

Deep Reinforcement Learning for Dialogue Generation

Jiwei Li¹, Will Monroe¹, Alan Ritter², Michel Galley³, Jianfeng Gao³ and Dan Jurafsky¹

¹Stanford University, Stanford, CA, USA

²Ohio State University, OH, USA

³Microsoft Research, Redmond, WA, USA

{jiweil, wmonroe4, jurafsky}@stanford.edu, ritter.1492@osu.edu

{mgalley, jfgao}@microsoft.com

Abstract

Recent neural models of dialogue generation offer great promise for generating responses for conversational agents, but tend to be short-sighted, predicting utterances one at a time while ignoring their influence on future outcomes. Modeling the future direction of a dialogue is crucial to generating coherent, interesting dialogues, a need which led traditional NLP models of dialogue to draw on reinforcement learning. In this paper, we show how to integrate these goals, applying deep reinforcement learning to model future reward in chatbot dialogue. The model simulates dialogues between two virtual agents, using policy gradient methods to reward sequences that display three useful conversational properties: informativity, coherence, and ease of answering (related to forward-looking function). We evaluate our model on diversity, length as well as with human judges, showing that the proposed algorithm generates more interactive responses and manages to foster a more sustained conversation in dialogue simulation. This work marks a first step towards learning a neural conversational model based on the long-term success of dialogues.

1 Introduction

Neural response generation (Sordoni et al., 2015; Shang et al., 2015; Vinyals and Le, 2015; Li et al., 2016a; Wen et al., 2015; Yao et al., 2015; Luan et al., 2016; Xu et al., 2016; Wen et al., 2016; Li et al., 2016b; Su et al., 2016) is of growing interest. The LSTM sequence-to-sequence (SEQ2SEQ) model (Sutskever et al., 2014) is one type of neural generation model that maximizes the probability of generating a response given the previous dialogue turn. This approach enables the incorporation of rich

context when mapping between consecutive dialogue turns (Sordoni et al., 2015) in a way not possible, for example, with MT-based dialogue models (Ritter et al., 2011).

Despite the success of SEQ2SEQ models in dialogue generation, two problems emerge: First, SEQ2SEQ models are trained by predicting the next dialogue turn in a given conversational context using the maximum-likelihood estimation (MLE) objective function. However, it is not clear how well MLE approximates the real-world goal of chatbot development: teaching a machine to converse with humans, while providing interesting, diverse, and informative feedback that keeps users engaged. One concrete example is that SEQ2SEQ models tend to generate highly generic responses such as “*I don’t know*” regardless of the input (Sordoni et al., 2015; Serban et al., 2016; Li et al., 2016a). This can be ascribed to the high frequency of generic responses found in the training set and their compatibility with a diverse range of conversational contexts. Yet “*I don’t know*” is apparently not a good action to take, since it closes the conversation down.

Another common problem, illustrated in the two sample conversations on the left of Table 1, is that the system becomes stuck in an infinite loop of repetitive responses. This is due to MLE-based SEQ2SEQ models’ inability to account for repetition. In example 2 (bottom left), the dialogue falls into an infinite loop after three turns, with both agents generating dull, generic utterances like *i don’t know what you are talking about* and *you don’t know what you are saying*. Looking at the entire conversation, utterance (4) turns out to be a bad action to take because it offers no way of continuing the conversation.¹

¹A similar rule is often suggested in improvisational comedy: https://en.wikipedia.org/wiki/Yes,_and...

Baseline mutual information model (Li et al. 2015)	Proposed reinforcement learning model
A: Where are you going? (1) B: I'm going to the restroom. (2) A: See you later. (3) B: See you later. (4) A: See you later. (5) B: See you later. (6)	A: Where are you going? (1) B: I'm going to the police station. (2) A: I'll come with you. (3) B: No, no, no, no, you're not going anywhere. (4) A: Why? (5) B: I need you to stay here. (6) A: I don't know what you are talking about. (7)
A: how old are you? (1) B: I'm 16. (2) A: 16? (3) B: I don't know what you are talking about. (4) A: You don't know what you are saying. (5) B: I don't know what you are talking about. (6) A: You don't know what you are saying. (7)	A: How old are you? (1) B: I'm 16. Why are you asking? (2) A I thought you were 12. (3) B: What made you think so? (4) A: I don't know what you are talking about. (5) B: You don't know what you are saying. (6)

Table 1: **Left Column:** Dialogue simulation between two agents using a 4-layer LSTM encoder-decoder trained on the OpenSubtitles dataset. The first turn (index 1) is input by the authors. Then the two agents take turns conversing, taking as input the other agent’s prior generated turn. The output is generated using the mutual information model (Li et al., 2015) in which an N-best list is first obtained using beam search based on $p(t|s)$ and reranked by linearly combining the backward probability $p(s|t)$, where t and s respectively denote targets and sources. **Right Column:** Dialogue simulated using the proposed reinforcement learning model. The new model has more forward-looking utterances (questions like “Why are you asking?” and offers like “I’ll come with you”) and lasts longer before it falls into conversational black holes.

These challenges suggest we need a conversation framework that has the ability to (1) integrate developer-defined rewards that better mimic the true goal of chatbot development and (2) model the long-term influence of a generated response in an ongoing dialogue.

To achieve these goals, we draw on the insights of reinforcement learning, which have been widely applied in MDP and POMDP dialogue systems (see Related Work section for details). We introduce a neural reinforcement learning (RL) generation method, which can optimize long-term rewards designed by system developers. Our model uses the encoder-decoder architecture as its backbone, and simulates conversation between two virtual agents to explore the space of possible actions while learning to maximize expected reward. We define simple heuristic approximations to rewards that characterize good conversations: good conversations are forward-looking (Allwood et al., 1992) or interactive (a turn suggests a following turn), informative, and coherent. The parameters of an encoder-decoder RNN define a policy over an infinite action space consisting of all possible

utterances. The agent learns a policy by optimizing the long-term developer-defined reward from ongoing dialogue simulations using policy gradient methods (Williams, 1992), rather than the MLE objective defined in standard SEQ2SEQ models.

Our model thus integrates the power of SEQ2SEQ systems to learn compositional semantic meanings of utterances with the strengths of reinforcement learning in optimizing for long-term goals across a conversation. Experimental results (sampled results at the right panel of Table 1) demonstrate that our approach fosters a more sustained dialogue and manages to produce more interactive responses than standard SEQ2SEQ models trained using the MLE objective.

2 Related Work

Efforts to build statistical dialog systems fall into two major categories.

The first treats dialogue generation as a source-to-target transduction problem and learns mapping rules between input messages and responses from a massive amount of training data. Ritter et al. (2011) frames the response generation problem as a statisti-

cal machine translation (SMT) problem. Sordoni et al. (2015) improved Ritter et al.’s system by rescore the outputs of a phrasal SMT-based conversation system with a neural model that incorporates prior context. Recent progress in SEQ2SEQ models inspire several efforts (Vinyals and Le, 2015) to build end-to-end conversational systems which first apply an encoder to map a message to a distributed vector representing its semantics and generate a response from the message vector. Serban et al. (2016) propose a hierarchical neural model that captures dependencies over an extended conversation history. Li et al. (2016a) propose mutual information between message and response as an alternative objective function in order to reduce the proportion of generic responses produced by SEQ2SEQ systems.

The other line of statistical research focuses on building task-oriented dialogue systems to solve domain-specific tasks. Efforts include statistical models such as Markov Decision Processes (MDPs) (Levin et al., 1997; Levin et al., 2000; Walker et al., 2003; Pieraccini et al., 2009), POMDP (Young et al., 2010; Young et al., 2013; Gašić et al., 2013a; Gašić et al., 2014) models, and models that statistically learn generation rules (Oh and Rudnicky, 2000; Ratnaparkhi, 2002; Banchs and Li, 2012; Nio et al., 2014). This dialogue literature thus widely applies reinforcement learning (Walker, 2000; Schatzmann et al., 2006; Gasic et al., 2013b; Singh et al., 1999; Singh et al., 2000; Singh et al., 2002) to train dialogue policies. But task-oriented RL dialogue systems often rely on carefully limited dialogue parameters, or hand-built templates with state, action and reward signals designed by humans for each new domain, making the paradigm difficult to extend to open-domain scenarios.

Also relevant is prior work on reinforcement learning for language understanding - including learning from delayed reward signals by playing text-based games (Narasimhan et al., 2015; He et al., 2016), executing instructions for Windows help (Branavan et al., 2011), or understanding dialogues that give navigation directions (Vogel and Jurafsky, 2010).

Our goal is to integrate the SEQ2SEQ and reinforcement learning paradigms, drawing on the advantages of both. We are thus particularly inspired by recent work that attempts to merge these paradigms, including Wen et al. (2016)— training an end-to-end

task-oriented dialogue system that links input representations to slot-value pairs in a database— or Su et al. (2016), who combine reinforcement learning with neural generation on tasks with real users, showing that reinforcement learning improves dialogue performance.

3 Reinforcement Learning for Open-Domain Dialogue

In this section, we describe in detail the components of the proposed RL model.

The learning system consists of two agents. We use p to denote sentences generated from the first agent and q to denote sentences from the second. The two agents take turns talking with each other. A dialogue can be represented as an alternating sequence of sentences generated by the two agents: $p_1, q_1, p_2, q_2, \dots, p_i, q_i$. We view the generated sentences as actions that are taken according to a policy defined by an encoder-decoder recurrent neural network language model.

The parameters of the network are optimized to maximize the expected future reward using policy search, as described in Section 4.3. Policy gradient methods are more appropriate for our scenario than Q-learning (Mnih et al., 2013), because we can initialize the encoder-decoder RNN using MLE parameters that already produce plausible responses, before changing the objective and tuning towards a policy that maximizes long-term reward. Q-learning, on the other hand, directly estimates the future expected reward of each action, which can differ from the MLE objective by orders of magnitude, thus making MLE parameters inappropriate for initialization. The components (states, actions, reward, etc.) of our sequential decision problem are summarized in the following sub-sections.

3.1 Action

An action a is the dialogue utterance to generate. The action space is infinite since arbitrary-length sequences can be generated.

3.2 State

A state is denoted by the previous two dialogue turns $[p_i, q_i]$. The dialogue history is further transformed to a vector representation by feeding the concatenation of p_i and q_i into an LSTM encoder model as

described in Li et al. (2016a).

3.3 Policy

A policy takes the form of an LSTM encoder-decoder (i.e., $p_{RL}(p_{i+1}|p_i, q_i)$) and is defined by its parameters. Note that we use a stochastic representation of the policy (a probability distribution over actions given states). A deterministic policy would result in a discontinuous objective that is difficult to optimize using gradient-based methods.

3.4 Reward

r denotes the reward obtained for each action. In this subsection, we discuss major factors that contribute to the success of a dialogue and describe how approximations to these factors can be operationalized in computable reward functions.

Ease of answering A turn generated by a machine should be easy to respond to. This aspect of a turn is related to its *forward-looking function*: the constraints a turn places on the next turn (Schegloff and Sacks, 1973; Allwood et al., 1992). We propose to measure the ease of answering a generated turn by using the negative log likelihood of responding to that utterance with a dull response. We manually constructed a list of dull responses \mathbb{S} consisting 8 turns such as “I don’t know what you are talking about”, “I have no idea”, etc., that we and others have found occur very frequently in SEQ2SEQ models of conversations. The reward function is given as follows:

$$r_1 = -\frac{1}{N_{\mathbb{S}}} \sum_{s \in \mathbb{S}} \frac{1}{N_s} \log p_{\text{seq2seq}}(s|a) \quad (1)$$

where $N_{\mathbb{S}}$ denotes the cardinality of \mathbb{S} and N_s denotes the number of tokens in the dull response s . Although of course there are more ways to generate dull responses than the list can cover, many of these responses are likely to fall into similar regions in the vector space computed by the model. A system less likely to generate utterances in the list is thus also less likely to generate other dull responses.

p_{seq2seq} represents the likelihood output by SEQ2SEQ models. It is worth noting that p_{seq2seq} is different from the stochastic policy function $p_{RL}(p_{i+1}|p_i, q_i)$, since the former is learned based on the MLE objective of the SEQ2SEQ model while the latter is the policy optimized for long-term future

reward in the RL setting. r_1 is further scaled by the length of target \mathbb{S} .

Information Flow We want each agent to contribute new information at each turn to keep the dialogue moving and avoid repetitive sequences. We therefore propose penalizing semantic similarity between consecutive turns from the same agent. Let h_{p_i} and $h_{p_{i+1}}$ denote representations obtained from the encoder for two consecutive turns p_i and p_{i+1} . The reward is given by the negative log of the cosine similarity between them:

$$r_2 = -\log \cos(h_{p_i}, h_{p_{i+1}}) = -\log \cos \frac{h_{p_i} \cdot h_{p_{i+1}}}{\|h_{p_i}\| \|h_{p_{i+1}}\|} \quad (2)$$

Semantic Coherence We also need to measure the adequacy of responses to avoid situations in which the generated replies are highly rewarded but are ungrammatical or not coherent. We therefore consider the mutual information between the action a and previous turns in the history to ensure the generated responses are coherent and appropriate:

$$r_3 = \frac{1}{N_a} \log p_{\text{seq2seq}}(a|q_i, p_i) + \frac{1}{N_{q_i}} \log p_{\text{seq2seq}}^{\text{backward}}(q_i|a) \quad (3)$$

$p_{\text{seq2seq}}(a|p_i, q_i)$ denotes the probability of generating response a given the previous dialogue utterances $[p_i, q_i]$. $p_{\text{seq2seq}}^{\text{backward}}(q_i|a)$ denotes the backward probability of generating the previous dialogue utterance q_i based on response a . $p_{\text{seq2seq}}^{\text{backward}}$ is trained in a similar way as standard SEQ2SEQ models with sources and targets swapped. Again, to control the influence of target length, both $\log p_{\text{seq2seq}}(a|q_i, p_i)$ and $\log p_{\text{seq2seq}}^{\text{backward}}(q_i|a)$ are scaled by the length of targets.

The final reward for action a is a weighted sum of the rewards discussed above:

$$r(a, [p_i, q_i]) = \lambda_1 r_1 + \lambda_2 r_2 + \lambda_3 r_3 \quad (4)$$

where $\lambda_1 + \lambda_2 + \lambda_3 = 1$. We set $\lambda_1 = 0.25$, $\lambda_2 = 0.25$ and $\lambda_3 = 0.5$. A reward is observed after the agent reaches the end of each sentence.

4 Simulation

The central idea behind our approach is to simulate the process of two virtual agents taking turns talking with each other, through which we can explore the

state-action space and learn a policy $p_{RL}(p_{i+1}|p_i, q_i)$ that leads to the optimal expected reward. We adopt an AlphaGo-style strategy (Silver et al., 2016) by initializing the RL system using a general response generation policy which is learned from a fully supervised setting.

4.1 Supervised Learning

For the first stage of training, we build on prior work of predicting a generated target sequence given dialogue history using the supervised SEQ2SEQ model (Vinyals and Le, 2015). Results from supervised models will be later used for initialization.

We trained a SEQ2SEQ model with attention (Bahdanau et al., 2015) on the OpenSubtitles dataset, which consists of roughly 80 million source-target pairs. We treated each turn in the dataset as a target and the concatenation of two previous sentences as source inputs.

4.2 Mutual Information

Samples from SEQ2SEQ models are often times dull and generic, e.g., “*i don’t know*” (Li et al., 2016a) We thus do not want to initialize the policy model using the pre-trained SEQ2SEQ models because this will lead to a lack of diversity in the RL models’ experiences. Li et al. (2016a) showed that modeling mutual information between sources and targets will significantly decrease the chance of generating dull responses and improve general response quality. We now show how we can obtain an encoder-decoder model which generates maximum mutual information responses.

As illustrated in Li et al. (2016a), direct decoding from Eq 3 is infeasible since the second term requires the target sentence to be completely generated. Inspired by recent work on sequence level learning (Ranzato et al., 2015), we treat the problem of generating maximum mutual information response as a reinforcement learning problem in which a reward of mutual information value is observed when the model arrives at the end of a sequence.

Similar to Ranzato et al. (2015), we use policy gradient methods (Sutton et al., 1999; Williams, 1992) for optimization. We initialize the policy model p_{RL} using a pre-trained $p_{SEQ2SEQ}(a|p_i, q_i)$ model. Given an input source $[p_i, q_i]$, we generate a candidate list $A = \{\hat{a}|\hat{a} \sim p_{RL}\}$. For each generated candi-

date \hat{a} , we will obtain the mutual information score $m(\hat{a}, [p_i, q_i])$ from the pre-trained $p_{SEQ2SEQ}(a|p_i, q_i)$ and $p_{SEQ2SEQ}^{backward}(q_i|a)$. This mutual information score will be used as a reward and back-propagated to the encoder-decoder model, tailoring it to generate sequences with higher rewards. We refer the readers to Zaremba and Sutskever (2015) and Williams (1992) for details. The expected reward for a sequence is given by:

$$J(\theta) = \mathbb{E}[m(\hat{a}, [p_i, q_i])] \quad (5)$$

The gradient is estimated using the likelihood ratio trick:

$$\nabla J(\theta) = m(\hat{a}, [p_i, q_i]) \nabla \log p_{RL}(\hat{a}|[p_i, q_i]) \quad (6)$$

We update the parameters in the encoder-decoder model using stochastic gradient descent. A curriculum learning strategy is adopted (Bengio et al., 2009) as in Ranzato et al. (2015) such that, for every sequence of length T we use the MLE loss for the first L tokens and the reinforcement algorithm for the remaining $T - L$ tokens. We gradually anneal the value of L to zero. A baseline strategy is employed to decrease the learning variance: an additional neural model takes as inputs the generated target and the initial source and outputs a baseline value, similar to the strategy adopted by Zaremba and Sutskever (2015). The final gradient is thus:

$$\nabla J(\theta) = \nabla \log p_{RL}(\hat{a}|[p_i, q_i]) [m(\hat{a}, [p_i, q_i]) - b] \quad (7)$$

4.3 Dialogue Simulation between Two Agents

We simulate conversations between the two virtual agents and have them take turns talking with each other. The simulation proceeds as follows: at the initial step, a message from the training set is fed to the first agent. The agent encodes the input message to a vector representation and starts decoding to generate a response output. Combining the immediate output from the first agent with the dialogue history, the second agent updates the state by encoding the dialogue history into a representation and uses the decoder RNN to generate responses, which are subsequently fed back to the first agent, and the process is repeated.

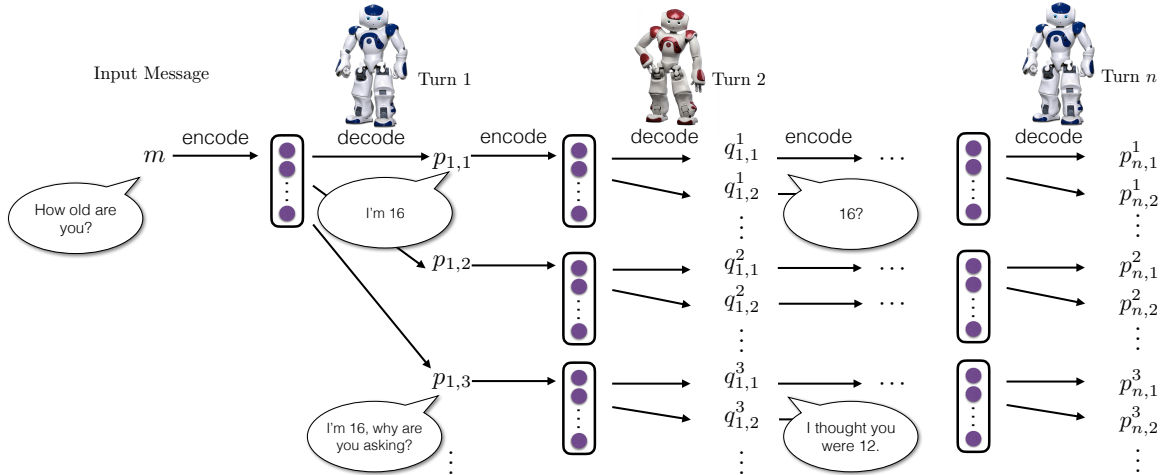


Figure 1: Dialogue simulation between the two agents.

Optimization We initialize the policy model p_{RL} with parameters from the mutual information model described in the previous subsection. We then use policy gradient methods to find parameters that lead to a larger expected reward. The objective to maximize is the expected future reward:

$$J_{RL}(\theta) = \mathbb{E}_{p_{RL}(a_{1:T})} \left[\sum_{i=1}^{i=T} R(a_i, [p_i, q_i]) \right] \quad (8)$$

where $R(a_i, [p_i, q_i])$ denotes the reward resulting from action a_i . We use the likelihood ratio trick (Williams, 1992; Glynn, 1990; Aleksandrov et al., 1968) for gradient updates:

$$\nabla J_{RL}(\theta) \approx \sum_i \nabla \log p(a_i | p_i, q_i) \sum_{i=1}^{i=T} R(a_i, [p_i, q_i]) \quad (9)$$

We refer readers to Williams (1992) and Glynn (1990) for more details.

4.4 Curriculum Learning

A curriculum Learning strategy is again employed in which we begin by simulating the dialogue for 2 turns, and gradually increase the number of simulated turns. We generate 5 turns at most, as the number of candidates to examine grows exponentially in the size of candidate list. Five candidate responses are generated at each step of the simulation.

5 Experimental Results

In this section, we describe experimental results along with qualitative analysis. We evaluate dialogue

generation systems using both human judgments and two automatic metrics: conversation length (number of turns in the entire session) and diversity.

5.1 Dataset

The dialogue simulation requires high-quality initial inputs fed to the agent. For example, an initial input of “why ?” is undesirable since it is unclear how the dialogue could proceed. We take a subset of 10 million messages from the OpenSubtitles dataset and extract 0.8 million sequences with the lowest likelihood of generating the response “*i don't know what you are taking about*” to ensure initial inputs are easy to respond to.

5.2 Automatic Evaluation

Evaluating dialogue systems is difficult. Metrics such as BLEU (Papineni et al., 2002) and perplexity have been widely used for dialogue quality evaluation (Li et al., 2016a; Vinyals and Le, 2015; Sordani et al., 2015), but it is widely debated how well these automatic metrics are correlated with true response quality (Liu et al., 2016; Galley et al., 2015). Since the goal of the proposed system is not to predict the highest probability response, but rather the long-term success of the dialogue, we do not employ BLEU or perplexity for evaluation².

²We found the RL model performs worse on BLEU score. On a random sample of 2,500 conversational pairs, single reference BLEU scores for RL models, mutual information models and vanilla SEQ2SEQ models are respectively 1.28, 1.44 and 1.17. BLEU is highly correlated with perplexity in generation tasks.

Model	# of simulated turns
SEQ2SEQ	2.68
mutual information	3.40
RL	4.48

Table 2: The average number of simulated turns from standard SEQ2SEQ models, mutual information model and the proposed RL model.

Length of the dialogue The first metric we propose is the length of the simulated dialogue. We say a dialogue ends when one of the agents starts generating dull responses such as “*i don’t know*”³ or two consecutive utterances from the same user are highly overlapping⁴.

The test set consists of 1,000 input messages. To reduce the risk of circular dialogues, we limit the number of simulated turns to be less than 8. Results are shown in Table 2. As can be seen, using mutual information leads to more sustained conversations between the two agents. The proposed RL model is first trained based on the mutual information objective and thus benefits from it in addition to the RL model. We observe that the RL model with dialogue simulation achieves the best evaluation score.

Diversity We report degree of diversity by calculating the number of distinct unigrams and bigrams in generated responses. The value is scaled by the total number of generated tokens to avoid favoring long sentences as described in Li et al. (2016a). The resulting metric is thus a type-token ratio for unigrams and bigrams.

For both the standard SEQ2SEQ model and the proposed RL model, we use beam search with a beam size 10 to generate a response to a given input message. For the mutual information model, we first generate n -best lists using $p_{\text{SEQ2SEQ}}(t|s)$ and then linearly re-rank them using $p_{\text{SEQ2SEQ}}(s|t)$. Results are presented in Table 4. We find that the proposed RL model generates more diverse outputs when com-

Since the RL model is trained based on future reward rather than MLE, it is not surprising that the RL based models achieve lower BLEU score.

³We use a simple rule matching method, with a list of 8 phrases that count as dull responses. Although this can lead to both false-positives and -negatives, it works pretty well in practice.

⁴Two utterances are considered to be repetitive if they share more than 80 percent of their words.

pared against both the vanilla SEQ2SEQ model and the mutual information model.

Model	Unigram	Bigram
SEQ2SEQ	0.0062	0.015
mutual information	0.011	0.031
RL	0.017	0.041

Table 4: Diversity scores (type-token ratios) for the standard SEQ2SEQ model, mutual information model and the proposed RL model.

Human Evaluation We explore three settings for human evaluation: the first setting is similar to what was described in Li et al. (2016a), where we employ crowdsourced judges to evaluate a random sample of 500 items. We present both an input message and the generated outputs to 3 judges and ask them to decide which of the two outputs is better (denoted as *single-turn general quality*). Ties are permitted. Identical strings are assigned the same score. We measure the improvement achieved by the RL model over the mutual information model by the mean difference in scores between the models.

For the second setting, judges are again presented with input messages and system outputs, but are asked to decide which of the two outputs is easier to respond to (denoted as *single-turn ease to answer*). Again we evaluate a random sample of 500 items, each being assigned to 3 judges.

For the third setting, judges are presented with simulated conversations between the two agents (denoted as *multi-turn general quality*). Each conversation consists of 5 turns. We evaluate 200 simulated conversations, each being assigned to 3 judges, who are asked to decide which of the simulated conversations is of higher quality.

Setting	RL-win	RL-lose	Tie
single-turn general quality	0.40	0.36	0.24
single-turn ease to answer	0.52	0.23	0.25
multi-turn general quality	0.72	0.12	0.16

Table 5: RL gains over the mutual information system based on pairwise human judgments.

Results for human evaluation are shown in Table 5. The proposed RL system does not introduce a significant boost in single-turn response quality (winning 40 percent of time and losing 36 percent of time).

Input Message	Mutual Information Model	RL
How old are you?	I'm 16.	I'm 16, why are you asking?
What's your full name?	I have no idea.	What's yours?
I don't want to go home tonight.	Really?	Why?
Do you have any feelings for me?	I don't know what you are talking about.	I don't want to hurt your feelings.
How much time do you have here?	Not long enough. Sorry, sir.	Ten seconds.
Shall we get started?	Of course!	Yes. We've got a lot of work to do here.
Do you play football?	No, i don't	Yes. I love football!
We'd have to talk to him.	I mean, he's a good guy	About what ?
How come you never say it?	Because I don't want to hurt you.	I don't think it's a good idea to say it.

Table 3: Sampled responses generated from the mutual information models and the proposed RL model.

This is in line with our expectations, as the RL model is not optimized to predict the next utterance, but rather to increase long-term reward. The RL system produces responses that are significantly easier to answer than does the mutual information system, as demonstrated by the *single-turn ease to answer* setting (winning 52 percent of time and losing 23 percent of time), and also significantly higher quality multi-turn dialogues, as demonstrated by the *multi-turn general quality* setting (winning 72 percent of time).

Qualitative Analysis and Discussion We show a random sample of generated responses in Table 3 and simulated conversations in Table 1 at the beginning of the paper. From Table 3, we can see that the RL based agent indeed generates more interactive responses than the other baselines. We also find that the RL model has a tendency to end a sentence with another question and hand the conversation over to the user. From Table 1, we observe that the RL model manages to produce more interactive and sustained conversations than the mutual information model.

During error analysis, we found that although we penalize repetitive utterances in consecutive turns, the dialogue sometimes enters a cycle with length greater than one, as shown in Table 6. This can be ascribed to the limited amount of conversational history we consider. Another issue observed is that the model sometimes starts a less relevant topic during the conversation. There is a tradeoff between relevance and less repetitiveness, as manifested in the reward function we define in Eq 4.

The fundamental problem, of course, is that the manually defined reward function can't possibly cover the crucial aspects that define an ideal conversation. While the heuristic rewards that we defined are amenable to automatic calculation, and do capture

A: What's your name ?
B: Daniel.
A: How old are you ?
B: Twelve. What's your name ?
A: Daniel.
B: How old are you ?
A: Twelve. What's your name ?
B: Daniel.
A: How old are you ?
B ...

Table 6: An simulated dialogue with a cycle longer than one.

some aspects of what makes a good conversation, ideally the system would instead receive real rewards from humans. Another problem with the current model is that we can only afford to explore a very small number of candidates and simulated turns since the number of cases to consider grow exponentially.

6 Conclusion

We introduce a reinforcement learning framework for neural response generation by simulating dialogues between two agents, integrating the strengths of neural SEQ2SEQ systems and reinforcement learning for dialogue. Like earlier neural SEQ2SEQ models, our framework captures the compositional models of the meaning of a dialogue turn and generates semantically appropriate responses. Like reinforcement learning dialogue systems, our framework is able to generate utterances that optimize future reward, successfully capturing global properties of a good conversation. Despite the fact that our model uses very simple, operationable heuristics for capturing these global properties, the framework generates more diverse, interactive responses that foster a more sustained conversation.

Acknowledgement

We would like to thank Chris Brockett, Bill Dolan and other members of the NLP group at Microsoft Research for insightful comments and suggestions. We also want to thank Kelvin Guu, Percy Liang, Chris Manning, Sida Wang, Ziang Xie and other members of the Stanford NLP groups for useful discussions. Jiwei Li is supported by the Facebook Fellowship, to which we gratefully acknowledge. This work is partially supported by the NSF via Awards IIS-1514268, IIS-1464128, and by the DARPA Communicating with Computers (CwC) program under ARO prime contract no. W911NF-15-1-0462. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of NSF, DARPA, or Facebook.

References

- V. M. Aleksandrov, V. I. Sysoyev, and V. V. Shemeneva. 1968. Stochastic optimization. *Engineering Cybernetics*, 5:11–16.
- Jens Allwood, Joakim Nivre, and Elisabeth Ahlsén. 1992. On the semantics and pragmatics of linguistic feedback. *Journal of Semantics*, 9:1–26.
- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2015. Neural machine translation by jointly learning to align and translate. In *Proc. of ICLR*.
- Rafael E Banchs and Haizhou Li. 2012. IRIS: a chat-oriented dialogue system based on the vector space model. In *Proceedings of the ACL 2012 System Demonstrations*, pages 37–42.
- Yoshua Bengio, Jérôme Louradour, Ronan Collobert, and Jason Weston. 2009. Curriculum learning. In *Proceedings of the 26th annual international conference on machine learning*, pages 41–48. ACM.
- SRK Branavan, David Silver, and Regina Barzilay. 2011. Learning to win by reading manuals in a monte-carlo framework. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pages 268–277.
- Michel Galley, Chris Brockett, Alessandro Sordani, Yangfeng Ji, Michael Auli, Chris Quirk, Margaret Mitchell, Jianfeng Gao, and Bill Dolan. 2015. deltaBLEU: A discriminative metric for generation tasks with intrinsically diverse targets. In *Proc. of ACL-IJCNLP*, pages 445–450, Beijing, China, July.
- Milica Gašić, Catherine Breslin, Matthew Henderson, Dongho Kim, Martin Szummer, Blaise Thomson, Pirros Tsiakoulis, and Steve Young. 2013a. Pomdp-based dialogue manager adaptation to extended domains. In *Proceedings of SIGDIAL*.
- Milica Gasic, Catherine Breslin, Mike Henderson, Dongkyu Kim, Martin Szummer, Blaise Thomson, Pirros Tsiakoulis, and Steve Young. 2013b. On-line policy optimisation of bayesian spoken dialogue systems via human interaction. In *Proceedings of ICASSP 2013*, pages 8367–8371. IEEE.
- Milica Gašić, Dongho Kim, Pirros Tsiakoulis, Catherine Breslin, Matthew Henderson, Martin Szummer, Blaise Thomson, and Steve Young. 2014. Incremental on-line adaptation of pomdp-based dialogue managers to extended domains. In *Proceedings on InterSpeech*.
- Peter W Glynn. 1990. Likelihood ratio gradient estimation for stochastic systems. *Communications of the ACM*, 33(10):75–84.
- Ji He, Jianshu Chen, Xiaodong He, Jianfeng Gao, Lihong Li, Li Deng, and Mari Ostendorf. 2016. Deep reinforcement learning with a natural language action space. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1621–1630, Berlin, Germany, August.
- Esther Levin, Roberto Pieraccini, and Wieland Eckert. 1997. Learning dialogue strategies within the markov decision process framework. In *Automatic Speech Recognition and Understanding, 1997. Proceedings., 1997 IEEE Workshop on*, pages 72–79. IEEE.
- Esther Levin, Roberto Pieraccini, and Wieland Eckert. 2000. A stochastic model of human-machine interaction for learning dialog strategies. *IEEE Transactions on Speech and Audio Processing*, 8(1):11–23.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016a. A diversity-promoting objective function for neural conversation models. In *Proc. of NAACL-HLT*.
- Jiwei Li, Michel Galley, Chris Brockett, Georgios Spithourakis, Jianfeng Gao, and Bill Dolan. 2016b. A persona-based neural conversation model. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 994–1003, Berlin, Germany, August.
- Chia-Wei Liu, Ryan Lowe, Iulian V Serban, Michael Noseworthy, Laurent Charlin, and Joelle Pineau. 2016. How not to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. *arXiv preprint arXiv:1603.08023*.
- Yi Luan, Yangfeng Ji, and Mari Ostendorf. 2016. LSTM based conversation models. *arXiv preprint arXiv:1603.09457*.
- Volodymyr Mnih, Koray Kavukcuoglu, David Silver, Alex Graves, Ioannis Antonoglou, Daan Wierstra, and Martin Riedmiller. 2013. Playing Atari with deep reinforcement learning. *NIPS Deep Learning Workshop*.

- Karthik Narasimhan, Tejas Kulkarni, and Regina Barzilay. 2015. Language understanding for text-based games using deep reinforcement learning. *arXiv preprint arXiv:1506.08941*.
- Lasguido Nio, Sakriani Sakti, Graham Neubig, Tomoki Toda, Mirna Adriani, and Satoshi Nakamura. 2014. Developing non-goal dialog system based on examples of drama television. In *Natural Interaction with Robots, Knowbots and Smartphones*, pages 355–361. Springer.
- Alice H Oh and Alexander I Rudnicky. 2000. Stochastic language generation for spoken dialogue systems. In *Proceedings of the 2000 ANLP/NAACL Workshop on Conversational systems-Volume 3*, pages 27–32.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318.
- Roberto Pieraccini, David Suendermann, Krishna Dayanidhi, and Jackson Liscombe. 2009. Are we there yet? Research in commercial spoken dialog systems. In *Text, Speech and Dialogue*, pages 3–13. Springer.
- Marc’Aurelio Ranzato, Sumit Chopra, Michael Auli, and Wojciech Zaremba. 2015. Sequence level training with recurrent neural networks. *arXiv preprint arXiv:1511.06732*.
- Adwait Ratnaparkhi. 2002. Trainable approaches to surface natural language generation and their application to conversational dialog systems. *Computer Speech & Language*, 16(3):435–455.
- Alan Ritter, Colin Cherry, and William B Dolan. 2011. Data-driven response generation in social media. In *Proceedings of EMNLP 2011*, pages 583–593.
- Jost Schatzmann, Karl Weilhammer, Matt Stuttle, and Steve Young. 2006. A survey of statistical user simulation techniques for reinforcement-learning of dialogue management strategies. *The knowledge engineering review*, 21(02):97–126.
- Emanuel A. Schegloff and Harvey Sacks. 1973. Opening up closings. *Semiotica*, 8(4):289–327.
- Iulian V Serban, Alessandro Sordoni, Yoshua Bengio, Aaron Courville, and Joelle Pineau. 2016. Building end-to-end dialogue systems using generative hierarchical neural network models. In *Proceedings of AAAI*, February.
- Lifeng Shang, Zhengdong Lu, and Hang Li. 2015. Neural responding machine for short-text conversation. In *Proceedings of ACL-IJCNLP*, pages 1577–1586.
- David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. 2016. Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489.
- Satinder P Singh, Michael J Kearns, Diane J Litman, and Marilyn A Walker. 1999. Reinforcement learning for spoken dialogue systems. In *Nips*, pages 956–962.
- Satinder Singh, Michael Kearns, Diane J Litman, Marilyn A Walker, et al. 2000. Empirical evaluation of a reinforcement learning spoken dialogue system. In *AAAI/IAAI*, pages 645–651.
- Satinder Singh, Diane Litman, Michael Kearns, and Marilyn Walker. 2002. Optimizing dialogue management with reinforcement learning: Experiments with the nj-fun system. *Journal of Artificial Intelligence Research*, pages 105–133.
- Alessandro Sordoni, Michel Galley, Michael Auli, Chris Brockett, Yangfeng Ji, Meg Mitchell, Jian-Yun Nie, Jianfeng Gao, and Bill Dolan. 2015. A neural network approach to context-sensitive generation of conversational responses. In *Proceedings of NAACL-HLT*.
- Pei-Hao Su, Milica Gasic, Nikola Mrksic, Lina Rojas-Barahona, Stefan Ultes, David Vandyke, Tsung-Hsien Wen, and Steve Young. 2016. Continuously learning neural dialogue management. *arxiv*.
- Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information processing systems*, pages 3104–3112.
- Richard S Sutton, David A McAllester, Satinder P Singh, Yishay Mansour, et al. 1999. Policy gradient methods for reinforcement learning with function approximation. In *NIPS*, volume 99, pages 1057–1063.
- Oriol Vinyals and Quoc Le. 2015. A neural conversational model. In *Proceedings of ICML Deep Learning Workshop*.
- Adam Vogel and Dan Jurafsky. 2010. Learning to follow navigational directions. In *Proceedings of ACL 2010*, pages 806–814.
- Marilyn A Walker, Rashmi Prasad, and Amanda Stent. 2003. A trainable generator for recommendations in multimodal dialog. In *Proceedings of INTERSPEECH 2003*.
- Marilyn A. Walker. 2000. An application of reinforcement learning to dialogue strategy selection in a spoken dialogue system for email. *Journal of Artificial Intelligence Research*, pages 387–416.
- Tsung-Hsien Wen, Milica Gasic, Nikola Mrkšić, Pei-Hao Su, David Vandyke, and Steve Young. 2015. Semantically conditioned LSTM-based natural language generation for spoken dialogue systems. In *Proceedings of EMNLP*, pages 1711–1721, Lisbon, Portugal.
- Tsung-Hsien Wen, Milica Gasic, Nikola Mrksic, Lina M Rojas-Barahona, Pei-Hao Su, Stefan Ultes, David Vandyke, and Steve Young. 2016. A network-based end-to-end trainable task-oriented dialogue system. *arXiv preprint arXiv:1604.04562*.

- Ronald J Williams. 1992. Simple statistical gradient-following algorithms for connectionist reinforcement learning. *Machine learning*, 8(3-4):229–256.
- Zhen Xu, Bingquan Liu, Baoxun Wang, Chengjie Sun, and Xiaolong Wang. 2016. Incorporating loose-structured knowledge into LSTM with recall gate for conversation modeling. *arXiv preprint arXiv:1605.05110*.
- Kaisheng Yao, Geoffrey Zweig, and Baolin Peng. 2015. Attention with intention for a neural network conversation model. In *NIPS workshop on Machine Learning for Spoken Language Understanding and Interaction*.
- Steve Young, Milica Gašić, Simon Keizer, François Mairesse, Jost Schatzmann, Blaise Thomson, and Kai Yu. 2010. The hidden information state model: A practical framework for pomdp-based spoken dialogue management. *Computer Speech & Language*, 24(2):150–174.
- Steve Young, Milica Gasic, Blaise Thomson, and Jason D Williams. 2013. Pomdp-based statistical spoken dialog systems: A review. *Proceedings of the IEEE*, 101(5):1160–1179.
- Wojciech Zaremba and Ilya Sutskever. 2015. Reinforcement learning neural Turing machines. *arXiv preprint arXiv:1505.00521*.