# Final Report - SS4850

Emiliano Penaloza
Anthony Rinaldi

April 19, 2022

# Table of Contents

## 1   Introduction

The way in which most people/agents negotiate and come to a decision based on differing preferences and goals is through the use of natural language. This is a complex problem because one must determine their intent, find a way to communicate this to another agent, and understand the intent of the other agent. What makes this scenario even more complex is the use of deceit, which is difficult to ascertain, but also very common in human negotiations [6]. Negotiation dialogues can take on many different forms depending on the agents involved (e.g., cooperative or aggressive) and each agent must come up with *utterances* (sequences of words) to help them achieve their goal in the negotiation.

We take a similar approach to Lewis et al. [5] using a very large one-on-one negotiation dataset to train neural negotiation agents. These agents learn to negotiate by maximizing the likelihood of copying human actions and are further refined by using reinforcement learning (RL) when negotiating against one another. Reinforcement learning forces the agent to maximize their reward, rather than simply copying human actions; possibly giving the agents skills that never would be exhibited by human participants. We supplement these models with the use of the *Monte Carlo tree search* (MCTS) algorithm to estimate the expected reward of different utterances. MCTS has been widely successful in strategic games which is why we expect it to perform well in strategic negotiations [2].

We study these negotiations by using a dataset of semi-cooperative dialogues between participants. Participants are shown a set of objects (books, balls, and hats), each with a specific value and are tasked with maximizing their value by determining how to divide the objects between themselves and another participant. Participants are unaware of the value function of the other agent.

The remainder of the paper is organized as follows: §2 gives an overview of the negotiation dataset used in the paper. §3 introduces both the exploratory and main methodology of the paper. §4 describes the main results of the paper. §5 concludes the paper with discussion about the implications of the results and §6 proposes future works.

## 2   Dataset

The dataset consists of 5,808 dialogues collected using Amazon's Mechanical Turk, a crowdsourcing service that gathers freelance workers to complete on-demand tasks that computers are unable to do. Each dialogue consists of two human agents engaging in a negotiation aiming to maximize their total reward. The two agents negotiate over 5-7 total items each of which belongs to one of three categories: hats, balls and books. The negotiation begins by providing each agent with a different randomly generated value function where the total value of all items is 10 for each agent and no item has a zero reward for both users. The agents then engage in conversation until one of them declares that an agreement has been made which marks the end of the negotiation. Thereafter, each agent declares the agreed-upon number of items assigned to each agent, if both agents agree on the decision, the appropriate rewards are assigned, if they disagree, they are both assigned a total reward of zero. This methodology of gathering dialogues yielded 2,236 different scenarios (unique value functions and number of items). An example scenario and the corresponding dialogue can be seen in Table 1 and Figure 1, respectively. As previously mentioned, the data is publicly available in the following repository: https://github.com/facebookresearch/end-to-end-negotiator.

|                          | Hat | Ball | Book |
|--------------------------|-----|------|------|
| Amount of items in pool  | 1   | 2    | 3    |
| Agent 1 value function   | 4   | 0    | 2    |
| Agent 2 value function   | 1   | 3    | 1    |

Table 1: Example Scenario

Dialogue:

**Agent 1:** I want the books and the hats, you get the ball
**Agent 2:** Give me a book too and we have a deal
**Agent 1:** Ok, deal
**Agent 2:** <choose>

Figure 1: Example Dialogue

## 3   Methods

### 3.1   Exploratory Data Analysis

The data used in this report is quite unique and traditional data analysis methods will not be very useful. Most of the data are English sentences, which we cannot easily visualize. One approach we will use is plotting the frequency of words used, however, this is not used in modelling and will only be used to understand the language in the dialogues. The other visualizations we will consider involve the score of each agent and their value functions. We will plot the distribution of the final score for the player negotiating first and the player negotiating second. We will view the number of Pareto optimal outcomes.[1] We will view the distribution of number of turns and number of words per turn across all negotiations. Finally, we will consider the relationship between some of these variables, such as the score of a game versus the number of turns in a game. Not all of the exploratory data analysis will be useful in this report, thus we will only include significant findings in the final report.

### 3.2   Existing Methods

Lewis et al. [5] proposed four model architectures to build a conversational negotiation agent. They first propose a model that is trained by optimizing the log-likelihood of the predicted token plus the log-likelihood of the output choice. Thereafter, they propose a model that is, first trained as the previous log-likelihood model, but then fine-tuned by selecting the utterance with the maximum expected reward. Finally, they proposed utterance rollouts for action selection, which can be combined with either of the stated models. We propose to improve these models by using *Monte Carlo Tree Search* (MCTS) and *Transformers*. We propose eight new methods.

### 3.3   Transformers vs GRUs

In the original work, Lewis et al. [5] propose a GRU (Gated Recurrent Unit) model architecture. We propose using Transformers in preference to GRUs. Given the large amount of empirical evidence showing that Transformers outperform GRU-based architectures [4], it is an intuitive suggestion. Transformers, initially proposed by Vaswani et al. [7] are neural network architectures that aim to replace traditional convolutional and recurrent networks. They have been shown to outperform said networks in language and vision tasks by exploiting multi-head self-attention. Unlike traditional recurrent or convolutional layers, through the use of multi-head attention, transformers are able to predict the output sequence by attending to the most relevant previous information.

---

[1]A solution is Pareto optimal if neither agent's score can be improved without lowering the other's score.

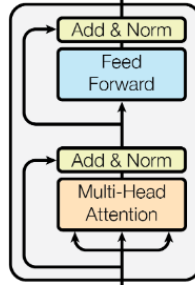| Lewis et al. [5] GRU Architecture | Proposed Transformer Architecture |
|---|---|
| $h^g = GRU_g(g)$ <br> $h_t = GRU_w(h_{t-1}, [Ex_{t-1}, h^g])$ <br> $h_t^{o\leftrightarrows} = GRU_{o\leftrightarrows}(h_{t\pm1}^{o\leftrightarrows}, [Ex_t, h_t])$ <br> $h_t^a = W[tanh(W'h_t^{o\leftrightarrows})]$ <br> $\alpha_t = \frac{exp(w \cdot h_t^a)}{\sum_t exp(w \cdot h_t^a)}$ <br> $h^s = tanh(W^s[h^g, \sum_t \alpha_t h_t])$ | $h^g = GRU_g(g)$ <br> $h_t = GRU_w(h_{t-1}, [Ex_{t-1}, h^g])$ <br> $h^o = Trans_o([Ex, h_t])$ <br> $h^a = W[tanh(W'h^o)]$ <br> $\alpha_t = \frac{exp(w \cdot h^a[t])}{exp(w \cdot h^a)}$ <br> $h^s = tanh(W^s[h^g, \sum_t \alpha_t h_t])$ |

Table 2: Model Architectures



Figure 2: Transformer Encoder Architecture

This method not only, on average, yields improved performance but allows for parallelization when training the network. This is because each *head* in the layer is summative rather than a product, allowing for faster gradient calculations. For a comparison between our proposed model architecture and the one proposed by Lewis et al. [5] please refer to Table 2. We expect that by introducing a Transformer, we will not only achieve a higher average score but produce more human-like utterances.

We will not be using a complete Transformer architecture, rather will be using just a single encoder layer. The network layout for the encoder can be seen in Figure 2. This architecture includes the powerful multi-head attention, along with a feed-forward layer, making this attention global. This means that after attention scores are calculated, a fully connected layer is used so that every point in the sequence can be connected to one another (i.e., global). We do not use the decoder part of the model because this is typically used for sentence-to-sentence modelling. Since our approach takes in a complete dialogue and generates the agent's item choices, there is no sequence-to-sequence modelling so the decoder is not necessary. The code used to create the transformer model and make predictions with it can be found in Code 11.

### 3.4  Monte Carlo Tree Search

This algorithm is similar to the rollouts algorithm that Lewis et al. [5] propose in the report. The key difference is that MCTS builds a game tree and simulates many rollouts to determine the optimal action given the current state/node. Once the actions are proposed, only those that are the most successful given the current state are expanded upon and simulated (rolled out) further. Since the algorithm uses an *Upper Confidence Bound* (UCB) to weight which node to expand, it is able to explore a much larger state space than a traditional rollout algorithm. Additionally,

MCTS explores a larger state space in a more computationally efficient matter. Since MCTS does not require expansion of every node but uses UCB to weight its options, it only explores the most promising nodes, marginally reducing the computational cost. Yet even with the reduced expansion, MCTS still achieves higher performance by searching into a deeper state space. Due to this, we expect MCTS to out-perform the traditional rollout method, and produce more relevant and fluent utterances. The MCTS description can be found in Algorithm 1.

The UCB formula is shown below in equation (1). The UCB gives us a number that represents the trade-off between exploring (trying new actions) and exploiting (reusing actions that we know have a good result). In the equation, $s_i$ represents the total score of node $i$, $n_i$ represents the number of times node $i$ has been visited, $C$ is a constant that we set to its typical value of 2, and $N_i$ represents the number of times the parent of node $i$ has been visited.

$$UCB_i = \frac{s_i}{n_i} + C \cdot \sqrt{\frac{\log N_i}{n_i}} \tag{1}$$

The class we implement for MCTS can be found in Code 12. This class has the ability to understand its opponent's turn, and thereafter perform MCTS to determine its best response. After the opponent gives their response the class will generate five unique possible responses. From there, it will iterate through 50 simulations of the MCTS algorithm, determining which of the five responses will result in the highest score. It then chooses the response that results in the highest score as its next action.

### 3.5  Proposed Models

We propose eight additional models expanding upon Lewis et al.'s [5] initial four. These models will be evaluated on the same dataset as the original four. In order to be computationally efficient, all hyperparameters will be selected using a validation dataset similar to Lewis et al. [5].

1. Supervised model with probabilistic sampling with addition of transformers

2. Supervised model, rollout dialogue (likelihood evaluation) with addition of transformers

3. RL, likelihood dialogue (one-step expected reward) with addition of transformers

4. RL, rollout dialogue (expected reward) with addition of transformers

5. Supervised model, MCTS dialogue (likelihood evaluation) with addition of transformers

6. Supervised-model, MCTS dialogue (likelihood evaluation) base GRUs

7. RL, MCTS dialogue (expected reward) base GRUs

8. RL, MCTS dialogue (expected reward) with addition of transformers

All hyper-parameters; number of layers, layer depth, discount factor, output importance $\alpha$, learning rate, and batch size were determined using cross-validation on a smaller validation dataset.

---

**Algorithm 1:** MCTS For Negotiation Dialogue

**Result:** Get the next best word based on Monte Carlo simulation

**Input**   : Dialogue History (i.e., $x_{1,\dots,k}$), items, value function, number of iterations

**Output:** Utterance to finish turn (i.e., $x_{k+1,\dots,k+n}$)

---

**1** Initialize the root node of the game tree

**2** **while** *current iteration < number of iterations* **do**

**3**  | Traverse the tree downwards until reaching a leaf node, selecting subsequent nodes that have the largest UCF score

**4**  | **if** *node has been visited before* **then**

**5**  |  | Expand the node to have children (each child is a dialogue for a full turn)

**6**  |  | Select the first child as the current node

**7**  | **else**

**8**  |  | Select the leaf node as the current node

**9**  | **end**

**10**  | Rollout to the end of the dialogue on the current node by sampling the likelihood function of tokens

**11**  | Make selections for the items based on the dialogue that occurs during the rollout

**12**  | Calculate the score of the game given the value function of the player

**13**  | Backpropogate the results:

**14**  |   Add the score of the game recursively to every parent node until reaching the root node

**15**  |   Increment the *times visited* variable of every parent recursively until reaching the root node

**16**  | Increment current iteration

**17** **end**

**18** Return the utterance with the highest average score

---

### 3.6   Training Setting

Model training will consist of supervised learning based on a training set of recorded dialogues and RL against another fitted model. This will give us two models: (1) supervised and (2) supervised + RL. With these models, we can apply rollouts and MCTS when competing against another player to see how the models perform. Supervised training takes the typical steps of deep learning: pass a batch through the network to calculate the output, calculate a loss function, backpropagate the error, calculate the new parameter estimates, and repeat. We found poor model performance using the same hyperparameters suggested by Lewis et al. [5], so implemented a Bayesian parameter tuning approach. This uses the `optuna` package in `Python` to determine the best parameters given priors and the observation likelihood. The code for the hyperparameter tuning can be found in Code 13. Thereafter we fit model (1) using the entire training dataset and the optimal hyperparameters. We then perform RL on model (1) to generate model (2). This is done by allowing model (1) to play a version of itself and work on optimizing the reward it receives from the negotiation.[2] For parameter stability, only one version of the model is trained with RL at a time, while the other is held constant. Further, we intertwine supervised training into the process as well so that the model does not deviate from the English language. We do so after every four iterations of RL, with a small learning rate to not take away from the RL process. The code for

---

[2]Code 14 contains the code that generates a dialogue between two models, which is used in this process.
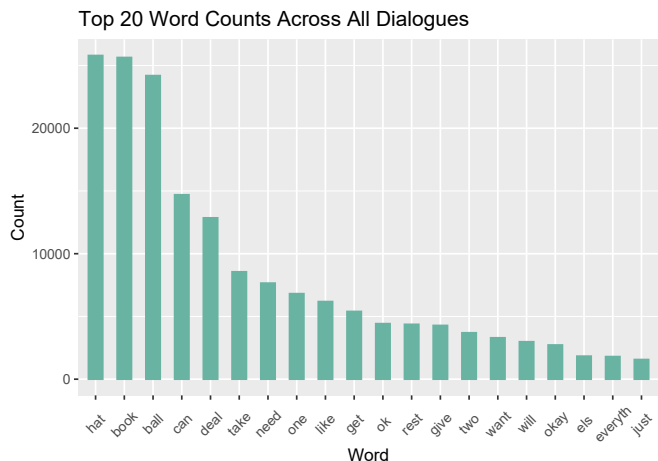
Figure 3: Top 20 Words and their Count

The figure shows the most common words used across both players throughout all negotiation dialogues.

the RL training process can be found in Code 15. The code for the agent/player that trains itself by playing in the mentioned RL training process can be found in Code 16.
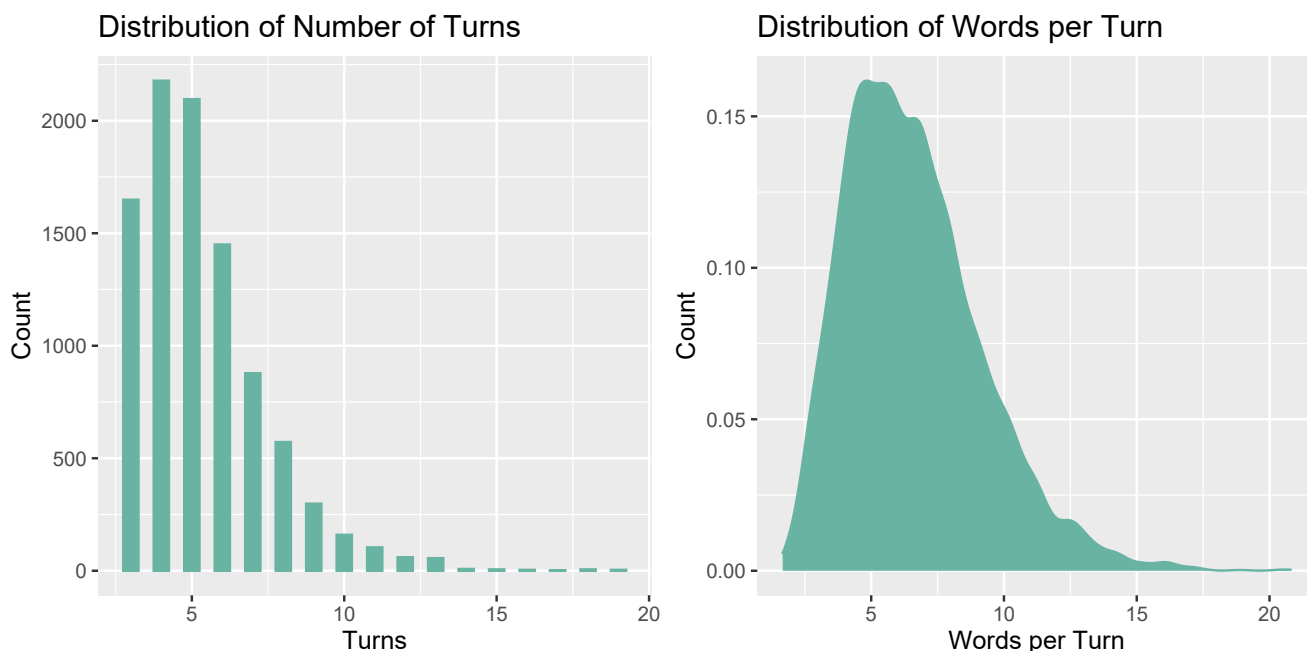
## 3.7   Testing Setting

Testing the models is very different than a typical supervised learning approach (e.g., classifier accuracy). If we were to evaluate our models by how well they perform on the dialogues we have, this would simply give a measure of similarity between our model's predictions and a human's sentences. However, we do not wish for our model to simply imitate humans, we want it to understand how to negotiate and maximize its reward. This gives it more of an RL evaluation than a supervised learning evaluation. To evaluate each model's performance we put it head-to-head against the most basic model we have: the supervised Transformer model. These models are given 8,172 different scenarios (i.e., items and value functions) and engage in dialogue. The scores are recorded, along with if there was agreement, and if the outcome is Pareto optimal. These scores are then aggregated across all scenarios to see how the proposed model performs relative to a baseline model. Here we can also choose to apply rollouts, MCTS, or neither to supplement the model that is undergoing evaluation. This is the same method that Lewis et al. [5] used to evaluate their models.

## 4   Results

### 4.1   Exploratory Data Analysis

We first look at the most common words used in the dialogues. This helps us determine the level of sophistication of the negotiations. Figure 3 shows the top 20 words used in the dialogues along with their respective word count (see Code 7 for the `R` code). Quite logically we see the top three words are the items that the players are trying to divide. The remaining words all seem typical of a negotiation dialogue with words like *need* and *deal*. The level of sophistication of these dialogues is quite low since we see no complex words; it is all very simple English. This means the model should be able to learn how to negotiate quite easily since the complexity of the language is low.

To understand the level of sophistication further, we view the distribution of the number of turns and the number of words used per turn. Figure 4 shows that most negotiations terminate in less than the maximum 20 turns allowed (see Code 8 for the R code). The majority of negotiations
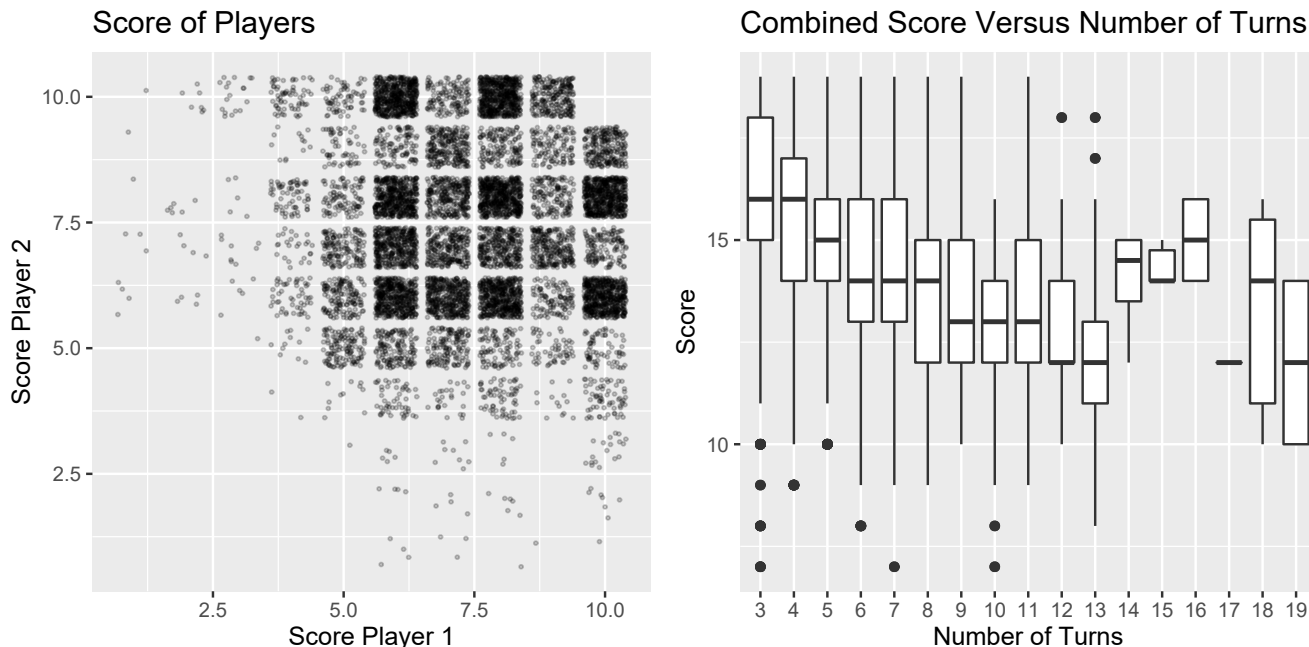
**Left:** number of turns taken to complete a negotiation. **Right:** average number of words used per turn in a dialogue.

Figure 4: Distribution of Turns and Words per Turn

finish in less than seven turns. This supports the previous conclusions that the negotiations cannot be extremely complex if they only take a few turns to arrive at an agreement. Further, Figure 4 also shows that there are very few words used per turn. With less than 10 words used per turn in most cases, it is unlikely that these negotiations are highly sophisticated. This supports the previous conclusion and indicates that the model should be able to learn the dialogues relatively easily.

Next, we try and understand the score of the agents better. First, we compare the score of each player in Figure 5 where we added *jitter* to the points so we can see the density at each score combination (see Code 9 for the R code). Players generally do not need to sacrifice their score for the other player to achieve a higher score. The plot shows that it is possible for both players to achieve a high score simultaneously (although not once did both players score a perfect 10). The most common scores are when both players score between six and eight. We also try and understand if the score of the game will increase as there are more turns. We expect this relationship since the players can better communicate their needs with more turns and get higher scores. The right side of Figure 5 shows the combined scores of both players versus the length of the dialogue. We actually see the inverse relationship than we expected; as there are more turns, the average combined score is lower. This could possibly be explained by irrational human behaviour such as the *endowment effect* where people think that the things they perceive to own are worth more and are less willing to give them up [3]. This result is counterintuitive and has implications for our model. We may seek to develop a model that attempts to end the negotiation as quickly as possible since that seems to produce better results.

Finally, we compare the score of a player to their positions in the dialogue (i.e., first or second). We expect that the player who goes first will have a higher score since they can set a reference point in the dialogue that all discussions will be based around (known as an anchoring point in

**Left:** comparison of scores of each player. **Right:** how scores vary as the length of dialogue varies.

Figure 5: Score Distributions

behavioural economics [1]). Figure 6 shows that the first player performs better in most cases, having a higher chance of achieving a higher score (see Code 10 for the R code). Further, the first player tends to achieve a perfect score at twice the rate of the second player. This finding agrees with our prior understanding of one-on-one negotiations.

## 4.2    Main Data Analysis

We carry out the training and testing steps mentioned in §3. This results in four models: (1) Transformer, (2) Transformer + RL, (3) Base GRU, and (4) Base GRU + RL. The baseline model we will use as our *default* agent to negotiate against other models is model (1). The default agent is compared to both model (1) and model (2) while using rollouts, MCTS, or no modifications on the latter two models. This gives us six groups of dialogues to analyze. Thereafter we compare model (3) and model (4), both with MCTS, to the default agent. This gives an additional two groups of dialogues to analyze, giving eight groups in total. These eight groups correspond to the eight new models we proposed, each compared to the default agent.

These groups of dialogues are individually analyzed to get the average score of each agent, the average score of each agent conditioned that the agents agreed, the agreement rate, and the proportion of Pareto optimal outcomes. These results are presented Table 3.

The model that performed the best in almost all metrics is the Transformer model that uses rollouts when competing against other agents. This is an unexpected result since we thought the best model would utilize MCTS given the algorithm's success in other domains.

We note that the RL is not successful at all and seems to be working in favour of the opponent instead of itself. Comparing any model without RL to its equivalent with RL shows that there is a significant score decrease. This is quite strange and we propose some reasons for this phenomenon
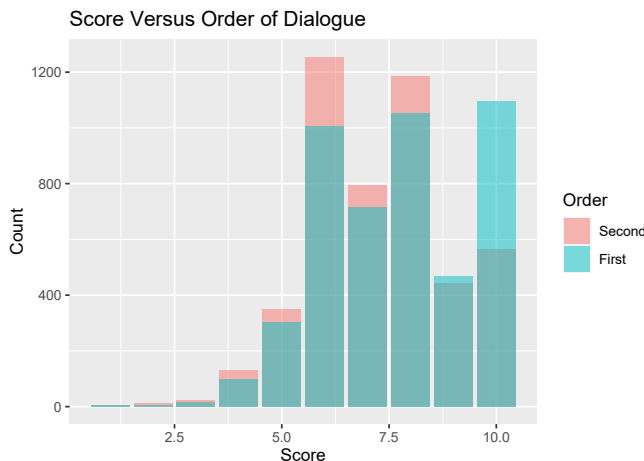
9

Score Versus Order of Dialogue

Figure 6: Distribution of Score versus Starting Player

This chart compares the number of times a player received a certain score based on if they were the first player to initiate the dialogue or the second player.

shortly. Also note that for all the Transformer models, the addition of RL increases the agreement rate significantly. It seems that the model is choosing to agree easier, at the expense of its score. This is a unique result but not the intention of using RL.

A possible flaw in the RL version of the model is the lack of hyperparameter tuning. Hyperparameters are just as important as in the supervised case but training takes much longer, making it difficult to tune. Each supervised model takes around 10-20 minutes to train, whereas an RL model takes around three hours. It is not feasible to train many RL models with different hyperparameters and assess which is the best. With more powerful devices, we may have been able to tune the parameters better for RL and get a better performing RL model, but we, unfortunately, did not have access to powerful computing resources.

We also note that all the agreement rates are much higher than those of the models proposed by Lewis et al. [5]. Since the part of the model responsible for generating sentences (encoder and decoder) is the same in our architecture and theirs, this difference can be attributed to the part of the model that makes the item selection. We can think of this as every dialogue being the exact same between our model and their model, but the item selections are the only difference. The difference in agreement rates is only attributable to the model being able to understand the dialogue and come up with the right item selection numbers. The Transformer is better able to predict the item selections given the dialogue compared to the GRU since it has a significantly higher agreement rate.

## 5    Conclusion and Discussion

We find that our model works well in dialogues and produces real English sentences. Along with generating coherent sentences, the model also negotiates quite well and can achieve a high score while not sacrificing the opponent's score. There is quite a high ratio of Pareto optimal outcomes, indicating the negotiations are equally beneficial for both parties, which is a preferred outcome. The model is also very agreeable, improving the rate greatly compared to previous approaches.

The second major finding is that MCTS does not provide a significant benefit over rollouts. Comparing the TRANSFORMER + ROLLOUTS to TRANSFORMER + MCTS we see that the rollouts perform better in most cases. We hypothesize that this is the case because the length of the game is quite short, averaging around seven turns per dialogue. MCTS was proposed to sample a large and deep action space as a way to approximate the value of a node in a dense game tree. Since these

| Model | vs. LIKELIHOOD | | | |
| --- | --- | --- | --- | --- |
| | Score (all) | Score (agreed) | % Agreed | % Pareto Optimal |
| TRANSFORMER LIKELIHOOD | 4.9 vs. 4.8 | 5.2 vs. 5.2 | 93.9 | 30.2 |
| TRANSFORMER + ROLLOUTS | **6.5 vs. 5.0** | **6.8 vs. 5.2** | 95.9 | **42.0** |
| TRANSFORMER RL | 3.7 vs. 6.6 | 3.7 vs. 6.5 | 98.9 | 36.6 |
| TRANSFORMER + RL + ROLLOUTS | 4.0 vs. 6.5 | 4.1 vs. 6.6 | **99.1** | 39.4 |
| TRANSFORMER + MCTS | 6.0 vs. 4.8 | 6.4 vs. 5.1 | 94.2 | 37.1 |
| BASE RNN + MCTS | 6.0 vs. 4.9 | 6.3 vs. 5.2 | 95.1 | 39.3 |
| BASE RNN + RL + MCTS | 6.0 vs. 5.0 | 6.3 vs. 5.2 | 94.9 | 40.8 |
| TRANSFORMER + RL + MCTS | 4.0 vs. 6.5 | 4.0 vs. 6.6 | 98.6 | 39.1 |

Table 3: Model Performance

dialogues are short the game tree is not too dense, meaning there are not many actions to sample from. Compared to rollouts there isn't much of a difference because rolling out a few candidates will cover a majority of the possible actions, making it more equivalent to MCTS.

MCTS is also meant to be used in scenarios where there are many potential (equally likely) moves. In our scenario, the dialogue for a given agent is essentially unimodal given the previous turns in the negotiation. This means that there are not many possible actions, rather there is one very likely action, with all other actions being very unlikely. This makes the dialogue somewhat deterministic. If we use rollouts, we will be sampling the most likely actions and this will inform us more about the likely outcome, compared to using MCTS and forcing the model to consider different actions when they are not all likely. MCTS considers many unlikely responses, meaning it is collecting unnecessary information about the action space and potentially misinforming itself about the most rewarding action.

Since we choose to sample the distribution of actions/responses to determine possible next moves, we are not actually considering all possible next turns. A possible turn would include any combination of any amount of the words in the game dictionary, meaning the possible action space would be massive. Even if combinations of words are not actual English, we should still consider them by MCTS theory, but we do not do this. We limit our game tree to have only the most probable combination of words, given the probability distribution. As mentioned previously, since the model is unimodal and almost deterministic there are not many possible actions to consider from the probability distribution. We found it almost impossible for the model to create five possible responses, given the previous turns, by randomly sampling the action space. Generating possible responses at random typically resulted in the same one or two responses being repeated over and over. This weakens MCTS since there are not many possible actions and makes it more similar to performing a rollout on a few different candidate responses (i.e., what Lewis et al. [5] implemented).

In Table 3 we also see that the TRANSFORMER + MCTS slightly outperforms the BASE RNN + MCTS in the agreed score. This shows that the Transformer has only slightly better performance than the original GRU network, but not as much of an improvement as we were hoping for. The two also

have similar agree rates, and the BASE RNN + MCTS has slightly higher Pareto optimality likely because it is willing to sacrifice its score to benefit its partner's score. This may be preferred in some cases, but here we hope to gain as much reward as possible, making the Transformer model preferred.

## 6    Future Works

The main adjustment we hope to make in this area in the future is using the Transformer model for understanding the dialogue better. Currently, we use the Transformer to generate the output selection, given the dialogue output from the GRU. We expected an improvement given the inherent sequence-to-sequence nature of the problem. Given that the dialogue generation is also a sequence-to-sequence problem, we believe the transition from a GRU to a Transformer would be beneficial. Further, one can also use Long Short-Term Memory (LSTM) cells instead of GRU cells. LSTMs have one more internal *gate* than a GRU unit, and thus tend to work better with longer sequences. The use of LSTMs or a Transformer might result in longer and more sophisticated dialogue than with GRUs.

We also suggest a different use for the MCTS algorithm be attempted in the future. Since the number of possible full-turn dialogues is limited and almost deterministic given the previous turns, we suggest performing MCTS with individual words, rather than with entire turns. This will generate many more possible turns for the model to consider, giving it a better evaluation of the action space. This will also increase the computational complexity though since the number of possible words to use at any point is quite large (above 400). We suggest limiting the scope of words to the 20 most likely words, rather than using the entire dictionary.

A final task we hope to take on in the future is comparing the model to humans. This is a real test of performance and will show if our model can negotiate better than humans can. We will create a web application for this paper, allowing people to compete against the different models we proposed. We will alternate the model that people play against and aggregate the scores for each model to see how well they perform against real people.

## 7   Appendix

The `Python` code used in this report is extensive so cannot be entirely included here. We include the most important `R` and `Python` code here and refer the reader here for the entirety of the code.

**List of Code**

Code 1: Import Packages

```r
library(knitr)
library(kableExtra)
library(ggplot2)
library(tm)
library(dplyr)
library(tidytext)
library(grid)
library(gridExtra)
```

Code 2: Load Data

```r
raw_data = scan("data.txt", what = character(), sep = "\n")
```

Code 3: Data Scraping Functions

```r
# function to get the score of a given dialogue
get_score = function(values, selections, agree){
  if (! agree) {return(0)}
  return(values %*% selections)
}


# function to get the following information from a dialogue:
#   number of items
#   value function for each player
#   Boolean variable indicating the players came to an agreement
#   the division of items between players
#   the score of each player
```

```r
13  #   the words in the dialogue
14  #   Boolean variable indicating the player had the first turn
15  #
16  # the function also checks for disconnected sessions and sessions where
17  # the players agreed to disagree
18  get_info = function(row){
19    split_vec = strsplit(row, " ")[[1]]
20    # get the number of items
21    item0 = split_vec[1]
22    item1 = split_vec[3]
23    item2 = split_vec[5]
24    # get the value function
25    value0 = split_vec[2]
26    value1 = split_vec[4]
27    value2 = split_vec[6]
28
29    n = length(split_vec)
30    # get the other player's value function
31    othervalue0 = split_vec[n-4]
32    othervalue1 = split_vec[n-2]
33    othervalue2 = split_vec[n]
34
35    # check if the players agreed
36    agree = split_vec[n-6] == "agree"
37
38    # check if the players agreed to disagree
39    no_agree = regmatches(row, regexpr("(?<=<selection>)(.*)(?=<eos>)", row, perl
      =T)) == " no agreement "
40    # check if the players disconnected
41    disconnect = regmatches(row, regexpr("(?<=<selection>)(.*)(?=<eos>)", row,
       perl=T)) == " disconnect "
42
43    if (no_agree | disconnect){
44      # give 0 score and no selection for disagreement and disconnection
45      selection = NA
46      score = 0
47      other_score = 0
48      agree = FALSE
49    } else {
50      # get the selection of player 1
51      selection0 = strsplit(regmatches(row, regexpr("item0=[0-9]+", row)), "=")
      [[1]][2]
52      selection1 = strsplit(regmatches(row, regexpr("item1=[0-9]+", row)), "=")
      [[1]][2]
53      selection2 = strsplit(regmatches(row, regexpr("item2=[0-9]+", row)), "=")
      [[1]][2]
54
55      selection = as.numeric(c(selection0, selection1, selection2))
56
57      # get the score of player 1
58      score = get_score(
59        as.numeric(c(value0,value1,value2)),
60        as.numeric(c(selection0, selection1, selection2)),
61        agree
62      )
63      # get the score of player 2
```

```
64    other_score = get_score(
65      as.numeric(c(othervalue0, othervalue1, othervalue2)),
66      as.numeric(c(item0, item1, item2)) - as.numeric(c(selection0, selection1,
        selection2)),
67      agree
68    )
69  }
70  # get all the words in the dialogue
71  dialogue = regmatches(row, regexpr("(THEM|YOU)(.+)(<selection>)", row))
72  # check if player 1 went first or second
73  went_first = split_vec[7] == "YOU:"
74
75  return(list(
76    item=as.numeric(c(item0, item1, item2)),
77    value=as.numeric(c(value0,value1,value2)),
78    othervalue=as.numeric(c(othervalue0, othervalue1, othervalue2)),
79    agree = agree,
80    selection = selection,
81    dialogue = dialogue,
82    score = score,
83    other_score = other_score,
84    went_first = went_first
85  ))
86 }
```

Code 4: Dataframe Creation

```
1 # use the get_info function to turn the dialogues into a dataframe
2 df = as.data.frame(t(sapply(raw_data, get_info, USE.NAMES = F)))
3 # format columns to be vectors instead of lists
4 df$agree = as.logical(df$agree)
5 df$dialogue = as.character(df$dialogue)
6 df$score = as.numeric(df$score)
7 df$other_score = as.numeric(df$other_score)
8 df$went_first = as.logical(df$went_first)
```

Code 5: Supplemental Functions

```
1 # get number of turns
2 get_turns = function(dialogue){
3   people = c("YOU:", "THEM:")
4   split_vec = strsplit(dialogue, " ")[[1]]
5   turns = 0
6   for (i in 1:length(split_vec)){
7     if (split_vec[i] %in% people) {
8       turns = turns + 1
9     }
10  }
11  return(turns)
12 }
13
14 # get the total number of words
15 get_words = function(dialogue){
16   ignore = c("YOU:", "THEM:", ".", "<eos>")
17   split_vec = strsplit(dialogue, " ")[[1]]
18   words = 0
19   for (i in 1:length(split_vec)){
```

```r
20      if (! split_vec[i] %in% ignore) {
21        words = words + 1
22      }
23    }
24    return(words)
25  }
26
27  # get possible item divisions
28  get_choices = function(items){
29    ind = 0
30    res = list()
31    for (i in 0:items[1]){
32      for (j in 0:items[2]){
33        for(k in 0:items[3]){
34          ind = ind + 1
35          res[[ind]] = c(i, j, k)
36        }
37      }
38    }
39    return(res)
40  }
41
42  # get pareto optimal
43  is_pareto = function(row){
44    agree = row$agree
45    items = row$item
46    score = row$score
47    other_score = row$other_score
48    values = row$value
49    other_values = row$othervalue
50
51    if (! agree) {return(FALSE)}
52
53    choices = get_choices(items)
54    # iterate through every possible item division in the game
55    for (i in 1:length(choices)){
56      # calculate the scores for the current division
57      curr_choice = choices[[i]]
58      curr_opp_choice = items - curr_choice
59      curr_score = get_score(values, curr_choice, T)
60      curr_opp_score = get_score(other_values, curr_opp_choice, T)
61      # check if both players could be better off
62      if ((curr_score > score & curr_opp_score >= other_score) |
63          (curr_score >= score & curr_opp_score > other_score)){
64        return(FALSE)
65      }
66    }
67    return(TRUE)
68  }
```

Code 6: Supplement Dataframe

```r
1  # use new functions to add new columns to the data frame
2  df$n_turns = sapply(df$dialogue, get_turns, USE.NAMES = F)
3  df$n_words = sapply(df$dialogue, get_words, USE.NAMES = F)
4  df$avg_words = df$n_words / df$n_turns
5  df$pareto = apply(df, 1, is_pareto)
```

Code 7: Word Analysis

```
1  # use corpus to manipulate the words
2  corpus = VCorpus(VectorSource(df$dialogue))
3  corpus = tm_map(corpus, removeNumbers)
4  corpus = tm_map(corpus, content_transformer(tolower))
5  # remove stopwords and end of sentence and selection marker
6  corpus = tm_map(corpus, removeWords, c(stopwords('english'), "eos", "selection"
       ))
7  corpus = tm_map(corpus, removePunctuation)
8  corpus = tm_map(corpus, stripWhitespace)
9  corpus = tm_map(corpus, stemDocument)
10
11 # turn the cleaned corpus into a dataframe
12 dialogue_df = data.frame(text = sapply(corpus, as.character), stringsAsFactors=
       F)
13
14 # get the count of each word in every dialogue
15 word_df = dialogue_df %>% unnest_tokens(word, text) %>% count(word, sort = TRUE
       )
16 word_df$word = factor(word_df$word, levels = word_df$word[order(word_df$n,
       decreasing = T)])
17
18 # plot the count of the top 20 words in all dialogues
19 ggplot(word_df[1:20, ], aes(y=n, x=word)) +
20   geom_bar(colour="#69b3a2", fill="#69b3a2", stat = "identity", width=0.5) +
21   ggtitle("Top 20 Word Counts Across All Dialogues") +
22   xlab("Word") +
23   ylab("Count") +
24   theme(axis.text.x = element_text(angle = 45, vjust = 0.5))
```

Code 8: Plot of Turns and Words per Turn

```
1  # create a separate dataframe for plotting
2  plot_data = df
3  # convert the went_first column to a factor
4  plot_data$went_first = factor(plot_data$went_first, labels = c("Second", "First
       "))
5
6  # plot the distribution of number of turns
7  p1 = plot_data[plot_data$agree, ] %>%
8    ggplot(aes(x=n_turns)) +
9    geom_bar(colour="#69b3a2", fill="#69b3a2", width = 0.5) +
10   ggtitle("Distribution of Number of Turns") +
11   xlab("Turns") +
12   ylab("Count")
13
14 # plot the distribution of words per turn
15 p2 = plot_data[plot_data$agree, ] %>%
16   ggplot(aes(x=avg_words)) +
17   geom_density(colour="#69b3a2", fill="#69b3a2") +
18   ggtitle("Distribution of Words per Turn") +
19   xlab("Words per Turn") +
20   ylab("Count")
21
22 grid.arrange(p1,p2, ncol=2)
```

Code 9: Plot of Score per Player and Score versus Turns

```r
# plot the score of each player
# adding scatter so we can see the overlap
p1 = plot_data[plot_data$agree, ] %>%
  ggplot(aes(x=score, y=other_score)) +
  geom_jitter(alpha = 0.2, cex = 0.5) +
  ggtitle("Score of Players") +
  xlab("Score Player 1") +
  ylab("Score Player 2")

# plot the distribution of score versus number of turns
p2 = plot_data[plot_data$agree, ] %>%
  ggplot(aes(x=as.factor(n_turns), y = score+other_score)) +
  geom_boxplot() +
  ggtitle("Combined Score Versus Number of Turns") +
  xlab("Number of Turns") +
  ylab("Score")

grid.arrange(p1,p2, ncol=2)
```

Code 10: Plot of Score versus Dialogue Order

```r
# plot the score versus the dialogue order
plot_data[plot_data$agree, ] %>%
  ggplot(aes(x = score, fill = went_first))+
  geom_bar(alpha=0.5, position = "identity")+
  ggtitle("Score Versus Order of Dialogue")+
  ylab("Count") +
  xlab("Score")+
  scale_fill_discrete(name="Order")
```

Code 11: Transformer Model Class

```python
import torch
import torch.nn as nn
import torch.nn.functional as F
import torch.nn.init
from torch.autograd import Variable
import torch.nn.functional as F

import data
from engines.rnn_engine import RnnEngine
from domain import get_domain
from models.utils import *
from models.ctx_encoder import MlpContextEncoder

class RnnModel(nn.Module):
    corpus_ty = data.WordCorpus
    engine_ty = RnnEngine

    def __init__(self, word_dict, item_dict, context_dict, count_dict, args):
        super(RnnModel, self).__init__()
        domain = get_domain(args.domain)
        self.word_dict = word_dict
        self.item_dict = item_dict
        self.context_dict = context_dict
        self.count_dict = count_dict
```

```
25        self.args = args
26
27        self.word_encoder = nn.Embedding(len(self.word_dict), args.nembed_word)
28        self.word_encoder_dropout = nn.Dropout(args.dropout)
29
30        ctx_encoder_ty = MlpContextEncoder
31        self.ctx_encoder = nn.Sequential(
32            ctx_encoder_ty(
33                len(self.context_dict),
34                domain.input_length(),
35                args.nembed_ctx,
36                args.nhid_ctx,
37                args.dropout,
38                args.init_range,
39            ),
40            nn.Dropout(args.dropout),
41        )
42
43        self.reader = nn.GRU(
44            args.nhid_ctx + args.nembed_word, args.nhid_lang, bias=True
45        )   # h_t
46        self.reader_dropout = nn.Dropout(args.dropout)
47
48        self.decoder = nn.Sequential(
49            nn.Linear(args.nhid_lang, args.nembed_word), nn.Dropout(args.
    dropout)
50        )
51
52        self.writer = nn.GRUCell(
53            input_size=args.nhid_ctx + args.nembed_word,
54            hidden_size=args.nhid_lang,
55            bias=True,
56        )
57
58        # Tie the weights of reader and writer
59        self.writer.weight_ih = self.reader.weight_ih_l0
60        self.writer.weight_hh = self.reader.weight_hh_l0
61        self.writer.bias_ih = self.reader.bias_ih_l0
62        self.writer.bias_hh = self.reader.bias_hh_l0
63
64        self.sel_rnn = nn.TransformerEncoderLayer(
65            d_model=args.nhid_lang + args.nembed_word, nhead=4, dropout=args.
    dropout
66        )
67
68        self.sel_dropout = nn.Dropout(args.dropout)
69
70        # Mask for disabling special tokens when generating sentences
71        self.special_token_mask = torch.FloatTensor(len(self.word_dict))
72
73        self.sel_encoder = nn.Sequential(
74            torch.nn.Linear(
75                args.nhid_lang + args.nembed_word + args.nhid_ctx, args.
    nhid_sel
76            ),
77            nn.Tanh(),
```

```
78                nn.Dropout(args.dropout),
79            )
80
81            self.attn = nn.Sequential(
82                torch.nn.Linear(args.nhid_lang + args.nembed_word, args.nhid_attn),
83                nn.Tanh(),
84                torch.nn.Linear(args.nhid_attn, 1),
85            )
86
87            self.sel_decoders = nn.ModuleList()
88            for i in range(domain.selection_length()):
89                self.sel_decoders.append(nn.Linear(args.nhid_sel, len(self.
        item_dict)))
90
91            self.init_weights()
92
93            self.special_token_mask = make_mask(
94                len(word_dict),
95                [word_dict.get_idx(w) for w in ["<unk>", "YOU:", "THEM:", "<pad>"
        ]],
96            )
97
98        def flatten_parameters(self):
99            self.reader.flatten_parameters()
100            self.sel_rnn.flatten_parameters()
101
102        def zero_h(self, bsz, nhid=None, copies=None):
103            nhid = self.args.nhid_lang if nhid is None else nhid
104            copies = 1 if copies is None else copies
105            h = torch.Tensor(copies, bsz, nhid).fill_(0)
106            return Variable(h)
107
108        def word2var(self, word):
109            x = torch.Tensor(1).fill_(self.word_dict.get_idx(word)).long()
110            return Variable(x)
111
112        def init_weights(self):
113            init_rnn(self.reader, self.args.init_range)
114            init_cont(self.decoder, self.args.init_range)
115            self.word_encoder.weight.data.uniform_(
116                -self.args.init_range, self.args.init_range
117            )
118
119            init_cont(self.attn, self.args.init_range)
120            init_cont(self.sel_encoder, self.args.init_range)
121            init_cont(self.sel_decoders, self.args.init_range)
122
123        def forward_context(self, ctx):
124            ctx_h = self.ctx_encoder(ctx).unsqueeze(0)
125            return ctx_h
126
127        def forward_lm(self, inpt_emb, lang_h, ctx_h):
128            # append the context embedding to every input word embedding
129            ctx_h_rep = ctx_h.narrow(0, ctx_h.size(0) - 1, 1).expand(
130                inpt_emb.size(0), ctx_h.size(1), ctx_h.size(2)
131            )
```

20

```python
132            inpt_emb = torch.cat([inpt_emb, ctx_h_rep], 2)
133            # n_words x (ctx_h size {64} + inpt size {256})
134            # combine word embeddings and context (g) embeddings
135
136            lang_hs, _ = self.reader(inpt_emb, lang_h)
137            lang_hs = self.reader_dropout(lang_hs)  # n_words x 128 (nhid_lang)
138
139            decoded = self.decoder(lang_hs.view(-1, lang_hs.size(2)))
140            # n_words x 256 (nembed_word)
141
142            out = F.linear(decoded, self.word_encoder.weight)
143            # n_words x 463 (encoder size)
144
145            return out, lang_hs
146
147     def forward_selection(self, inpt_emb, lang_h, ctx_h):
148            # run a birnn over the concatenation of the input embeddings and
      language model hidden states
149            h = torch.cat([lang_h, inpt_emb], 2)
150            # n_words x (256 + 128)
151            # combine word imbeddings and output of second GRU
152
153            attn_h = self.zero_h(h.size(1), self.args.nhid_attn, copies=2)
154            # 2 x 64 (nhid_attn)
155            # initial hidden state of the third GRU
156            # n_words x 128 (nhid_attn x 2) b/c bidirectional
157
158            h = self.sel_rnn(h)
159
160            # h = self.sel_dropout(h)
161
162            h = h.transpose(0, 1).contiguous()
163            # batch_size x n_words x 128
164
165            logit = self.attn(h.view(-1, self.args.nhid_lang + self.args.
      nembed_word)).view(
166                h.size(0), h.size(1)
167            )
168            # first removes the batch dimension
169            # then calculates attn (Linear, Tanh, Linear)
170            # batch_size x n_words
171            # this is h_t^a
172
173            prob = F.softmax(logit, dim=1).unsqueeze(2).expand_as(h)
174            # batch_size x n_words x 128
175            # these are all the same along the last dimension
176            # these are the alpha_t
177
178            attn = (
179                torch.sum(torch.mul(h, prob), 1, keepdim=True).transpose(0, 1).
      contiguous()
180            )
181            # 1 x 1 x 128
182
183            h = torch.cat([attn, ctx_h], 2).squeeze(0)
184            # concatenate the context vector (g) and the attention scores
```

```
185         # 1 x 1 x 192
186
187         h = self.sel_encoder.forward(h)
188         # this is h_s
189         # batch_size x 128
190
191         # there are 6 decoders in self.sel_decoders
192         # each is a linear with output size of 18
193
194         outs = [decoder.forward(h) for decoder in self.sel_decoders]
195         out = torch.cat(outs, 0)
196         # this is equation (9) in the paper
197         # 6 x 18
198         return out
199
200     def forward(self, inpt, ctx):
201         # ctx is 6 x 1
202         # inpt is n_words x 1
203         ctx_h = self.forward_context(ctx)  # 1 x 1 x 64
204
205         lang_h = self.zero_h(
206             ctx_h.size(1), self.args.nhid_lang
207         )  # 1 x 1 x 128 all zeros
208         # initial hidden state of the second GRU
209
210         inpt_emb = self.word_encoder(inpt)  # n_words x 256
211         inpt_emb = self.word_encoder_dropout(inpt_emb)
212
213         out, lang_hs = self.forward_lm(inpt_emb, lang_h, ctx_h)
214         sel_out = self.forward_selection(inpt_emb, lang_hs, ctx_h)
215         return out, sel_out
```

Code 12: Monte Carlo Tree Search Agent Class

```
1  class RnnMCTSAgent(RnnAgent):
2      def __init__(self, model,args, name="Alice", train=False, diverse=False):
3          super(RnnMCTSAgent, self).__init__(model, args, name)
4          # add variables to store the number of MCTS
5          self.nsim = 50
6          self.rollout_len = 100
7          self.n_tries = 5
8
9      def pickNode(self, curr_node: Node):
10         '''
11         picks max ucb node
12         '''
13         max_ucb = -math.inf
14         selection = None
15         for child in curr_node.children:
16             curr_ucb = child.get_ucb()
17             if curr_ucb > max_ucb:
18                 max_ucb = curr_ucb
19                 selection = child
20         return selection
21
22     def expansion(self, curr_node: Node):
23         """
```

```python
24          goes through all nodes until it reaches a leaf
25          """
26          if not curr_node.children:
27              return curr_node
28          selection = self.pickNode(curr_node)
29          return self.expansion(selection)
30
31
32      def backprop(self, curr_node: Node, score: float):
33          """
34          backpropagate scores throughout the tree
35          """
36          while True:
37              curr_node.score += score
38              curr_node.n += 1
39              if not curr_node.parent:
40                  break
41              curr_node = curr_node.parent
42          return curr_node
43
44      def get_all_children_states(self, curr_node: Node, n_tries: int):
45          """
46          get a few possible dialogues given the current state
47          """
48          children = []
49          prev_moves = []
50          for _ in range(n_tries):
51              _, move, move_lang_h, move_lang_hs = self.model.write(
52                  curr_node.lang_h, self.ctx_h, self.rollout_len, self.args.
    temperature)
53              is_selection = len(move) == 1 and \
54                  self.model.word_dict.get_word(move.data[0][0]) == '<selection>'
55              if not any([torch.equal(move,x) for x in prev_moves]):
56                  children.append(Node(parent = curr_node, lang_h=move_lang_h,
    lang_hs = move_lang_hs, move = move, sel=is_selection))
57                  prev_moves.append(move)
58          return children
59
60      def prepare_words(self, combined_words):
61          """
62          get words in the right format to be passed to _choose1
63          """
64          res = []
65          start = 0
66          end = 0
67          for i in range(combined_words.size(0)):
68              if combined_words[i, :] == 0:
69                  end = i + 1
70                  res.append(combined_words[start:end, :])
71                  start = end
72          return res
73
74      def _choose1(self, sents, sample=False):
75          lens, rev_idxs, hid_idxs = self._make_idxs(sents)
76          sel_out = self.sel_model.forward(sents, lens, rev_idxs, hid_idxs,
    Variable(self.ctx))
```

```python
77
78         choices = self.domain.generate_choices(self.context, with_disagreement=
      True)
79
80         choices_logits = []
81         for i in range(self.domain.selection_length()):
82             idxs = [self.sel_model.item_dict.get_idx(c[i]) for c in choices]
83             idxs = Variable(torch.Tensor(idxs).long())
84             choices_logits.append(torch.gather(sel_out[i], 0, idxs).unsqueeze
      (1))
85
86         choice_logit = torch.sum(torch.cat(choices_logits, 1), 1, keepdim=True)
      .squeeze(1)
87         choice_logit = choice_logit.sub(choice_logit.max(0)[0].item())
88         prob = F.softmax(choice_logit, dim=0)
89
90         if sample:
91             idx = prob.multinomial(1).detach()
92             logprob = F.log_softmax(choice_logit, dim=0).gather(0, idx)
93         else:
94             _, idx = prob.max(0, keepdim=True)
95             logprob = None
96
97         p_agree = prob[idx.item()]
98
99         # Pick only your choice
100        return choices[idx.item()][:self.domain.selection_length()], logprob,
      p_agree.item()
101
102    def write(self, max_words):
103        """
104        new write method for generating turn
105        """
106        root = Node(lang_h = self.lang_h)
107        root.children = self.get_all_children_states(root, self.n_tries)
108
109        for _ in range(self.nsim):
110            score = 0
111            '''
112            get max node for traversal
113            '''
114            candidate = self.expansion(root)
115
116            if candidate.n != 0:
117                candidate.children = self.get_all_children_states(candidate,
      self.n_tries)
118                candidate = candidate.children[0]
119
120
121            combined_words = self.words + [self.model.word2var('YOU:'),
      candidate.outs]
122            # print(combined_words)
123            if not candidate.is_selection:
124                """
125                do full rollout and get the score if its not a terminal node
126                """
```

24

```
127                     _, rollout, _, rollout_lang_hs = self.model.write(
128                             candidate.lang_h, self.ctx_h, self.rollout_len, self.
    args.temperature,
129                             stop_tokens=['<selection>'], resume=True)
130                 combined_words += [rollout]
131
132             combined_words = [(lambda x: x.unsqueeze(1) if len(x.size()) == 1
    else x)(x) for x in combined_words]
133             combined_words = torch.cat(combined_words, dim = 0)
134             combined_words = self.prepare_words(combined_words)
135             rollout_choice, _, p_agree = self._choose1(sents=combined_words,
    sample=False)
136             rollout_score = self.domain.score(self.context, rollout_choice)
137             # score += p_agree * rollout_score
138             score += rollout_score
139
140             self.backprop(candidate, score)
141
142         max_score = -math.inf
143         for child in root.children:
144             if child.score / child.n > max_score:
145                 bestChild = child
146                 max_score = child.score / child.n
147         self.lang_h = bestChild.lang_h
148         self.lang_hs.append(bestChild.lang_hs)
149         self.words.append(self.model.word2var('YOU:'))
150         self.words.append(bestChild.outs)
151         self.sents.append(torch.cat([self.model.word2var('YOU:').unsqueeze(1),
    bestChild.outs], 0))
152         return self._decode(bestChild.outs, self.model.word_dict)
```

Code 13: Hyperparameter Tuning

```
1  # define function to maximize
2  def objective(trial: optuna.trial.Trial):
3      # add arguments to be read from the command line
4      parser = argparse.ArgumentParser(description='training script')
5      parser.add_argument('--data', type=str, default='data/negotiate',
6          help='location of the data corpus')
7      parser.add_argument('--nembed_word', type=int, default=256,
8          help='size of word embeddings')
9      parser.add_argument('--nembed_ctx', type=int, default=64,
10         help='size of context embeddings')
11     parser.add_argument('--nhid_lang', type=int, default=256,
12         help='size of the hidden state for the language module')
13     parser.add_argument('--nhid_cluster', type=int, default=256,
14         help='size of the hidden state for the language module')
15     parser.add_argument('--nhid_ctx', type=int, default=64,
16         help='size of the hidden state for the context module')
17     parser.add_argument('--nhid_strat', type=int, default=64,
18         help='size of the hidden state for the strategy module')
19     parser.add_argument('--nhid_attn', type=int, default=64,
20         help='size of the hidden state for the attention module')
21     parser.add_argument('--nhid_sel', type=int, default=64,
22         help='size of the hidden state for the selection module')
23     parser.add_argument('--lr', type=float, default=20.0,
24         help='initial learning rate')
```

```
25     parser.add_argument('--min_lr', type=float, default=1e-5,
26         help='min threshold for learning rate annealing')
27     parser.add_argument('--decay_rate', type=float,  default=9.0,
28         help='decrease learning rate by this factor')
29     parser.add_argument('--decay_every', type=int,  default=1,
30         help='decrease learning rate after decay_every epochs')
31     parser.add_argument('--momentum', type=float, default=0.0,
32         help='momentum for sgd')
33     parser.add_argument('--clip', type=float, default=0.2,
34         help='gradient clipping')
35     parser.add_argument('--dropout', type=float, default=0.5,
36         help='dropout rate in embedding layer')
37     parser.add_argument('--init_range', type=float, default=0.1,
38         help='initialization range')
39     parser.add_argument('--max_epoch', type=int, default=30,
40         help='max number of epochs')
41     parser.add_argument('--num_clusters', type=int, default=50,
42         help='number of clusters')
43     parser.add_argument('--bsz', type=int, default=25,
44         help='batch size')
45     parser.add_argument('--unk_threshold', type=int, default=20,
46         help='minimum word frequency to be in dictionary')
47     parser.add_argument('--temperature', type=float, default=0.1,
48         help='temperature')
49     parser.add_argument('--partner_ctx_weight', type=float, default=0.0,
50         help='selection weight')
51     parser.add_argument('--sel_weight', type=float, default=0.6,
52         help='selection weight')
53     parser.add_argument('--seed', type=int, default=1,
54         help='random seed')
55     parser.add_argument('--cuda', action='store_true', default=False,
56         help='use CUDA')
57     parser.add_argument('--model_file', type=str,  default='',
58         help='path to save the final model')
59     parser.add_argument('--prediction_model_file', type=str,  default='',
60         help='path to save the prediction model')
61     parser.add_argument('--selection_model_file', type=str,  default='',
62         help='path to save the selection model')
63     parser.add_argument('--cluster_model_file', type=str,  default='',
64         help='path to save the cluster model')
65     parser.add_argument('--lang_model_file', type=str,  default='',
66         help='path to save the language model')
67     parser.add_argument('--visual', action='store_true', default=False,
68         help='plot graphs')
69     parser.add_argument('--skip_values', action='store_true', default=False,
70         help='skip values in ctx encoder')
71     parser.add_argument('--model_type', type=str, default='rnn_model',
72         help='model type', choices=models.get_model_names())
73     parser.add_argument('--domain', type=str, default='object_division',
74         help='domain for the dialogue')
75     parser.add_argument('--clustering', action='store_true', default=False,
76         help='use clustering')
77     parser.add_argument('--sep_sel', action='store_true', default=False,
78         help='use separate classifiers for selection')
79
80     # grab the arguments from the command line
```

```
81      args = parser.parse_args()
82
83      # use GPU and set seed
84      utils.use_cuda(args.cuda)
85      utils.set_seed(args.seed)
86
87      # set the possible range for variables to tune
88      args.clip = trial.suggest_float("clip",0.25,0.75)
89      args.decay_every = trial.suggest_int("decay_every", 1, 5)
90      args.decay_rate = trial.suggest_int("decay_rate", 2, 10)
91      args.dropout = trial.suggest_float("dropout", 0.1, 0.9)
92      args.init_range = trial.suggest_float("init_range", 0.1, 0.5)
93      args.lr = trial.suggest_float("initial_learning_rate", 1e-5, 1e-2)
94      args.min_lr = trial.suggest_float("min_learning_rate", 1e-9,1e-6)
95      args.momentum = trial.suggest_float("momentum", 0, 1)
96      args.nembed_ctx = trial.suggest_categorical("ctx_embeding", [64,128,256])
97      args.nembed_word = trial.suggest_categorical("word_embeding", [64,128,256])
98      args.nhid_attn = trial.suggest_categorical("hidden_attn_size",
        [64,128,256])
99      args.nhid_ctx = trial.suggest_categorical("hidden_ctx_size", [64,128,256])
100     args.nhid_lang = trial.suggest_categorical("hidden_lang_size",
        [64,128,256])
101     args.nhid_sel = trial.suggest_categorical("hidden_selection_size",
        [64,128,256])
102     args.sel_weight = trial.suggest_float("selection_weight", 0.2, 0.8)
103
104     # set the domain of the game (i.e., dividing objects between players)
105     domain = get_domain(args.domain)
106     # get the model we use for training
107     model_ty = models.get_model_type(args.model_type)
108     # get the dialogues we train on
109     corpus = model_ty.corpus_ty(domain, args.data, freq_cutoff=args.
        unk_threshold,
110         verbose=True, sep_sel=args.sep_sel)
111     # initialize the model
112     model = model_ty(corpus.word_dict, corpus.item_dict_old,
113         corpus.context_dict, corpus.count_dict, args)
114     if args.cuda:
115         model.cuda()
116     # initialize the engine (the object that actually trains the model)
117     engine = model_ty.engine_ty(model, args, verbose=True)
118     # train the model
119     train_loss, valid_loss, select_loss, extra = engine.train(corpus)
120     # save the model
121     utils.save_model(engine.get_model(), args.model_file)
122
123     return valid_loss
124
125
126 # define how we will optimize our objective function
127 study = optuna.create_study(direction = 'minimize', sampler = optuna.samplers.
    TPESampler(seed=4850))
128 # find the best hyperparameters
129 study.optimize(objective, n_trials=1000)
130 # print the best parameters
131 print(study.best_trial)
```

Code 14: Class to Create Dialogues Between Two Players

```python
class Dialog(object):
    def __init__(self, agents, args):
        # For now we only suppport dialog of 2 agents
        assert len(agents) == 2
        self.agents = agents
        self.args = args
        self.domain = domain.get_domain(args.domain)
        self.metrics = MetricsContainer()
        self._register_metrics()

    def _register_metrics(self):
        self.metrics.register_average('dialog_len')
        self.metrics.register_average('sent_len')
        self.metrics.register_percentage('agree')
        self.metrics.register_moving_percentage('moving_agree')
        self.metrics.register_average('advantage')
        self.metrics.register_moving_average('moving_advantage')
        self.metrics.register_time('time')
        self.metrics.register_average('comb_rew')
        self.metrics.register_average('agree_comb_rew')
        for agent in self.agents:
            self.metrics.register_average('%s_rew' % agent.name)
            self.metrics.register_moving_average('%s_moving_rew' % agent.name)
            self.metrics.register_average('agree_%s_rew' % agent.name)
            self.metrics.register_percentage('%s_sel' % agent.name)
            self.metrics.register_uniqueness('%s_unique' % agent.name)
        # text metrics
        if self.args.ref_text:
            ref_text = ' '.join(data.read_lines(self.args.ref_text))
            self.metrics.register_ngram('full_match', text=ref_text)

    def _is_selection(self, out):
        return len(out) == 1 and (out[0] in ['<selection>', '<no_agreement>'])

    def show_metrics(self):
        return ' '.join(['%s=%s' % (k, v) for k, v in self.metrics.dict().items
    ()])

    def run(self, ctxs, logger, max_words=5000):
        self.agents[0].model.train()
        assert len(self.agents) == len(ctxs)
        for agent, ctx, partner_ctx in zip(self.agents, ctxs, reversed(ctxs)):
            agent.feed_context(ctx)
            agent.feed_partner_context(partner_ctx)
            logger.dump_ctx(agent.name, ctx)
        logger.dump('-' * 80)

        # Choose who goes first by random
        if np.random.rand() < 0.5:
            writer, reader = self.agents
        else:
            reader, writer = self.agents

        conv = []
        self.metrics.reset()
```

```python
55
56        #words_left = np.random.randint(50, 200)
57        words_left = max_words
58        length = 0
59        expired = False
60        while True:
61            out = writer.write(max_words=words_left)
62            words_left -= len(out)
63            length += len(out)
64
65            self.metrics.record('sent_len', len(out))
66            if 'full_match' in self.metrics.metrics:
67                self.metrics.record('full_match', out)
68            self.metrics.record('%s_unique' % writer.name, out)
69
70            conv.append(out)
71            reader.read(out)
72            if not writer.human:
73                logger.dump_sent(writer.name, out)
74
75            if self._is_selection(out):
76                self.metrics.record('%s_sel' % writer.name, 1)
77                self.metrics.record('%s_sel' % reader.name, 0)
78                break
79
80            if words_left <= 1:
81                break
82
83            writer, reader = reader, writer
84
85
86        choices = []
87        for agent in self.agents:
88            choice = agent.choose()
89            choices.append(choice)
90            logger.dump_choice(agent.name, choice[: self.domain.
    selection_length() // 2])
91
92        agree, rewards = self.domain.score_choices(choices, ctxs)
93        if expired:
94            agree = False
95        logger.dump('-' * 80)
96        logger.dump_agreement(agree)
97        for i, (agent, reward) in enumerate(zip(self.agents, rewards)):
98            logger.dump_reward(agent.name, agree, reward)
99            j = 1 if i == 0 else 0
100           agent.update(agree, reward, choice=choices[i],
101               partner_choice=choices[j], partner_input=ctxs[j],
    partner_reward=rewards[j])
102
103       if agree:
104           self.metrics.record('advantage', rewards[0] - rewards[1])
105           self.metrics.record('moving_advantage', rewards[0] - rewards[1])
106           self.metrics.record('agree_comb_rew', np.sum(rewards))
107           for agent, reward in zip(self.agents, rewards):
108               self.metrics.record('agree_%s_rew' % agent.name, reward)
```

```
109
110          self.metrics.record('time')
111          self.metrics.record('dialog_len', len(conv))
112          self.metrics.record('agree', int(agree))
113          self.metrics.record('moving_agree', int(agree))
114          self.metrics.record('comb_rew', np.sum(rewards) if agree else 0)
115          for agent, reward in zip(self.agents, rewards):
116              self.metrics.record('%s_rew' % agent.name, reward if agree else 0)
117              self.metrics.record('%s_moving_rew' % agent.name, reward if agree
      else 0)
118
119          logger.dump('-' * 80)
120          logger.dump(self.show_metrics())
121          logger.dump('-' * 80)
122          for ctx, choice in zip(ctxs, choices):
123              logger.dump('debug: %s %s' % (' '.join(ctx), ' '.join(choice)))
124
125          return conv, agree, rewards
```

Code 15: Class for Reinforcement Learning

```
1  class Reinforce(object):
2      def __init__(self, dialog, ctx_gen, args, engine, corpus, logger=None):
3          self.dialog = dialog
4          self.ctx_gen = ctx_gen
5          self.args = args
6          self.engine = engine
7          self.corpus = corpus
8          self.logger = logger if logger else DialogLogger()
9
10     def run(self):
11         validset, validset_stats = self.corpus.valid_dataset(self.args.bsz)
12         trainset, trainset_stats = self.corpus.train_dataset(self.args.bsz)
13
14         n = 0
15         for ctxs in self.ctx_gen.iter(self.args.nepoch):
16             n += 1
17             if self.args.sv_train_freq > 0 and n % self.args.sv_train_freq ==
      0:
18                 batch = random.choice(trainset)
19                 self.engine.model.train()
20                 self.engine.train_batch(batch)
21                 self.engine.model.eval()
22             self.logger.dump('=' * 80)
23             self.engine.model.train()
24             self.dialog.run(ctxs, self.logger)
25             self.logger.dump('=' * 80)
26             self.logger.dump('')
27             if n % 100 == 0:
28                 self.logger.dump('%d: %s' % (n, self.dialog.show_metrics()),
      forced=True)
29
30         def dump_stats(dataset, stats, name):
31             loss, select_loss = self.engine.valid_pass(N, dataset, stats)
32             self.logger.dump('final: %s_loss %.3f %s_ppl %.3f' % (
33                 name, float(loss), name, np.exp(float(loss))),
34                 forced=True)
```

```
35            self.logger.dump('final: %s_select_loss %.3f %s_select_ppl %.3f' %
       (
36                 name, float(select_loss), name, np.exp(float(select_loss))),
37                 forced=True)
38
39        dump_stats(trainset, trainset_stats, 'train')
40        dump_stats(validset, validset_stats, 'valid')
41
42        self.logger.dump('final: %s' % self.dialog.show_metrics(), forced=True)
```

Code 16: Agent Class for Reinforcement Learning

```
1  class RlAgent(RnnAgent):
2      def __init__(self, model, args, name='Alice', train=False):
3          self.train = train
4          super(RlAgent, self).__init__(model, args, name=name)
5          self.model.train()
6          self.sel_model.train()
7          self.opt = optim.RMSprop(
8              self.model.parameters(),
9              lr=args.rl_lr,
10             momentum=self.args.momentum)
11
12         self.all_rewards = []
13
14         if self.args.visual:
15             self.model_plot = vis.ModulePlot(self.model, plot_weight=False,
       plot_grad=True)
16             self.agree_plot = vis.Plot(['agree',], 'agree', 'agree')
17             self.reward_plot = vis.Plot(
18                 ['reward', 'partner_reward'], 'reward', 'reward')
19             self.loss_plot = vis.Plot(['loss',], 'loss', 'loss')
20             self.agree_reward_plot = vis.Plot(
21                 ['reward', 'partner_reward'], 'agree_reward', 'agree_reward')
22         self.t = 0
23
24     def feed_context(self, ctx):
25         super(RlAgent, self).feed_context(ctx)
26         self.logprobs = []
27
28     def write(self, max_words):
29         logprobs, outs, self.lang_h, lang_hs = self.model.write(self.lang_h,
       self.ctx_h,
30             100, self.args.temperature)
31         self.logprobs.extend(logprobs)
32         self.lang_hs.append(lang_hs)
33         self.words.append(self.model.word2var('YOU:').unsqueeze(0))
34         self.words.append(outs)
35         assert (torch.cat(self.words).size()[0] == torch.cat(self.lang_hs).size
       ()[0])
36         return self._decode(outs, self.model.word_dict)
37
38     def choose(self):
39         if self.args.eps < np.random.rand():
40             choice, _, _ = self._choose(sample=False)
41         else:
42             choice, logprob, _ = self._choose(sample=True)
```

```
43              self.logprobs.append(logprob)
44          return choice
45
46   def update(self, agree, reward, choice=None, partner_choice=None,
     partner_input=None, partner_reward=None):
47          if not self.train:
48              return
49
50          self.t += 1
51          if len(self.logprobs) == 0:
52              return
53          reward_agree = reward
54          partner_reward_agree = partner_reward
55
56          reward = reward if agree else 0
57          partner_reward = partner_reward if agree else 0
58
59          diff = reward - partner_reward
60          self.all_rewards.append(diff)
61          #self.all_rewards.append(reward)
62          r = (diff - np.mean(self.all_rewards)) / max(1e-4, np.std(self.
     all_rewards))
63          g = Variable(torch.zeros(1, 1).fill_(r))
64          rewards = []
65          for _ in self.logprobs:
66              rewards.insert(0, g)
67              g = g * self.args.gamma
68
69          loss = 0
70          for lp, r in zip(self.logprobs, rewards):
71              loss -= lp * r
72
73          self.opt.zero_grad()
74          loss.backward()
75          nn.utils.clip_grad_norm(self.model.parameters(), self.args.rl_clip)
76          if self.args.visual and self.t % 10 == 0:
77              self.model_plot.update(self.t)
78              self.agree_plot.update('agree', self.t, int(agree))
79              self.reward_plot.update('reward', self.t, reward)
80              self.reward_plot.update('partner_reward', self.t, partner_reward)
81              self.agree_reward_plot.update('reward', self.t, reward_agree)
82              self.agree_reward_plot.update('partner_reward', self.t,
     partner_reward_agree)
83              self.loss_plot.update('loss', self.t, loss.data[0][0])
84
85          self.opt.step()
```

## 8    References

[1] Adrian Furnham and Hua Chu Boo. A literature review of the anchoring effect. *The Journal of Socio-Economics*, 40(1):35–42, 2011.

[2] Sylvain Gelly, Levente Kocsis, Marc Schoenauer, Michèle Sebag, David Silver, Csaba Szepesvári, and Olivier Teytaud. The grand challenge of computer go. *Communications of the ACM*, 55(3):106–113, 2012.

[3] Daniel Kahneman, Jack L Knetsch, and Richard H Thaler. Anomalies: The endowment effect, loss aversion, and status quo bias. *Journal of Economic Perspectives*, 5(1):193–206, 1991.

[4] Shigeki Karita, Nanxin Chen, Tomoki Hayashi, Takaaki Hori, Hirofumi Inaguma, Ziyan Jiang, Masao Someki, Nelson Enrique Yalta Soplin, Ryuichi Yamamoto, Xiaofei Wang, Shinji Watanabe, Takenori Yoshimura, and Wangyou Zhang. A comparative study on transformer vs RNN in speech applications. *CoRR*, abs/1909.06317, 2019.

[5] Mike Lewis, Denis Yarats, Yann N. Dauphin, Devi Parikh, and Dhruv Batra. Deal or no deal? end-to-end learning for negotiation dialogues. *CoRR*, abs/1706.05125, 2017.

[6] Alan Strudler. On the ethics of deception in negotiation. *Business Ethics Quarterly*, 5(4):805–822, Oct 1995.

[7] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Conference on Neural Information Processing Systems*, 31, 2017.