of clip art objects in a semantically meaningful configuration, while the Drawer sees a drawing canvas that is initially empty. The goal of the game is to have both players communicate so that the Drawer can reconstruct the image of the Teller, without ever seeing it.

Our task requires effective communication because the two players cannot see each other's scenes. The Teller must describe the scene in sufficient detail for the Drawer to reconstruct it, which will necessitate grounded language. Moreover, the Drawer will need to carry out a series of actions from a rich action space to position, orient, and resize all of the clip art pieces required for the reconstruction. Note how clip art pieces form a representation that is perceived visually by humans but is easy to manipulate in a structured manner, in contrast to lower-level pixel-based image representations. The performance of a pair of agents is judged based on the quality of reconstructed scenes, where high-quality reconstructions result from successful communication.

We collect a CoDraw dataset[1] of ∼10K variable-length dialogs consisting of ∼138K messages with the drawing history at each step of the dialog. We also define a similarity metric for clip art scenes, which allows us to automatically evaluate the effectiveness of agent communication at the end of a dialog and at intermediate states. We evaluate several Drawer and Teller models[2] automatically as well as by pairing them with humans, and show that long-term planning and contextual reasoning are key challenges of the CoDraw task.

As we developed models and protocols for CoDraw, we found it critical to train the Teller and the Drawer separately on disjoint subsets of the training data. Otherwise, the two machine agents may conspire to successfully achieve the goal while communicating using a shared "codebook" that bears little resemblance to natural language. We call this separate-training, joint-evaluation protocol *crosstalk*, which prevents learning of mutually agreed upon codebooks, while still checking for goal completion at test time. We highlight crosstalk as one of our contributions, and believe it can be generally applicable to other related tasks (Sukhbaatar et al., 2016; Foerster et al., 2016; de Vries et al., 2016; Das et al., 2017b; Lewis et al., 2017).

## 2 Related work

**Language grounded in environments.** Learning language games in a grounded environment has been studied recently (Wang et al., 2016, 2017). While language in these works is tied to actions that modify the environment, the tasks do not involve multiple agents that need to cooperate. Other work on grounded instruction following relies on datasets of pre-generated action sequences annotated with human descriptions, rather than using a single end goal (Long et al., 2016). Generation models for these tasks are only evaluated based on their ability to describe an action sequence that is given to them (Fried et al., 2018a), whereas Teller models for CoDraw also need to select in a goal-driven manner the action sequence to describe to the Drawer. Language grounding has been studied for robot navigation, manipulation, and environment mapping (Tellex et al., 2011; Mei et al., 2015; Daniele et al., 2016). However, these works manually pair each command with robot actions and lack end-to-end training (Tellex et al., 2011), dialog (Mei et al., 2015; Daniele et al., 2016), or both (Walter et al., 2014). Compared to work on navigation (Vogel and Jurafsky, 2010; Anderson et al., 2018; Fried et al., 2018b) where an agent must follow instructions to move itself in a static environment, CoDraw involves a structured action space for manipulating clip art pieces to form a semantically meaningful configuration.

**End-to-end goal-driven dialog.** Traditional goal-driven agents are often based on 'slot filling' (Lemon et al., 2006; Wang and Lemon, 2013; Yu et al., 2015), in which the structure of the dialog is pre-specified but the individual slots are replaced by relevant information. Recently, end-to-end neural models are also proposed for goal-driven dialog (Bordes et al., 2017; Li et al.,

---

2017a,b; He et al., 2017), as well as goal-free dialog or 'chit-chat' (Shang et al., 2015; Sordoni et al., 2015; Vinyals and Le, 2015; Li et al., 2016; Dodge et al., 2016). Unlike CoDraw, in these approaches, symbols in the dialog are not grounded into visual objects.

**Emergent communication.** Building on the seminal works by Lewis (1969, 1975), a number of recent works study cooperative games between agents where communication protocols emerge as a consequence of training the agents to accomplish shared goals (Sukhbaatar et al., 2016; Foerster et al., 2016). These methods have typically been applied to learn to communicate small amounts of information, rather than the complete, semantically meaningful scenes used in the CoDraw task. In addition, the learned communication protocols are usually not natural (Kottur et al., 2017) or interpretable, whereas the CoDraw task is designed to develop agents that use human language.

**Language and vision.** The proposed CoDraw game is related to several well-known language and vision tasks that study grounded language understanding (Karpathy and Fei-Fei, 2015; Donahue et al., 2015; de Vries et al., 2016). For instance, in contrast to image captioning (Vinyals et al., 2017; Xu et al., 2015; Chen and Zitnick, 2015; Lu et al., 2017), visual question answering (Antol et al., 2015; Zhang et al., 2016; Goyal et al., 2016; Gao et al., 2015; Krishna et al., 2017; Malinowski and Fritz, 2014; Ren et al., 2015; Tapaswi et al., 2016; Yu et al., 2015; Zhu et al., 2016) and recent embodied extensions (Das et al., 2018), CoDraw involves multiple rounds of interactions between two agents. Both agents hold their own partially observable states and may need to build a model of their partner's state to collaborate effectively. Compared to past work on generating abstract scenes from single captions (Zitnick et al., 2013), scenes in CoDraw are reconstructed over multiple rounds, and the task requires Teller models to generate coherent and precise descriptions over the course of a full dialog. Compared to visual dialog (Das et al., 2017a,b; Strub et al., 2017; Mostafazadeh et al., 2017) tasks, agents need to additionally cooperate to change the environment with actions (*e.g.*, move pieces around). Thus, the agents have to possess the ability to adapt and hold a dialog about partially-constructed scenes that will occur over the course of their interactions. In addition, we also want to highlight that CoDraw

has a well-defined communication goal, which facilitates objective measurement of success and enables end-to-end goal-driven learning.

# 3 CoDraw task and dataset

In this section, we first detail our task, then present the CoDraw dataset, and finally propose a scene similarity metric which allows automatic evaluation of the reconstructed and original scene.

## 3.1 Task

**Abstract scenes.** To enable people to easily draw semantically rich scenes on a canvas, we leverage the Abstract Scenes dataset of Zitnick et al. (2013) and Zitnick and Parikh (2013). This dataset consists of 10,020 semantically consistent scenes created by human annotators. An example scene is shown in the left portion of Figure 1. Most scenes contain 6 objects (min 6, max 17, mean 6.67). These scenes depict children playing in a park, and are made from a library of 58 clip arts, including a boy (Mike) and a girl (Jenny) in one of 7 poses and 5 expressions, and various other objects including trees, toys, hats, animals, food, etc. An abstract scene is created by dragging and dropping multiple clip art objects to any $(x, y)$ position on the canvas. Spatial transformations can be applied to each clip art, including sizes (small, normal, large) and two orientations (facing left or right). The clip art serve simultaneously as a high-level visual representation and as a mechanism by which rich drawing actions can be carried out.

**Interface.** We built a drag-and-drop interface based on the Visual Dialog chat interface (Das et al., 2017a) (see Figures 5 and 6 in Appendix A for screen shots of the interface). The interface allows real-time interaction between two people. During the conversation, the Teller describes the scene and answers any questions from the Drawer on the chat interface, while Drawer "draws" or reconstructs the scene based on the Teller's descriptions and instructions. Each side is only allowed to send one message at a time, and must wait for a reply before continuing. The maximum length of a single message is capped at 140 characters: this prevents excessively verbose descriptions and gives the Drawer more chances to participate in the dialog by encouraging the Teller to pause more frequently. Both participants were asked to submit the task when they are both confident that Drawer has accurately reconstructed the scene of Teller.
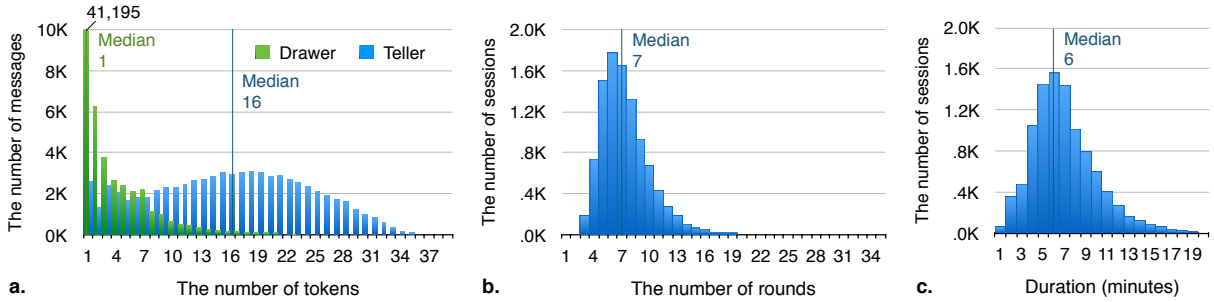
Figure 2: Statistics of the CoDraw dataset. **(a)** The distribution of the number of tokens in Teller (blue) and Drawer (green) messages. Note that the number of single-token messages by Drawers is 41,195 (62.06%). The median token counts for Tellers and Drawers are 16 and 1, respectively. **(b)** The distribution of the numbers of conversation rounds. The median is 7 rounds. **(c)** The distribution of the duration of dialog sessions. The median is 6 minutes.

To focus the natural language on the high-level semantics of the scene rather than instructions calling for the execution of low-level clip art manipulation actions, the Teller is not able to observe the Drawer's canvas while communicating.

### 3.2 Dataset

We collect 9,993[3] dialogs where pairs of people complete the CoDraw task, consisting of one dialog per scene in the Abstract Scenes dataset. The dialogs contain of a total of 138K utterances and include snapshots of the intermediate state of the Drawer's canvas after each round of each conversation. See Section 5 for a description of how we split the data into training, validation, and test sets. **Messages.** Figure 2a shows the distribution of message lengths for both Drawers and Tellers. The message length distribution for the Drawer is skewed toward 1 with passive replies like *"ok"*, *"done"*, etc. There does exist a heavy tail, which shows that Drawers ask clarifying questions about the scene like *"where is trunk of second tree, low or high"*. On the other hand, Teller utterances have a median length of 16 tokens and a vocabulary size of 4,555. Due to the limited number of clip arts, the vocabulary is smaller than it would be for real images. However, humans still use compositional language to describe clip art configurations and attributes, and make references to previous discourse elements in their messages. **Rounds.** Figure 2b shows the distribution of the numbers of conversational rounds for dialog sessions. Most interactions are shorter than 20 rounds; the median number of rounds is 7. **Durations.** In Figure 2c we see that the median session duration is 6 minutes. We had placed a 20-minute maximum limit on each session.

---

[3]Excluding 27 empty scenes from the original dataset.

### 3.3 Scene similarity metric

The goal-driven nature of the CoDraw task naturally lends itself to evaluation by comparing the reconstructed scene to the original. For this purpose we define a scene similarity metric, which allows us to automatically evaluate communication effectiveness both at the end of a dialog and at intermediate states. We use the metric to compare how well different machine-machine, human-machine, and human-human pairs can complete the task.

We represent a scene $C$ as a set of clip art objects $c \in C$, each of which consists of an identifier $id(c)$ that denotes its type, and additional features such as size and $x, y$ position. We denote by $ids(C)$ the set of clip art types that occur in the scene. Given two scenes, the intersection-over-union measure computed over clip art types is:

$$IOU(C, \hat{C}) = \frac{n_{\text{int}}}{n_{\text{union}}} = \frac{\sum_i \mathbf{1}_{i \in ids(C) \wedge i \in ids(\hat{C})}}{\left| ids(C) \cup ids(\hat{C}) \right|}$$
(1)

where $n_{\text{int}}$ ($n_{\text{union}}$) is the numbers of clip art types in the intersection (union).

To also incorporate features such as size and position, we replace the indicator function in the numerator with a term $g(i, C, \hat{C})$ that measures attribute similarity for shared clip art types. We also introduce a pairwise similarity term $h(i, j, C, \hat{C})$. Overall, scene similarity is defined as:

$$s(C, \hat{C}) = \underbrace{\frac{\sum_i g(i, C, \hat{C})}{n_{\text{union}}}}_{\text{unary}} + \underbrace{\frac{\sum_{i<j} h(i, j, C, \hat{C})}{n_{\text{union}}(n_{\text{int}} - 1)}}_{\text{pairwise}}$$
(2)

The denominator terms normalize the metric to penalize missing or extra clip art, and we set $g$ and $h$ such that our metric is on a 0-5 scale. The exact terms $g$ and $h$ are described in Appendix B.

## 4   Models

We model both the Teller and the Drawer, and evaluate the agents using the metric described in the previous section. Informed by our analysis of the collected dataset (see Section 3.2), we make several modeling assumptions compared to the full generality of the setup that humans were presented with during data collection. These assumptions hold for all models studied in this paper.

**Assumption 1: Silent Drawer.** We choose to omit the Drawer's ability to ask clarification questions: our Drawer models will not generate any messages and our Teller models will not condition on the text of the Drawer replies. This is consistent with typical human replies such as *"ok"* or *"done"* (around 62% of human Drawer replies only use a single token) and the fact that the Drawer talking is not strictly required to resolve the information asymmetry inherent in the task. We note that this assumption does not reduce the number of modalities needed to solve the task: there is still language generation on the Teller side, in addition to language understanding, scene perception, and scene generation on the Drawer side. Drawer models that can detect when a clarification is required, and then generate a natural language clarification question is interesting future work.

**Assumption 2: Full clip art library.** The other assumption is that our drawer models can select from the full clip art library. Humans are only given access to a smaller set so that it can easily fit in the user interface (Zitnick and Parikh, 2013), while ensuring that all pieces needed to reconstruct the target scene are available. We choose to adopt the full-library condition as the standard for models because it is a stricter evaluation of whether the models are able to make correct grounding decisions.

### 4.1   Rule-based nearest-neighbor methods

Simple methods can be quite effective even for what appear to be challenging tasks, so we begin by building models based on nearest-neighbors and rule-based approaches. We split the recorded human conversations available for training into a set of conversation rounds $R$ (possibly from different dialogs), where at each round $r \in R$:

- Teller sends a message $m_r$
- Drawer removes clip art pieces $C_r^{(-)}$
- Drawer adds clip art pieces $C_r^{(+)}$
- Drawer replies or ends the conversation

**Rule-based nearest-neighbor Teller.** Our first Teller model uses a rule-based dialog policy where the Teller describes exactly one clip art each time it talks. The rule-based system determines which clip art to describe during each round of conversation, following a fixed order that roughly starts with objects in the sky (sun, clouds), followed by objects in the scene (trees, Mike, Jenny), ending with small objects (sunglasses, baseball bat). The message for each object $c$ is then copied from a nearest neighbor in the data:

$$R^{(\text{single})} = \left\{ r \in R : C_r^{(-)} = \emptyset, \left| C_r^{(+)} \right| = 1 \right\} \tag{3}$$

$$\hat{r}(c) = \underset{r \in R^{(\text{single})}}{\arg\max}\, s\left( \{c\}, C_r^{(+)} \right) \tag{4}$$

$$\hat{m}(c) = m_{\hat{r}(c)} \tag{5}$$

where $s$ is the scene similarity metric from Section 3.3. This baseline approach is based on the assumptions that the Drawer's action was elicited by the Teller utterance immediately prior, and that the Teller's utterance will have a similar meaning when copied verbatim into a new conversation and scene context.

**Rule-based nearest-neighbor Drawer.** This Drawer model is the complement to the rule-based nearest-neighbor Teller. It likewise follows a fixed rule that the response to each Teller utterance should be the addition of a single clip art, and uses a character-level string edit distance $d$ to select which clip art object to add to the canvas:

$$\hat{r}'(m) = \underset{r \in R^{(\text{single})}}{\arg\min}\, d\left( m, m_r \right) \tag{6}$$

$$\hat{C}(m) = C_{\hat{r}'(m)}^{(+)} \tag{7}$$

### 4.2   Neural Drawer

Our second Drawer model is based on the neural network architecture shown in the left portion of Figure 3. At each round of conversation, the Drawer conditions on the Teller's last message, which is encoded into a vector using a bidirectional LSTM. The Drawer also uses as input a vector that represents the current state of the canvas. These vectors are then processed by a dense feed-forward neural network to produce a vector that represents the Drawer's action, which consists of adding a (possibly empty) set of clip art pieces to the drawing. It is trained using a combination of cross-entropy losses (for categorical decisions
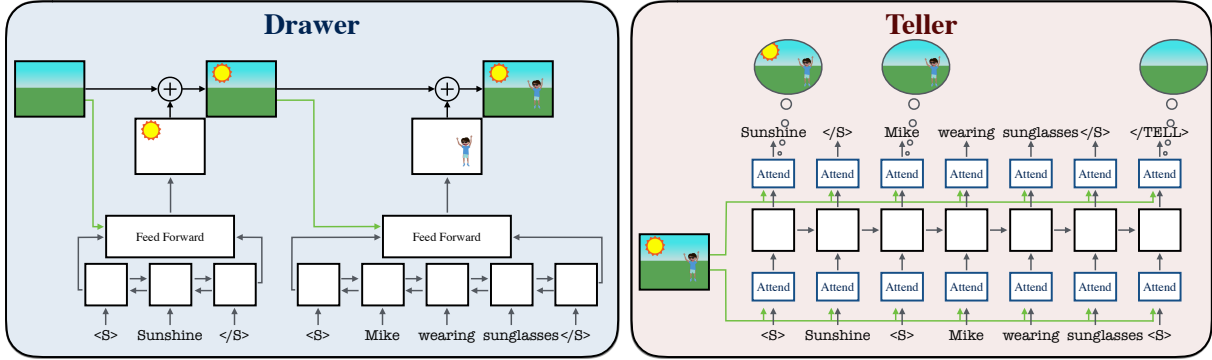
Figure 3: A sketch of our model architectures for the neural Drawer and Teller. The Drawer (left) conditions on the current state of the canvas and a BiLSTM encoding of the previous utterance to decide which clip art pieces to add to a scene. The Teller (right) uses an LSTM language model with attention to the scene (in blue) taking place before and after the LSTM. The "thought bubbles" represent intermediate supervision using an auxiliary task of predicting which clip art have not been described yet. In reinforcement learning, the intermediate scenes produced by the drawer are used to calculate rewards. Note that the language used here was constructed for illustrative purposes, and that the messages in our dataset are more detailed and precise.

such as which clip art pieces to add and what orientation to use) and $L_2$ losses that penalizes placing pieces at distant $(x, y)$ coordinates; see Appendix C for details.

## 4.3 Neural Teller: scene2seq

For our neural Teller models, we adopt an architecture that we call *scene2seq* (right portion of Figure 3). This architecture is a conditional language model over the Teller's side of the conversation with special *next-utterance tokens* to indicate when the Teller ends its current utterance and waits for a reply from the Drawer.[4] The language model is implemented using an LSTM, where information about the ground-truth scene is incorporated at both the input and output of each LSTM cell through the use of an attention mechanism. Attention occurs over individual clip art pieces: each clip art in the ground-truth scene is represented using a vector that is the sum of learned embeddings for different clip art attributes (e.g. $e_{\text{type=Mike}}$, $e_{\text{size=small}}$, etc.) At test time, the Teller's messages are constructed by decoding from the language model using greedy word selection.

To communicate effectively, the Teller must keep track of which parts of the scene it has and has not described, and also generate language that is likely to accomplish the task objective when interpreted by the Drawer. We found that training the *scene2seq* model using a maximum likeli-

hood objective did not result in long-term coherent dialogs for novel scenes. Rather than introducing a new architecture to address these deficiencies, we explore reducing them by using alternative training objectives. To better ensure that the model keeps track of which pieces of information it has already communicated, we take advantage of the availability of drawings at each round of the recorded human dialogs and introduce an auxiliary loss based on predicting these drawings. To select language that is more likely to lead to successful task completion, we further fine-tune our Teller models to directly optimize the end-task goal using reinforcement learning.

### 4.3.1 Intermediate supervision

We incorporate state tracking into the scene2seq architecture through the use of an auxiliary loss. This formulation maintains the end-to-end training procedure and keeps test-time decoding exactly the same; the only change is that during training, at each utterance separator token, the output from the LSTM is used to classify whether each clip art in the ground truth has been drawn already or not. Here we make use of the fact that the CoDraw dataset records human drawer actions at each round of the conversation, not just at the end. The network outputs a score for each clip art type, which is connected to a softmax loss for the clip art in the ground truth scene (the scores for absent clip arts do not contribute to the auxiliary loss). We find that adding such a supervisory signal reduces the Teller's propensity for repeating itself or omitting objects.

---

[4]Though none of the models in this paper handle language in the Drawer replies, these can be incorporated into the *scene2seq* framework similar to the approach of Lewis et al. (2017).

### 4.3.2 Reinforcement learning

The auxiliary loss helps the agent be more coherent throughout the dialog, but it is still an indirect proxy for the end goal of having the Drawer successfully reconstruct the scene. By training the agents using reinforcement learning (RL), it is possible to more directly optimize for the goal of the task. In this work we only train the Teller with RL, because the Teller has challenges maintaining a long-term strategy throughout a long dialog, whereas preliminary results showed that making local decisions is less detrimental for Drawers. The scene2seq Teller architecture remains unchanged, and each action from the agent is to output a word or one of two special tokens: a next-utterance token and a stop token. After each next-utterance token, our neural Drawer model is used to take an action in the scene and the resulting change in scene similarity metric is used as a reward. However, this reward scheme alone has an issue: once all objects in the scene are described, any further messages will not result in a change in the scene and have a reward of zero. As a result, there is no incentive to end the conversation. We address this by applying a penalty of 0.3 to the reward whenever the Drawer makes no changes to the scene. We train our Teller with RE-INFORCE (Williams, 1992), while the parameters of the Drawer are held fixed.

## 5 Training protocol and evaluation

To evaluate our models, we pair our models with other models, as well as with a human.

**Human-machine pairs.** We modified the interface used for data collection to have each trained model to play one game with a human per scene in the test set. We then compare the scene reconstruction quality between human-model pairs for various models and with human-human pairs.

**Script-based Drawer evaluation.** In addition to human evaluation, we would like to have automated evaluation protocols that can quickly estimate the quality of different models. Drawer models can be evaluated against a recorded human conversation from a script (a recorded dialog from the dataset) by measuring scene similarity at the end of the dialog. While this setup does not capture the full interactive nature of the task, the Drawer model still receives human descriptions of the scene and should be able to reconstruct it. Our modeling assumptions include not
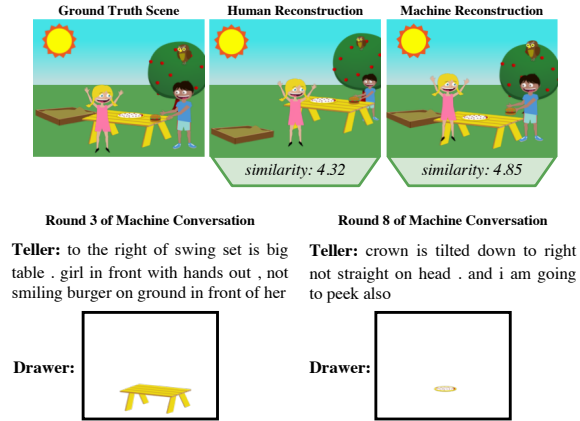


Figure 4: A rule-based nearest-neighbor Teller and Drawer pair "trained" on the same data outperforms humans for this scene according to the similarity metric, but the language used by the models doesn't always correspond in meaning to the actions taken. The top row shows a scene from the test set and corresponding human/model reconstructions. The bottom row shows the Teller message and Drawer action from two rounds of conversation by the machine agents.

giving Drawer models the ability to ask clarifying questions, which further suggests that script-based evaluation can reasonably measure model quality.

**Machine-machine evaluation.** To evaluate Teller models in a goal-driven manner, a "script" from the dataset is not sufficient. We instead consider an evaluation where a Teller model and Drawer model are paired, and their joint performance is evaluated using the scene similarity metric.

### 5.1 Crosstalk training protocol

Automatically evaluating agents, especially in the machine-machine paired setting, requires some care because a pair of agents can achieve a perfect score while communicating in a shared code that bears no resemblance to natural language. There are several ways such co-adaptation can develop. One is by overfitting to the training data to the extent that it's used as a codebook – we see this with the rule-based nearest-neighbor agents described in Section 4.1, where a Drawer-Teller pair "trained" on the same data outperforms humans on the CoDraw task. An examination of the language, however, reveals that only limited generalization has taken place (see Figure 4). Another way that agents can co-adapt is if they are trained jointly, for example using reinforcement learning. To limit these sources of co-adaptation, we propose a training protocol we call "crosstalk." In this setting, the training data is split in half, and the

|  | Teller | Drawer | Scene similarity |
|---|---|---|---|
| *Script-based* | Script (replays human messages) | Rule-Based Nearest Neighbor | 0.94 |
|  | Script (replays human messages) | Neural Network | 3.39 |
|  | Script (replays human messages) | Human | **3.83** |
| *Human-Machine* | Rule-based Nearest Neighbor | Human | 3.21 |
|  | Scene2seq (imitation learning) | Human | 2.69 |
|  | + auxiliary loss | Human | 3.04 |
|  | + RL fine-tuning | Human | **3.65** |
| *Machine-Machine* | Rule-based Nearest Neighbor | Neural Network | 3.08 |
|  | Scene2seq (imitation learning) | Neural Network | 2.67 |
|  | + auxiliary loss | Neural Network | 3.02 |
|  | + RL fine-tuning | Neural Network | **3.67** |
|  | Human | Human | **4.17** |

Table 1: Results for our models on the test set, using three types of evaluation: script-based (i.e. replaying Teller utterances from the dataset), human-machine, and machine-machine pair evaluation.

Teller and Drawer are trained separately on disjoint halves of the training data. When joint training of a Teller-Drawer pair is required (as with reinforcement learning), the training process is run separately for both halves of the training data, but evaluation pairs a Teller trained on the first partition with a Drawer trained on the second. This ensures that models can succeed only by communicating in a way that generalizes to new conversation partners, and not via a highly specialized codebook specific to model instances.

Taking the crosstalk training protocol into account, the dataset split we use for all experiments is: 40% Teller training data (3,994 scenes/dialogs), 40% Drawer training data (3,995), 10% development data (1,002) and 10% testing data (1,002).

## 6 Results

Results for our models are shown in Table 1. All numbers are scene similarities, averaged across scenes in the test set.

**Neural Drawer is the best Drawer model.** In the script setting, our neural Drawer is able to outperform the rule-based nearest-neighbor baseline (3.39 vs. 0.94) and close most of the gap between baseline (0.94) and human performance (4.17).

**Validity of script-based Drawer evaluation.** To test the validity of script-based Drawer evaluation – where a Drawer is paired with a Teller that recites the human script from the dataset corresponding to the test scenes – we include results from interactively pairing human Drawers with a Teller that recites the scripted messages. While average scene similarity is lower than when using live human Tellers (3.83 vs. 4.17), the scripts are sufficient to achieve over 91% of the effectiveness of the same Teller utterances when they were communicated live (according to our metric). The drop in similarity may be in part because the Teller can't answer clarifying questions specific to the Drawer's personal understanding of the instructions. Note that a human Drawer with a script-based Teller still outperforms our best Drawer model paired with a script-based Teller.

**Benefits of intermediate supervision and goal-driven training.** Pairing our models with humans shows that the *scene2seq* Teller model trained with imitation learning is worse than the rule-based nearest-neighbor baseline (2.69 vs. 3.21), but that the addition of an auxiliary loss followed by fine-tuning with reinforcement learning allow it to outperform the baseline (3.65 vs. 3.21). However, there is still a gap compared to human Tellers (3.65 vs. 4.17). Many participants in our human study noted that they received unclear instructions from the models they were paired with, or expressed frustration that their partners could not answer clarifying questions as a way of resolving such situations. Recall that our Teller models currently ignore any utterances from the Drawer.

**Correlation between fully-automated and human-machine evaluation.** We also report the result of paired evaluation for different Teller

models and our best Drawer, showing that the relative rankings of the different Teller types match those we see when models are paired with humans. This shows that automated evaluation while following the crosstalk training protocol is a suitable automated proxy for human-evaluation.

## 6.1 Typical errors

The errors made by Teller reflect two key challenges posed by the CoDraw task: reasoning about the context of the conversation and what has already been drawn so far, and planning ahead to fully and effectively communicate the required information. A common mistake the rule-based nearest-neighbor Teller makes is to reference objects that are not present in the current scene. Figure 4 shows an example (bottom left) where the Teller has copied a message referencing a "swing" that does not exist in the current scene. In a sample of 5 scenes from the test set, the rule-based nearest-neighbor Teller describes a non-existent object 11 times, compared to just 1 time for the scene2seq Teller trained with imitation learning. The scene2seq Teller, on the other hand, frequently describes clip art pieces multiple times or forgets to mention some of them: in the same sample of scenes, it re-describes an object 10 times (vs. 2 for the baseline) and fails to mention 11 objects (vs. 2.) The addition of an auxiliary loss and RL fine-tuning reduces these classes of errors while avoiding frequent descriptions of irrelevant objects (0 references to non-existent objects, 3 instances of re-describing an object, and 4 objects omitted.)

On the Drawer side, the most salient class of mistakes made by the neural network model is semantically inconsistent placement of multiple clip art pieces. Several instances of this can be seen in Figure 9 in Appendix D, where the Drawer places a hat in the air instead of on a person's head, or where the drawn clip art pieces overlap in a visually unnatural way.

Qualitative examples of both human and model behavior are provided in Appendix D.

## 7 Conclusion

In this paper, we introduce CoDraw: a collaborative task designed to facilitate learning of effective natural language communication in a grounded context. The task combines language, perception, and actions while permitting automated goal-driven evaluation both at the end and as a measure of intermediate progress. We introduce a dataset and models for this task, and propose a *crosstalk* training + evaluation protocol that is more generally applicable to studying emergent communication. The models we present in this paper show levels of task performance that are still far from what humans can achieve. Long-term planning and contextual reasoning as two key challenges for this task that our models only begin to address. We hope that the grounded, goal-driven communication setting that CoDraw is a testbed for can lead to future progress in building agents that can speak more naturally and better maintain coherency over a long dialog, while being grounded in perception and actions.

## Acknowledgments

## References

Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton van den Hengel. 2018. Vision-and-language navigation: Interpreting visually-grounded navigation instructions in real environments. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, volume 2.

Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C. Lawrence Zitnick, and Devi Parikh. 2015. VQA: Visual Question Answering. In *IEEE International Conference on Computer Vision*.

Lawrence W Barsalou. 1999. Perceptions of perceptual symbols. *Behavioral and brain sciences*, 22(4):637–660.

Steven Bird, Ewan Klein, and Edward Loper. 2009. *Natural language processing with Python*. "O'Reilly Media, Inc.".

Antoine Bordes, Y-Lan Boureau, and Jason Weston. 2017. Learning End-to-End Goal-Oriented Dialog. In *5th International Conference on Learning Representations*.

Xinlei Chen and C Lawrence Zitnick. 2015. Mind's eye: A recurrent visual representation for image caption generation. In *Proceedings of the IEEE con-*

*ference on computer vision and pattern recognition*, pages 2422–2431.

Andrea F. Daniele, Mohit Bansal, and Matthew R. Walter. 2016. Navigational instruction generation as inverse reinforcement learning with neural machine translation. *CoRR*, abs/1610.03164.

Abhishek Das, Samyak Datta, Georgia Gkioxari, Stefan Lee, Devi Parikh, and Dhruv Batra. 2018. Embodied question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.

Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José M. F. Moura, Devi Parikh, and Dhruv Batra. 2017a. Visual Dialog. In *IEEE Conference on Computer Vision and Pattern Recognition*.

Abhishek Das, Satwik Kottur, José M. F. Moura, Stefan Lee, and Dhruv Batra. 2017b. Learning Cooperative Visual Dialog Agents with Deep Reinforcement Learning. *arXiv preprint arXiv:1703.06585*.

Jesse Dodge, Andreea Gane, Xiang Zhang, Antoine Bordes, Sumit Chopra, Alexander Miller, Arthur Szlam, and Jason Weston. 2016. Evaluating Prerequisite Qualities for Learning End-to-End Dialog Systems. In *4th International Conference on Learning Representations*.

Jeffrey Donahue, Lisa Anne Hendricks, Sergio Guadarrama, Marcus Rohrbach, Subhashini Venugopalan, Kate Saenko, and Trevor Darrell. 2015. Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2625–2634.

Jakob Foerster, Yannis M Assael, Nando de Freitas, and Shimon Whiteson. 2016. Learning to Communicate with Deep Multi-Agent Reinforcement Learning. In *Advances in Neural Information Processing Systems 29*, pages 2137–2145.

Daniel Fried, Jacob Andreas, and Dan Klein. 2018a. Unified pragmatic models for generating and following instructions. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1951–1963. Association for Computational Linguistics.

Daniel Fried, Ronghang Hu, Volkan Cirik, Anna Rohrbach, Jacob Andreas, Louis-Philippe Morency, Taylor Berg-Kirkpatrick, Kate Saenko, Dan Klein, and Trevor Darrell. 2018b. Speaker-follower models for vision-and-language navigation. In *Proceedings of NIPS*.

Haoyuan Gao, Junhua Mao, Jie Zhou, Zhiheng Huang, Lei Wang, and Wei Xu. 2015. Are You Talking to a Machine? Dataset and Methods for Multilingual Image Question Answering. In *Advances in neural information processing systems 28*, pages 2296–2304.

Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2016. Making the V in VQA Matter: Elevating the Role of Image Understanding in Visual Question Answering. *arXiv preprint arXiv:1612.00837*.

Stevan Harnad. 1990. The symbol grounding problem. *Physica D: Nonlinear Phenomena*, 42(1-3):335–346.

He He, Anusha Balakrishnan, Mihail Eric, and Percy Liang. 2017. Learning symmetric collaborative dialogue agents with dynamic knowledge graph embeddings. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1766–1776. Association for Computational Linguistics.

Andrej Karpathy and Li Fei-Fei. 2015. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3128–3137.

Satwik Kottur, José Moura, Stefan Lee, and Dhruv Batra. 2017. Natural language does not emerge 'naturally' in multi-agent dialog. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2962–2967. Association for Computational Linguistics.

Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Fei-Fei Li. 2017. Visual Genome: Connecting Language and Vision Using Crowdsourced Dense Image Annotations. *International Journal of Computer Vision*, 123(1):32–73.

Oliver Lemon, Kallirroi Georgila, James Henderson, and Matthew Stuttle. 2006. An isu dialogue system exhibiting reinforcement learning of dialogue policies: generic slot-filling in the talk in-car system. In *Proceedings of the Eleventh Conference of the European Chapter of the Association for Computational Linguistics: Posters & Demonstrations*, pages 119–122. Association for Computational Linguistics.

David Lewis. 1969. Convention: A Philosophical Study. *Harvard University Press*.

David Lewis. 1975. Languages and language. In Keith Gunderson, editor, *Minnesota Studies in the Philosophy of Science*, pages 3–35. University of Minnesota Press.

Mike Lewis, Denis Yarats, Yann Dauphin, Devi Parikh, and Dhruv Batra. 2017. Deal or no deal? end-to-end learning of negotiation dialogues. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing*, pages 2443–2453. Association for Computational Linguistics.

Jiwei Li, Alexander H. Miller, Sumit Chopra, Marc'Aurelio Ranzato, and Jason Weston. 2017a. Dialogue Learning with Human-in-the-Loop. In *5th*

*International Conference on Learning Representations*.

Jiwei Li, Alexander H. Miller, Sumit Chopra, Marc'Aurelio Ranzato, and Jason Weston. 2017b. Learning through Dialogue Interactions by Asking Questions. In *5th International Conference on Learning Representations*.

Jiwei Li, Will Monroe, Alan Ritter, and Dan Jurafsky. 2016. Deep Reinforcement Learning for Dialogue Generation. In *2016 Conference on Empirical Methods in Natural Language Processing*, pages 1192–1202.

Reginald Long, Panupong Pasupat, and Percy Liang. 2016. Simpler context-dependent logical forms via model projections. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1456–1465. Association for Computational Linguistics.

Jiasen Lu, Caiming Xiong, Devi Parikh, and Richard Socher. 2017. Knowing When to Look: Adaptive Attention via A Visual Sentinel for Image Captioning. In *IEEE Conference on Computer Vision and Pattern Recognition*.

Mateusz Malinowski and Mario Fritz. 2014. A Multi-World Approach to Question Answering about Real-World Scenes based on Uncertain Input. In *Advances in Neural Information Processing Systems 27*, pages 1682–1690.

Hongyuan Mei, Mohit Bansal, and Matthew R. Walter. 2015. Listen, attend, and walk: Neural mapping of navigational instructions to action sequences. *CoRR*, abs/1506.04089.

Nasrin Mostafazadeh, Chris Brockett, Bill Dolan, Michel Galley, Jianfeng Gao, Georgios P. Spithourakis, and Lucy Vanderwende. 2017. Image-Grounded Conversations: Multimodal Context for Natural Question and Response Generation. *arXiv preprint arXiv:1701.08251*.

Mengye Ren, Ryan Kiros, and Richard Zemel. 2015. Exploring Models and Data for Image Question Answering. In *Advances in Neural Information Processing Systems 28*, pages 2935–2943.

Lifeng Shang, Zhengdong Lu, and Hang Li. 2015. Neural Responding Machine for Short-Text Conversation. In *53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing*, pages 1577–1586.

Alessandro Sordoni, Michel Galley, Michael Auli, Chris Brockett, Yangfeng Ji, Margaret Mitchell, Jian-Yun Nie, Jianfeng Gao, and William B. Dolan. 2015. A Neural Network Approach to Context-Sensitive Generation of Conversational Responses. In *2015 Annual Conference of the North American Chapter of the ACL*, pages 196–205.

Florian Strub, Harm de Vries, Jeremie Mary, Bilal Piot, Aaron Courville, and Olivier Pietquin. 2017. End-to-end optimization of goal-driven and visually grounded dialogue systems. *arXiv preprint arXiv:1703.05423*.

Sainbayar Sukhbaatar, Arthur Szlam, and Rob Fergus. 2016. Learning Multiagent Communication with Backpropagation. In *Advances in Neural Information Processing SystemsNeural Information Processing Systems 29*, pages 2244–2252.

Makarand Tapaswi, Yukun Zhu, Rainer Stiefelhagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler. 2016. MovieQA: Understanding Stories in Movies through Question-Answering. In *IEEE Conference on Computer Vision and Pattern Recognition*.

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

Orioi Vinyals and Quoc V. Le. 2015. A Neural Conversational Model. In *ICML Deep Learning Workshop 2015*.

Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2017. Show and Tell: Lessons learned from the 2015 MSCOCO Image Captioning Challenge. *IEEE transactions on pattern analysis and machine intelligence*, 39(4):652–663.

Adam Vogel and Daniel Jurafsky. 2010. Learning to follow navigational directions. In *Proceedings of the 48th Annual Meeting of the Association for Computational Linguistics*, pages 806–814. Association for Computational Linguistics.

Harm de Vries, Florian Strub, Sarath Chandar, Olivier Pietquin, Hugo Larochelle, and Aaron Courville. 2016. GuessWhat?! Visual object discovery through multi-modal dialogue. *arXiv preprint arXiv:1611.08481*.

Matthew R. Walter, Sachithra Hemachandra, Bianca Homberg, Stefanie Tellex, and Seth Teller. 2014. A framework for learning semantic maps from grounded natural language descriptions. *The International Journal of Robotics Research*, 33(9):1167–1190.

Sida I. Wang, Samuel Ginn, Percy Liang, and Christoper D. Manning. 2017. Naturalizing a Programming Language via Interactive Learning. In *55th Annual Meeting of the Association for Computational Linguistics*.

Sida I. Wang, Percy Liang, and Christopher D. Manning. 2016. Learning Language Games through Interaction. In *54th Annual Meeting of the Association for Computational Linguistics*, pages 2368–2378.

Zhuoran Wang and Oliver Lemon. 2013. A simple and generic belief tracking mechanism for the dialog state tracking challenge: On the believability of observed information. In *SIGDIAL 2013 Conference*, pages 423–432.

Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256.

Kelvin Xu, Aaron Courville, Richard S Zemel, and Yoshua Bengio. 2015. Show, Attend and Tell : Neural Image Caption Generation with Visual Attention. In *32nd International Conference on Machine Learning*.

Licheng Yu, Eunbyung Park, Alexander C Berg, and Tamara L. Berg. 2015. Visual Madlibs : Fill in the blank Description Generation and Question Answering. In *IEEE International Conference on Computer Vision*, pages 2461–2469.

Peng Zhang, Yash Goyal, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. 2016. Yin and Yang: Balancing and Answering Binary Visual Questions. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5014–5022.

Yuke Zhu, Oliver Groth, Michael Bernstein, and Li Fei-Fei. 2016. Visual7W: Grounded Question Answering in Images. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4995–5004.

C. Lawrence Zitnick and Devi Parikh. 2013. Bringing semantics into focus using visual abstraction. *Proceedings of the IEEE Computer Society Conference on Computer Vision and Pattern Recognition*, pages 3009–3016.

C. Lawrence Zitnick, Devi Parikh, and Lucy Vanderwende. 2013. Learning the visual interpretation of sentences. *Proceedings of the IEEE International Conference on Computer Vision*, pages 1681–1688.

# A  Interface and data collection

## A.1  Interface

Figure 5 shows the interface for the Teller, and Figure 6 shows the interface for the Drawer. Following previous works (Zitnick et al., 2013; Zitnick and Parikh, 2013), Drawers are given 20 clip art objects selected randomly from the 58 clip art objects in the library, while ensuring that all objects required to reconstruct the scene are available.
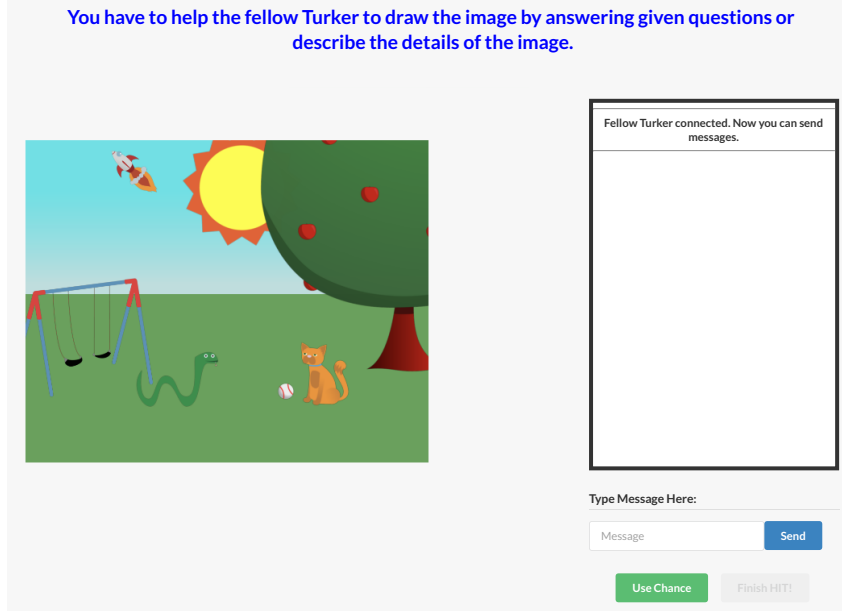


Figure 5: User interface for Teller. The left image is an abstract scene from Zitnick and Parikh (2013). The Teller sends messages using an input box. The Teller has a single chance to peek at the Drawer's canvas to correct mistakes. The Teller can decide when to finish the session.

## A.2  Additional interaction: a chance to peek

To make sure that the natural language focused on the high-level semantics of the scene rather than instructions calling for the execution of low-level clip art manipulation actions, we did not allow Teller to continuously observe Drawer's canvas. However, direct visual feedback may be necessary to get the all the details right. We also hypothesize that such feedback would help human participants calibrate themselves when they are new to the task (models do not have this issue because of the rich supervisory signal available in the collected dataset.)

To capture this idea, we give one chance for the Teller to look at the Drawer's canvas using a 'peek' button in the interface. Communication is only allowed after the peek window is closed.

Although we consider the ability to peek to be a part of the CoDraw task, we leave for future work the creation of models that can strategically reason about when to use this chance in a way that maximizes task effectiveness. We note that omitting this behavior from the Teller models described in this paper does not decrease the number of modalities needed to complete the task – our models still incorporate language understanding, language generation, perception, and action.

## A.3  Participant statistics

We found that approximately 13.6% of human participants disconnected early, prior to fully completing the task with their partner. We paid participants who stayed in the conversation and had posted at least three messages. However, we exclude those incomplete sessions in the dataset, and only use the completed sessions.

There are 616 unique participants represented in our collected data. Among these workers, the 5 most active have done 26.63% of all finished tasks (1,419, 1,358, 1,112, 1,110, and 1,068 tasks). Across all
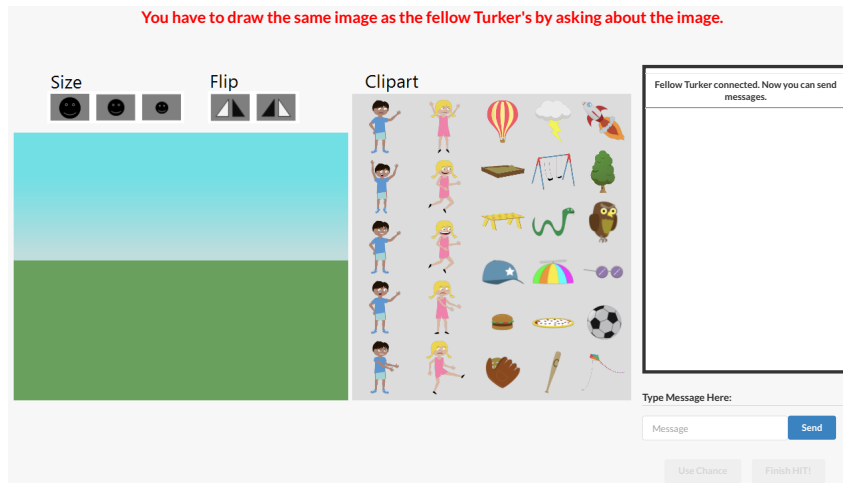
Figure 6: User interface for a Drawer. The Drawer has an empty canvas and a randomly generated drawing palette of Mike, Jenny, and 18 other objects, chosen from a library of 58 clip arts. We ensure that using the available objects, the Drawer can fully reproduce the scene. Using the library, the Drawer can draw on the canvas in a drag-and-drop fashion. The Drawer can also send messages using the provided input box. However, the peek button is disabled: only the Teller can use it.

workers, the maximum, median, and minimum numbers of tasks finished by a worker are 1,419, 3, and 1, respectively.

## A.4   Pre-processing

We pre-process all collected Teller and Drawer utterances using the Bing Spell Check API[5]. The text is then tokenized using the Python Natural Language Toolkit, nltk (Bird et al., 2009). We release the token sequences after pre-processing as part of the CoDraw dataset, so that different models may be compared in a standard set of data conditions. At the same time, raw (unprocessed) text is also made available, to allow revisiting the pre-processing decisions should the need arise.

---

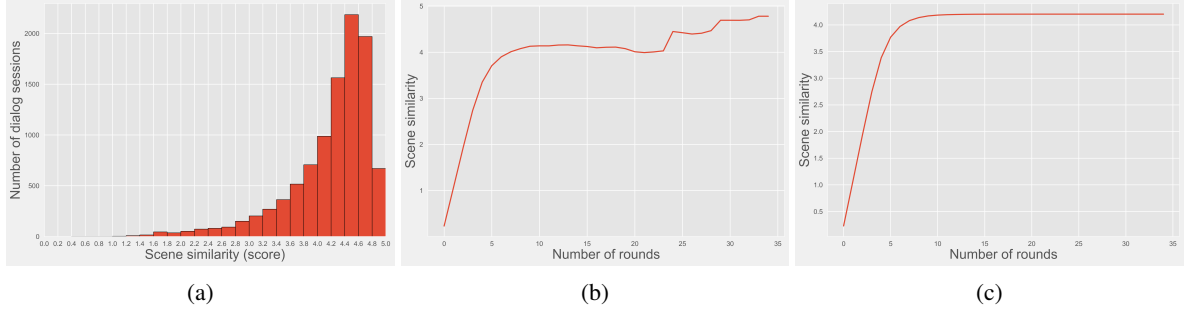[5]https://www.microsoft.com/cognitive-services/en-us/bing-spell-check-api

Figure 7: **(a)** The distribution of overall scores at the end of the dialog. **(b-c)** Average scene similarity plotted for different conversation rounds. In (b), only conversations that have reached the given number of rounds are included. In (c), conversations that end early are padded to 35 rounds through the addition of empty messages/actions.

## B  Scene similarity metric

The clip art library consists of 58 base clip art types (e.g. the sun, a cloud, Mike, Jenny, soccer ball, etc.) Each clip art object $c$ consists of an identifier $id(c)$ that denotes its type, an indicator feature vector $f(c)$ that determines properties such such as size and orientation (e.g. $\mathbf{1}_{\text{size=small}}$, $\mathbf{1}_{\text{size=medium}}$, etc. for a total of 41 binary features), and two real-valued features $x(c)$ and $y(c)$ that encode the $x$ and $y$ position on the canvas, normalized to the 0-1 range.

We represent a scene $C$ as a set of individual clip art objects $c \in C$. We denote by $ids(C)$ the set of clip art types that occurs in the scene. Following Zitnick et al. (2013), a given clip art type may occur at most once in the scene; let $C[i]$ be the clip art $c \in C$ such that $id(c) = i$.

Given a ground-truth scene $C$ and a predicted scene $\hat{C}$ scene similarity $s$ is defined as:

$$s(C, \hat{C}) = \underbrace{\frac{\sum_{i \in ids(C) \cap ids(\hat{C})} g(C[i], \hat{C}[i])}{|ids(C) \cup ids(\hat{C})|}}_{\text{unary}} + \underbrace{\frac{\sum_{i,j \in ids(C) \cap ids(\hat{C}), i<j} h(C[i], C[j], \hat{C}[i], \hat{C}[j])}{|ids(C) \cup ids(\hat{C})|(|ids(C) \cap ids(\hat{C})| - 1)}}_{\text{pairwise}}$$

where

$$g(c, \hat{c}) = w_0$$
$$- w_1 \mathbf{1}_{\text{clip art piece } \hat{c} \text{ faces the wrong direction}}$$
$$- w_2 \mathbf{1}_{\text{clip art piece } \hat{c} \text{ is Mike or Jenny and has the wrong facial expression}}$$
$$- w_3 \mathbf{1}_{\text{clip art piece } \hat{c} \text{ is Mike or Jenny and has the wrong body pose}}$$
$$- w_4 \mathbf{1}_{\text{clip art piece } \hat{c} \text{ has the wrong size}}$$
$$- w_5 \sqrt{(x(\hat{c}) - x(c))^2 + (y(\hat{c}) - y(c))^2}$$

$$h(c_i, c_j, \hat{c}_i, \hat{c}_j) = - w_6 \mathbf{1}_{(\hat{x}_{c_i} - \hat{x}_{c_j})(x_{c_i} - x_{c_j}) < 0}$$
$$- w_7 \mathbf{1}_{(\hat{y}_{c_i} - \hat{y}_{c_j})(y_{c_i} - y_{c_j}) < 0}$$

We use parameters $\boldsymbol{w} = [5, 1, 0.5, 0.5, 1, 1, 1, 1]$, which provides a balance between the different components and ensures that scene similarities are constrained to be between 0 and 5.

Figure 7a shows the distribution of scene similarity scores throughout the dataset. Figure 7b-c shows the progress of scene similarity scores over the rounds of a conversation. An average conversation is done improving the scene similarity after about 5 rounds, but for longer conversations that continue to 23 rounds, there is still room for improvement.

## C Neural Drawer architecture

In this section, we describe in greater detail our neural network architecture approach for the Drawer. Contextual reasoning is an important part of the CoDraw task: each message from the Teller can relate back to what the Drawer has previously heard or drawn, and the clip art pieces it places on the canvas must form a semantically coherent scene. To capture these effects, our model should condition on the past history of the conversation and use an action representation that is conducive to generating coherent scenes.

When considering past history, we make the Markovian assumption that the current state of the Drawer's canvas captures all information from the previous rounds of dialog. Thus, the Drawer need only consider the most recent utterance from the Teller and the current canvas to decide what to draw next. We experimented with incorporating additional context – such as previous messages from the Teller or the action sequence by which the Drawer arrived at its current canvas configuration – but did not observe any gains in performance.

We represent the state of the canvas with a vector $v_{canvas}$ that is the concatenation of feature vectors for each of the 58 possible clip art types:

$$v_{canvas}(C) = [v_0(C); v_1(C); \ldots; v_{57}(C)]$$

$$\text{where } v_i(C) = \begin{cases} [1; f(C[i]); x(C[i]); y(C[i])] & \text{if } i \in ids(C) \\ \mathbf{0} & \text{otherwise} \end{cases}$$

The individual feature vectors $v_i(C)$ represent binary and $(x, y)$ features of the clip art piece if it is present on the canvas, and are zeroed out if a clip art of the given type is not present on the canvas.

The most recent Teller utterance is encoded into a vector $v_{msg}$ using a bi-directional LSTM. A vector representing the Drawer's action is then computed using a feed-forward network with a rectified linear unit (ReLU) nonlinearity:

$$v_{action} = W_{out} relu(W_{canvas} v_{canvas} + W_{msg} v_{msg} + b_{in}) + b_{out}$$

The action representation $v_{action}$ has the form:

$$v_{action} = [a_0; a_1; \ldots; a_{57}]$$

$$\text{where } a_i = \begin{bmatrix} q(i \in ids(C)) \\ q(f_0(C[i]) = 1 | i \in ids(C)) \\ q(f_1(C[i]) = 1 | i \in ids(C)) \\ \ldots \\ \hat{x}(C[i]) \\ \hat{y}(C[i]) \end{bmatrix}$$

The values $\hat{x}(C[i])$ and $\hat{y}(C[i])$ are the predicted location for clipart $C[i]$ if it is placed on the canvas, and each quantity $q(event)$ is a logit corresponding to a particular event. The probability of adding a clip art piece to the scene is calculated using the sigmoid function:

$$p(i \in ids(C)) = \frac{1}{1 + \exp -q(i \in ids(C))}$$

while all other probabilities are calculated by applying softmax to each set of mutually-exclusive outcomes, e.g.:

$$p(size(C[i]) = \text{small} | i \in ids(C)) = \frac{\exp\left(q(size(C[i]) = \text{small} | i \in ids(C))\right)}{\sum_{s \in \{\text{small,medium,large}\}} \exp\left(q(size(C[i]) = s | i \in ids(C))\right)}$$

At inference time, the Drawer's action is chosen using greedy decoding. A clip art of type $i$ is added to the canvas if $p(i \in ids(C)) > 0.5$, in which case it is placed at location $(\hat{x}(C[i]), \hat{y}(C[i]))$ with its orientation, size, and other attributes set to their most probable values (as determined by the vector $a_i$.)

The model is trained using a combination of cross-entropy losses (that maximize the probability of the categorical decisions present in the human action) and an $L_2$ loss that compares the locations where the human placed each clip art piece with the model's estimate.

# D   Qualitative examples

Figure 8 shows some examples of scenes and dialogs from the CoDraw dataset. The behavior of our Drawer and Teller models on a few randomly-selected scenes is illustrated in Figures 9, 10, and 11.

Figure 8: Examples from the Collaborative Drawing (CoDraw) dataset, chosen at random from the test set. The images depict the Drawer's canvas after each round of conversation. From left to right, we show rounds one through four, then the last round, followed by the ground truth scene. The corresponding conversations between the Teller (T) and Drawer (D) are shown below the images. Note that there is no restriction on which of the two participants begins or ends the dialog.

Figure 9: Drawer model behavior where the Teller is simulated by replaying messages associated with three randomly-selected scenes from the test set. The images depict the Drawer's canvas after each round of conversation. For each of the three scenes, the behavior of the Rule-Based Nearest-Neighbor Drawer is shown in the upper row and the behavior of the Neural Network Drawer is shown in the lower row.

Scene A

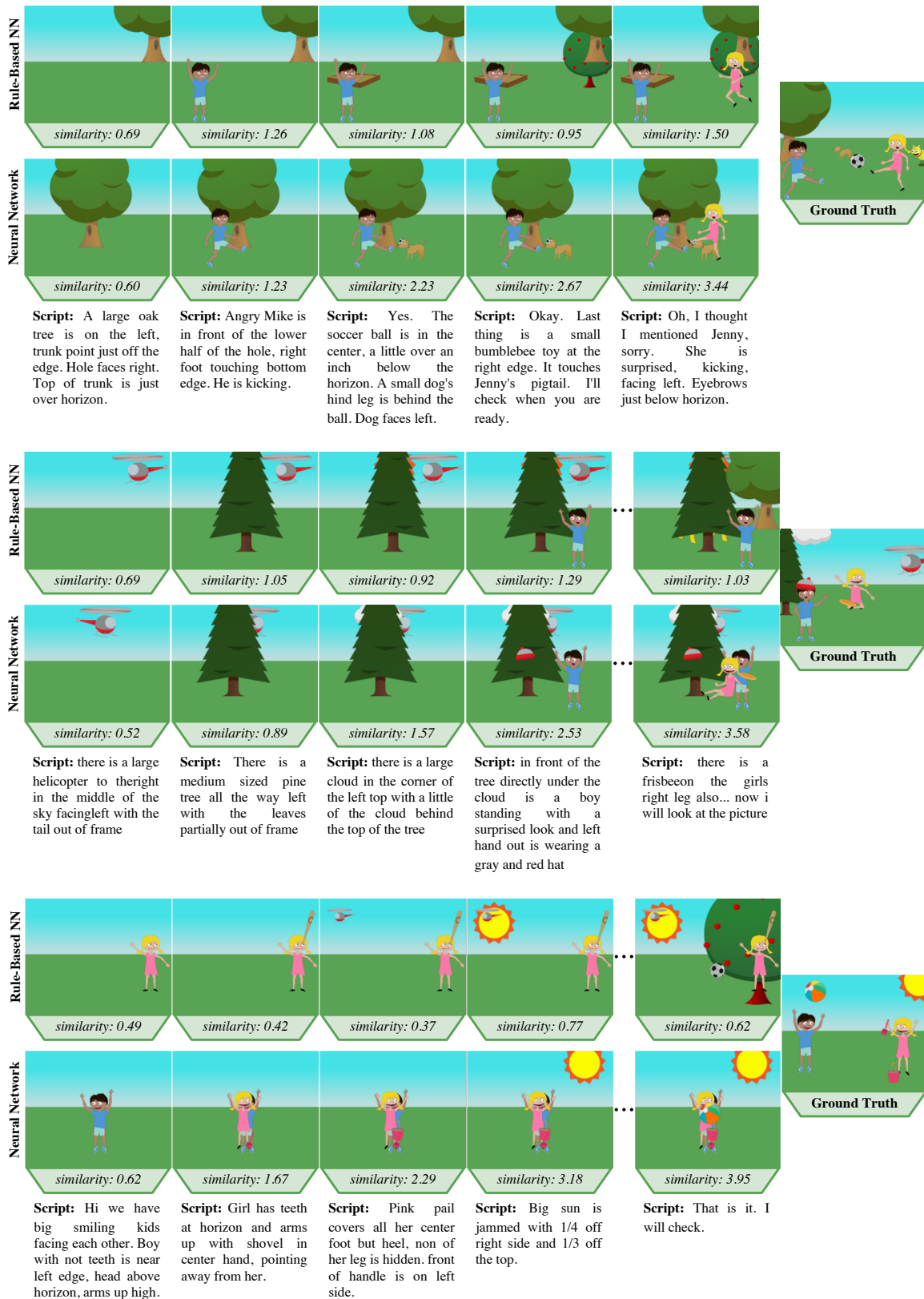| Scene, Round | Rule-Based Nearest-Neighbor | Scene2seq | | |
|---|---|---|---|---|
| | | Imitation learning | + Auxiliary Loss | + RL Fine-tuning |
| A1 | in the top left corner is a medium sun very close to the top and left edge | medium sun on left , 1 2 inch from top and side | on the right hand side is a medium bear facing left with a basketball in his left hand | medium sun top left |
| A2 | guy angry arms down facing right like his right knees touches the top right corner of the table | medium bear on right , facing left , 1 2 " from bottom , small basketball in left hand , racket in left hand | in the middle of the grass is a medium sized boy facing right with a racket in his left hand | medium boy on right angry facing right holding racket |
| A3 | right under it is medium girl , waving with a sad face , facing right . | small bear on right side , facing left , 1 2 " from bottom , small bear on left facing right , armpits on horizon , | he is holding a racket in his left hand and a tennis racket in the middle of the sky | medium girl on left facing right sad one arm up holding basketball |
| A4 | on the right side of screen , medium bear , chest on horizon , paw almost touching right side of screen | small bear on right facing left , chest at horizon , small basketball in left hand , 1 inch from left , 1 2 from bottom | on the left hand side is a medium girl with sad face facing right with a racket in her right hand . | medium bear on right facing left wearing rainbow hat basketball racket in right hand |
| A5 | he also has a rainbow hat , but it 's not snugly on him , it 's only slightly touching his head . | small bear on right facing left , chest at horizon , small basketball in left hand , racket in left hand , racket in left hand , racket in left | the bear is facing left . the bear is facing right . the bear is facing left . | <STOP> |
| A6 | sorry there 's a basketball above the boy 's head near the tree | <STOP> | <STOP> | |
| A7 | she is holding a racket in her left hand . | | | |
| A8 | <STOP> | | | |



Scene B

| Scene, Round | Rule-Based Nearest-Neighbor | Scene2seq | | |
|---|---|---|---|---|
| | | Imitation learning | + Auxiliary Loss | + RL Fine-tuning |
| B1 | two finger spaces to right of helicopter , small hot air balloon . top of balloon level with helicopter propeller | small hot air balloon , top right , 1 2 " from top , 2 " from right | small hot air balloon on right corner , fully visible . | small hot air balloon top right |
| B2 | in front of tree is boy , he is to the left part of tree and is covering the curve up . he is angry , standing , arms , out facing left | small girl , running , facing right , surprised , 1 " from bottom , 1 2 " from left | small hot balloon on right corner , half " from top . | large bear on left faced right |
| B3 | the head of surprised girl is on front the trunk . she is like running and faces right . | small pine tree behind her , bottom of trunk at horizon , bottom of trunk at horizon , small boy in front of tree , head touching bottom of tree , standing , smiling , facing right , holding a hot dog in left hand | on center , a mad mike with hands front facing left . | medium shocked girl running facing right holding hot dog in right hand |
| B4 | large bear on left side . facing right . arm cut off on left . feet half inch from bottom | he is at the edge , small size , 1 2 " from bottom , 1 2 " from bottom , small kite in left corner , facing right | on center , a small jenny , surprised , running facing right , shoulders at horizon . | boy on right angry facing left holding kite in right hand hot dog in right hand |
| B5 | place a hot dog in her left hand | <STOP> | on center , a small hot balloon , half finger from right side . | <STOP> |
| B6 | right top a kite , small size , tail is off screen . below kite dog medium size , facing left | | a small kite on right corner , only half visible . | |
| B7 | <STOP> | | <STOP> | |

Figure 10: A comparison of the descriptions generated by each of our Teller models for two randomly-sampled scenes from the test set.

**Row 1:**

*similarity: 0.77*
**T:** large bushy tree on left hole facing right
**D:** ok

*similarity: 1.55*
**T:** large cloud on right
**D:** ok

*similarity: 2.90*
**T:** large boy on left facing right sad one hand up wearing blue hat
**D:** ok

*similarity: 4.32*
**T:** large girl angry facing left holding hot dog in left hand
**D:** ok

*similarity: 4.32*
**T:** large hot dog on left facing left
**D:** ok
**T:** <stop>

**Ground Truth**

**Row 2:**

*similarity: 0.81*
**T:** large cloud on right
**D:** ok

*similarity: 2.17*
**T:** large boy on left facing right wearing pirate hat
**D:** ok

*similarity: 2.88*
**T:** soccer ball in middle
**D:** ok

*similarity: 3.60*
**T:** big dog on right facing left
**D:** ok

*similarity: 4.25*
**T:** girl sitting on right facing right smiling wearing pirate hat
**D:** ok
**T:** <stop>

**Ground Truth**

**Row 3:**

*similarity: 0.80*
**T:** large cloud on left
**D:** ok

*similarity: 1.60*
**T:** small bushy tree on right
**D:** ok

*similarity: 2.25*
**T:** large rocket in middle facing right
**D:** ok

*similarity: 3.73*
**T:** girl sitting surprised facing right wearing sunglasses
**D:** ok

*similarity: 3.73*
**T:** large rocket on left facing right
**D:** ok
**T:** <stop>

**Ground Truth**

**Row 4:**

*similarity: 0.82*
**T:** large cloud on left
**D:** ok

*similarity: 1.63*
**T:** large bushy tree on right hole facing right
**D:** ok

*similarity: 2.32*
**T:** large snake on left facing left
**D:** ok

*similarity: 3.50*
**T:** large boy on left facing left sad holding bat
**D:** ok

*similarity: 4.22*
**T:** girl surprised kneeling facing left holding bat in left hand
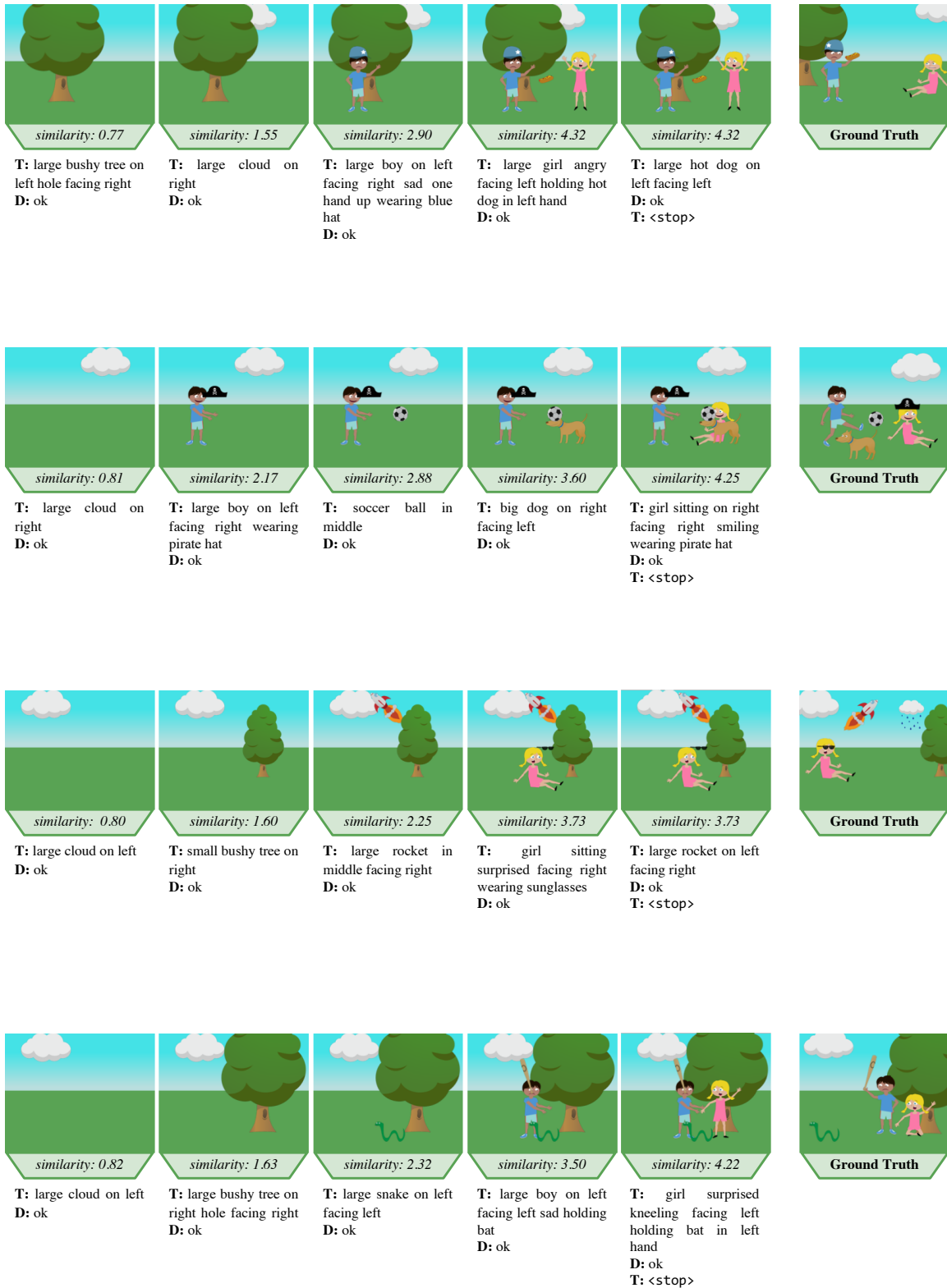**D:** ok
**T:** <stop>

**Ground Truth**

Figure 11: Dialogs from our best Teller model (scene2seq with an auxiliary loss and RL fine-tuning) communicating with our best Drawer model (Neural Network Drawer). The dialogs feature the same scenes as in Figure 8, which were sampled at random from the test set. From left to right, we show the first to the fifth rounds of conversations, followed by the ground truth scene. Our Teller model chose to use exactly five rounds for each of these four scenes. The corresponding conversations between Teller (T) and Drawer (D) are shown below the images.