Hierarchical Text Generation and Planning for Strategic Dialogue

Denis Yarats * 1 Mike Lewis * 1

Abstract

End-to-end models for strategic dialogue are challenging to train, because linguistic and strategic aspects are entangled in latent state vectors. We introduce an approach to generating latent representations of dialogue moves, by inducing sentence representations to maximize the likelihood of subsequent sentences and actions. The effect is to decouple much of the semantics of the utterance from its linguistic realisation. We then use these latent sentence representations for hierarchical language generation, planning and reinforcement learning. Experiments show that using our message representations increases the reward achieved by the model, improves the effectiveness of long-term planning using rollouts, and allows self-play reinforcement learning to improve decision making without diverging from human language. Our hierarchical latent-variable model outperforms previous work both linguistically and strategically.

1. Introduction

Word-by-word approaches to text generation have been successful in many tasks. However, they have limitations in under-constrained generation settings, such as dialogue response or summarization, where models have significant freedom in the semantics of the text to generate. In such cases, models are prone to overly generic responses that may be valid but suboptimal in many situations (Li et al., 2015), or generating utterances that are semantically inconsistent. Further, such models do not cleanly distinguish between the semantics of language and its surface realization. Entangling form and meaning is problematic for reinforcement learning, where backpropagating caused by semantic decisions can adversely affect the linguistic quality of text, and for candidate generation for longterm planning, as linguistically diverse text may lack semantic diversity.

We focus on a strategic dialogue setting, where the text

generated by the model have consequences than can be easily measured. Substitutions of similar words (for example substituting a 'one' for a 'two') can have a large impact on the reward achieved by the dialogue agent. We use a hierarchical generation approach for a strategic dialogue agent, where the agent first samples a short-term plan in the form of a latent sentence representation. The agent then conditions on this plan during generation, allowing precise and consistent generation of text to achieve a short-term goal. We introduce a method for learning discrete latent representations of sentences based on their effect on the continuation of the dialogue.

Recent work has explored hierarchical generation of dialogue responses, where a latent variable z_t is inferred to maximize the likelihood of a message x_t , given previous messages $x_{0:t-1}$ (Serban et al., 2016b; Wen et al., 2017; Cao & Clark, 2017), which has the effect of clustering similar message strings. Our approach differs in that the latent variable z_t is optimized to maximize the likelihood of messages and actions of the continuation of the dialogue, but not the message x_t itself. Hence, z_t learns to represent x_t 's effect on the dialogue, but not the words of x_t . The distinction is important because messages with similar words can have very different semantics; and conversely the same meaning can be conveyed with different sentences. We show empirically that our method for learning sentence representations leads both to better perplexities and end task rewards, and qualitatively that our representations group sentences that are more semantically coherent but linguistically diverse.

We use our message representations to improve the strategic decision making of our dialogue agent. We improve the model's ability to plan ahead by creating a set of semantically diverse candidate messages by sampling distinct z_t , and then use rollouts to identify the an expected reward for each. We also apply reinforcement learning based on end-task reward. Previous work has found that RL can adversely effect the fluency of the language generated by the model (Lewis et al., 2017). We instead show that simply fine-tuning the parameters for choosing z_t allows the model to substantially improve its rewards while maintaining human-like language.

^{*}Equal contribution ¹Facebook AI Research, Menlo Park, CA. Correspondence to: Mike Lewis <mikelewis@fb.com>.

2. Background

2.1. Natural Language Negotiations

We focus on the natural language negotiation task introduced by Lewis et al. (2017).

In this setting, agents A and B are initially given a space \mathcal{A} of possible agreements, and value functions v^A and v^B , which specify a non-negative reward for each agreement $a \in \mathcal{A}$. Agents cannot observe each other's value functions and can only infer it through a dialogue.

The agents sequentially exchange turns of natural language x_t , consisting of words $x_t^{0:n_t}$, until one agent enters a special turn that ends the dialogue. Then, both agents independently enter agreements $a^A, a^B \in \mathcal{A}$ respectively. If the agreements are compatible, both agents receive a reward based on their action and the value function. If the actions are incompatible, neither agent receives any reward.

Lewis et al. (2017) collected a corpus of human dialogues on a multi-issue bargaining task, where the agents must divide a collection of items of 3 different types (*books*, *hats* and *balls*) between them. Actions correspond to choosing a particular subset of the items, and agents choose compatible actions if each item is assigned to exactly one agent.

Training dialogues from an agent's perspective consist of agreement space A, value function v, messages $x_{0:T}$ and agreement a.

2.2. Challenges in Text Generation for Strategic Dialogue

We identify a number of challenges for end-to-end text generation for strategic dialogue. These problems have been identified in other text generation settings, but strategic dialogue makes an interesting test case, where decisions have measurable consequences.

- Lack of semantic diversity Multiple samples from a model are often paraphrases of the same intent. This lack of a diversity is a problem if samples are later re-ranked by a long-term planning model.
- Lack of linguistic diversity Neural language models are prone to capturing the head of the distribution, and provide much less varied language than people.
- Lack of internal coherence Messages generated by the model often lack self consistency—for example, *I'll take one hat, and give you all the hats*.
- Lack of contextual coherence Utterances may also lack coherence given the dialogue context so far. For example, Lewis et al. (2017) identify cases where a model starts a message by indicating agreement, but then proposes a counter offer.

Entanglement of linguistic and strategic parameters End-to-end approaches do not cleanly distinguish between what to say and how to say it. This is problematic as reinforcement learning aiming to improve decision making may adversely affect the quality of the generated language.

We argue that these limitations partly stem from the wordby-word sampling approach to generation, with no explicit plan in advance of generation for what the meaning of the sentence is to be. In section 8, we show our hierarchical approach to generation helps with these problems.

3. Baseline Model

As a baseline, we train a hierarchical encoder-decoder model to maximize the likelihood of training dialogue messages and actions, similarly to Serban et al. (2016a).

The model contains a sentence encoder GRU_x^e that embeds individual messages x_t as e_t ; a sentence level GRU_e^s that reads sentence embeddings $e_{0:t}$ to produce dialogue state s_t ; and a decoder GRU_s^x that produces message x_{t+1} , using s_t . The encoder and decoder share a word embedding matrix E.

$$e_t = \operatorname{GRU}_x^e(Ex_t^{0:n_t})$$

$$s_t = \operatorname{GRU}_e^s(e_{0:t})$$

$$p_x(x_t^i|x_t^{0:i-1}, x_{0:t-1}) \propto \exp(E^{\top} \operatorname{GRU}_s^x(s_{t-1}; x_t^{0:i-1}))$$

$$p_x(x_t|x_{0:t-1}) = \prod_{i=1}^{n_t} p_x(x_t^i|x_t^{0:i-1}, x_{0:t-1})$$

We optimize the following loss, over the training set $\mathbf{x}_{0:T}$.

$$\mathcal{L} = \sum_{\mathbf{x}} \sum_{t} \log p_x(x_t | x_{0:t-1})$$

We also train an action classifier $\pi_a(a|x_{0:T})$ that predicts the final action chosen at the end of the dialogue. This model first encodes the dialogue using a hierarchical GRU, then uses attention to select the encoding of the most relevant sentences, and finally uses a feed-forward network to output each component of the agreement.

4. Learning Latent Message Representations

The key part of our model is a method for encoding messages x_t as discrete latent variables z_t . The goal of this model is to learn message representations that reflect the message's effect on the dialogue, but abstract over semantically equivalent paraphrases. We show that such representations are helpful for planning and reinforcement learning.

Our representation learning model (Figure 1a) has a similar structure to that of $\S 3$, except that message embedding e_t is

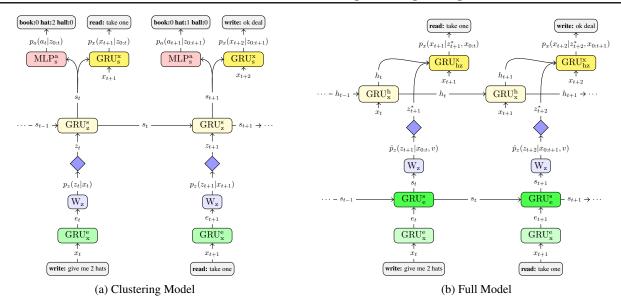


Figure 1: We pre-train a model to learn a discrete encoder for sentences, which bottlenecks the message x_t through a discrete representation z_t (Figure 1a; §4). This architecture forces z_t to capture the most relevant aspects of x_t for predicting future messages and actions. We then train our full model (Figure 1b): p_x is trained to translate representations z_t^* into messages x_t (§5.1), and \hat{p}_z is trained to predict a distribution over z_t given the dialogue history (§5.2).

used as input to a softmax with parameters W_z over latent states z_t . We use expectation maximization to learn how to assign messages to clusters to maximize the likelihood of future messages and actions.

After each message x_t , GRU_z^s is updated with representation z_t to give hidden state s_t . From s_t , we train the model to predict the next message x_{t+1} and an action a_t . In the training dialogues, there is only an action after the final turn x_T ; for other turns x_t , we use a soft proxy action by regressing to the distribution over actions predicted by $a_t = \pi_a(a|x_{0:t})$. Therefore, a_t is a distribution over what deal would be agreed if the dialogue stopped after message x_t . When predicting x_{t+1} and a_t , the model only has access to latent variables $z_{0:t}$, so z_t must contain useful information about the meaning of x_t .

We optimize the following loss:

$$\mathcal{L} = \sum_{\mathbf{x}} \sum_{t} \log p_x(x_{t+1}|z_{0:t}) + \\ D_{\text{KL}}(\pi_a(a|x_{0:t})||p_a(a_t|z_{0:t}))$$

encodings are passed through a discrete bottleneck:

$$e_{t} = \operatorname{GRU}_{x}^{e}(Ex_{t}^{0:n_{t}})$$

$$p_{z}(z_{t}|x_{t}) \propto \exp(W_{z}e_{t})$$

$$s_{t} = \operatorname{GRU}_{z}^{s}(z_{0:t})$$

$$p_{x}(x_{t+1}^{i}|x_{t+1}^{0:i-1}, z_{0:t}) \propto \exp(E^{\top}\operatorname{GRU}_{s}^{x}(s_{t}; x_{t+1}^{0:i-1}))$$

$$p_{x}(x_{t+1}|z_{0:t}) = \prod_{i=1}^{n_{t+1}} p_{x}(x_{t+1}^{i}|x_{t+1}^{0:i-1}, z_{0:t})$$

$$p_{a}(a_{t}|z_{0:t}) \propto \exp(\operatorname{MLP}_{s}^{a}(s_{t}))$$

We optimize latent variables z using minibatch Viterbi Expectation Maximisation (Dempster et al., 1977). For each minibatch, for each timestep t, we compute:

$$z_t^* = \operatorname*{argmax}_{z} p(x_{t+1}, a_t | z, z_{0:t-1}) p_z(z | x_t)$$

Computing the argmax requires a separate forward pass for each z.

We then advance to the next timestep using z_t^* to update GRU_z^s . Finally, we perform a gradient update maximizing:

$$\sum_{t} \log p(x_{t+1}, a_t | z_t^*, z_{0:t-1}) p_z(z_t^* | x_t)$$

The model employs a hierarchical GRU, in which message

At convergence, we extract message representations z_t^* .

5. Hierarchical Text Generation

We then train a new hierarchical dialogue model (Figure 1b), which uses pre-trained representations z_t^* to predict messages x_t .

First, we train a recurrent neural network to predict $p_x(x_{t+1}|z_{t+1}^*, x_{0:t})$. p_x learns how to translate the latent variables into fluent text in context. Then, we optimize a model $\hat{p}_z(z_{t+1}|x_{0:t})$ to maximize the marginal likelihood of training sentences.

5.1. Conditional Language Model

 p_x learns to translate pretrained representation z_t^* into a message x_t . p_x is implemented as a hierarchical RNN:

$$p_x(x_{t+1}^i|z_{t+1}^*, x_{0:t}, x_{t+1}^{0:i-1})$$

$$\propto \exp(E^{\top} \text{GRU}_{hz}^x(z_{t+1}^*, h_t; x_{t+1}^{0:i-1}))$$

where

$$h_t = \mathrm{GRU}_x^h(Ex_{0:t})$$

We optimize the following loss:

$$\mathcal{L}_x = \sum_{\mathbf{x}} \sum_{t} \log p_x(x_{t+1}|z_{t+1}^*, x_{0:t})$$

given

$$p_x(x_{t+1}|z_{t+1}^*, x_{0:t}) = \prod_{i=1}^{n_{t+1}} p_x(x_{t+1}^i|z_{t+1}^*, x_{0:t}, x_{t+1}^{0:i-1})$$

Note that, unlike the baseline model, text generation does not condition explicitly on the agent's value function v—meaning all knowledge of the goals and available actions is bottlenecked through the dialogue state. This restriction forces the text generation to depend strongly on the semantics expressed by z_t .

5.2. Latent Variable Prediction Model

At test time, z_t^* is not available, as it contains information about the future dialogue. Instead, we train a model \hat{p}_z to predict z_t conditioned on the current dialogue context:

$$e_t = \operatorname{GRU}_x^e(Ex_t^{0:n_t})$$
$$s_t = \operatorname{GRU}_e^s(e_{0:t})$$
$$\hat{p}_z(z_{t+1}|x_{0:t}, v) \propto \exp(W_z s_t)$$

We optimize \hat{p}_z to maximize the marginal likelihood of training messages, without updating p_x . The model learns to reconstruct the distribution over z_t that best explains message x_t .

$$\mathcal{L}_z = \sum_{\mathbf{x}} \sum_{t} \log \sum_{z} p_x(x_t|z, x_{0:t-1}) \hat{p}_z(z|x_{0:t-1}, v)$$

5.3. Decoding

To generate an utterance x_t , the model first samples a predicted plan z_t from \hat{p}_z :

$$z_t \sim \hat{p}_z(z|x_{0:t-1})$$

The model then sequentially generates tokens x_t^i based on plan z_t and context $x_{0:t}$:

$$x_t^{i+1} \sim p_x(x|z_t, x_{0:t-1}, x_t^{0:i})$$

6. Hierarchial Reinforcement Learning

Lewis et al. (2017) experiment with end-to-end reinforcement learning to fine-tune pre-trained supervised models. The model engages in a dialogue with another model, achieving reward R. This reward is then backpropagated using policy gradients.

One challenge is that because model parameters govern both strategic and linguistic aspects of generation, backpropagating errors can adversely affect the quality of the generated language.

To avoid divergence from human language, we experiment with fixing all model parameters, except for the parameters of \hat{p}_z . This allows reinforcement learning to improve decisions about what to say, without affecting language generation parameters. A similar approach was taken in a different dialogue setting by Wen et al. (2017).

7. Hierarchical Planning

Lewis et al. (2017) propose planning for strategic dialogue using rollouts, where first a set of K unique candidate messages $\{x_t^{(1)}, x_t^{(2)}, \dots, x_t^{(K)}\}$ are sampled from $p_x(x|x_{0:t-1})$. Then, multiple rollouts of the future dialogue are sampled from the model, and the outcomes a are scored according to the value function v, to estimate the expected reward $R(x_t)$:

$$R(x_t) = \mathbb{E}_{x_{t+1:T} \sim p_x, a \sim \pi_a}[r(a, v)\pi_a(a|x_{0:T})]$$
 (1)

The expectation is approximated by taking N samples. Then, the candidate x^{\ast} with the highest expected score can be returned.

$$x^* = \operatorname*{argmax}_{x} R(x) \tag{2}$$

One challenge is that even though the candidates $\{x_t^{(1)}, x_t^{(2)}, \dots, x_t^{(K)}\}$ can be constrained to be different strings, it is difficult to enforce semantic diversity. For example, if all the candidates are paraphrases of the same intent, then the choice makes little difference to the outcome of the dialogue.

In order to improve the diversity of candidate generation, we take a hierarchical approach of first sampling K unique latent intents $\{z_t^{(1)}, z_t^{(2)}, \dots, z_t^{(K)}\}$ from $\hat{p}_z(z|x_{0:t-1})$. Then, for each $z_t^{(i)}$, we choose a candidate turn conditioned on that state:

$$x_t^{(i)} = \underset{x}{\operatorname{argmax}} p_x(x|z_t^{(i)}, x_{0:t-1})$$

We then estimate the reward of the candidate message using Equation 1, and finally choose a message as in Equation 2.

8. Experiments

8.1. Training Details

The following hyper-parameters are used for training. We embed input words into 256-dimensional space. The hidden states of each GRU module in the network are of size 256. For each unique agreement space A we learn 50 discrete latent message representations. During training, we optimize the parameters using RMSProp (Tieleman & Hinton, 2012) with an initial learning rate of 0.0005 and momentum $\mu = 0.1$, we also employ clipping of gradients whose L^2 norm exceeds 1 to ease the training. We train the models for 15 epochs with mini-batch size of 16. We then pick the best snapshot according to validation perplexity and anneal the learning rate by a factor of 5 each epoch. For reinforcement learning, we use a smaller learning rate of 0.0001. The discount factor γ is set to 0.95. For predicting the final agreement given a dialogue, we use classifier π_a with all models.

A development set was used to select the hyper-parameters above. All the models were implemented using PyTorch.

8.2. Baselines

We compare the following models:

- RNN A simple word-by-word approach to generation, similar to Lewis et al. (2017).
- HIERARCHICAL Version of our model in which the two levels of RNNs are connected directly, without the discrete bottleneck z_t . This is effectively similar to Serban et al. (2016a).
- BASELINE CLUSTERS Our model (Figure 1b) without pretraining the sentence encoder. A latent representation z_t of message x_t is inferred to maximize the likelihood of $p(x_t|z_t,x_{0:t-1})p(z_t|x_{0:t-1})$. This model is closely related to the LATENT INTENTS DIALOGUE MODEL (Wen et al., 2017).
- FULL MODEL Our full model, where we first pre-train sentence representations z_t^* to maximize the log likelihood $\sum_t \log p(x_{t+1}, a_t|z_t^*, z_{0:t-1}) p_z(z_t^*|x_t)$, and

Model	Validation Perplexity	Test Perplexity
RNN	5.62	5.47
HIERARCHICAL	5.37	5.21
BASELINE CLUSTERS	5.61	5.46
FULL MODEL	5.37	5.24

Table 1: Perplexity results, showing the likelihood of human dialogues using different models. Our model with discrete message representations is able to achieve state-of-the-art performance, showing that the representations effectively capture relevant aspects of messages for predicting the future dialogue.

then we train models to predict $p_x(x_t|z_t^*)$ and $\hat{p}_z(z_t|x_{0:t-1})$.

To focus the evaluation on the linguistic and strategic aspects of the dialogue, all systems use the same model for predicting the final agreement represented by the dialogue, which is implemented as a bidirectional GRU with attention over the words of the dialogue.

8.3. Likelihood Models

First, we experiment with models using no reinforcement learning or rollouts.

8.3.1. Perplexity

Models were developed to maximize the likelihood of human dialogues, which is an indicator of how human-like the language is. Results are shown in Table 1.

The use of a hierarchical RNN model improves performance over a strong baseline from previous work.

Perhaps surprisingly, our hierarchical latent-variable model is also able to achieve state-of-the-art performance. This shows our model's discrete encodings of messages are as informative for predicting the future dialogue as the more-expressive embeddings used by the hierarchical baseline.

8.3.2. Coherence of Clusters

Table 4 shows random samples of messages generated by different clusters from our predicted state model, and the BASELINE CLUSTERS model.

Qualitatively, the states from our model show a higher degree of semantic coherence, and higher linguistic variability. Compared to the BASELINE CLUSTERS, our approach tends to generate more dissimilar surface strings, but with more similar semantics. Our clusters appear to capture *meaning* rather than *form*.

Model	Score vs. RNN	Score vs. HIERARCHICAL
RNN	5.33	5.17
HIERARCHICAL	5.37	5.08
BASELINE CLUSTERS	4.68	4.66
FULL MODEL	6.75	6.57

Table 2: Comparison of different models based on their end-task reward. Our clusters substantially improve reward, indicating that they make it easier for supervised learning to model strategic decision making.

Rollout Type	Score vs. No Rollouts	Score vs. BASELINE ROLLOUTS
No Rollouts	5.08	4.91
BASELINE	7.81	6.57
DIVERSE	8.41	7.36

Table 3: Comparison of different rollout strategies for the FULL MODEL. DIVERSE rollouts use distinct latent variables to create more semantic diversity in rollout candidates, significantly improving performance.

8.3.3. END TASK PERFORMANCE

Finally, we measure the performance of the different models on their end-task reward over 1000 negotiations in self-play. Results are shown in Table 2. We find that the use of our latent representations leads to a large improvement in the reward, indicating that our representations make it easier for the supervised model to learn the latent decision making process in the human dialogues it was trained on.

8.4. Hierarchical Planning

Next, we evaluate the effectiveness of different rollout strategies using our model:

- BASELINE ROLLOUTS following Lewis et al. (2017), where first K candidate sentences are sampled from the model, and then tokens are sampled iteratively from p_x until reaching the end of the dialogue.
- DIVERSE ROLLOUTS where we first choose the mostly likely K unique z_t from \hat{p}_z . By choosing unique z_t we aim to increase the semantic diversity of the candidates.

We evaluate compared to the baseline model and wordlevel rollouts and record the average score.

Results are shown in Table 3, and that the DIVERSE ROLL-OUTS that use our message representations lead to a large improvement over previous approaches.

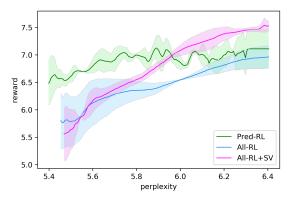


Figure 2: Plotting reward against language quality (lower perplexity is better) during reinforcement learning training, in dialogues with the HIERARCHICAL model. Our method (green) achieves higher rewards while maintaining human-like language (left of graph).

8.5. Finetuning with Reinforcement Learning

A challenge in using reinforcement learning for end-toend text generation models is that optimising for reward can adversely affect language generation. In selfplay, the model can learn to achieve a high reward by finding uninterpretable sequences of tokens that the baseline model was not exposed to at training time. We compare several reinforcement learning approaches:

- ALL-RL Reinforcement learning after pre-training with supervised learning.
- ALL-RL+SV Interleaved reinforcement learning and supervised learning updates, weighting supervised updates with an additional hyperparameter α, similarly to Lewis et al. (2017).
- PRED-RL Reinforcement learning only to fine-tune the intent prediction model \hat{p}_z , with all other model parameters fixed.

We measure both the average reward of the model (a measure of its ability to achieve its goals) and the perplexity of the model on human dialogues (a measure of how human-like the language is). After hyper-parameter grid search, we plot the reward of the best model whose perplexity is at most a.

Results are shown in Figure 2. Using RL on all parameters allows high rewards at the price of poor quality language. Only fine-tuning \hat{p}_z allows the model to improve its strategic decision making, while still assigning high likelihood to human language.

Cluster	BASELINE CLUSTERS	FULL MODEL
1	i can give you the books but, i would need the hat and the balls	i would like the hat and 1 book
	i can do that . i need both balls and one book	i can't give up the hat, but i can offer you the book and 2 balls
2	i need both books and the hat	i want the hat
	how about you get the hat and 1 ball	i need the hat . you can have all the books and the balls
3	i can not make that deal . i need the hat and one book	i can give you the hat and 1 ball
	i can give you the hat and 1 ball	i would like the books and a ball
4	i need two books and the hat	i need the books and the hat
	i need the hat, you can have the rest	i can give you the balls but i need the hat and books
5	i can give you the hat if i can have the rest	could i have the books and a ball ?
	i want one of each	i would like the books and one ball

Table 4: Sample messages that are probable under different clusters for specified context, in comparison to a previous approach to learning message representations. An agreement needs to be done over a set of **2 books**, **1 hat**, and **2 balls**. The clusters produced by our method are much more semantically coherent than the baseline, and correspond closely to different ways of proposing the same deal.

9. Analysis

Results in section 8 show quantitatively that our hierarchical model improves the likelihood of human generated language and the average score achieved by the agent. Here, we investigate specific issues that the model improved on, and identify remaining challenges. We analyzed 1000 dialogues between our FULL MODEL and the HIERARCHICAL baseline. These models achieve similar perplexity on human dialogues (Table 1).

9.1. Linguistic Diversity

First, we investigate the amount of variation in the language used by the agents.

RNN language models are known to prefer overly generic messages. In our task, this often manifests itself as short messages expressing agreement such as *deal* or *ok*. We measure the frequency of simple variations on these messages, and find that the HIERARCHICAL model uses very generic messages far more often than FULL MODEL (815 times vs. 245).

The messages sent by FULL MODEL are also longer on average (8.9 words vs. 6.7, ignoring the special final message that is a single token ending the dialogue), giving further evidence of greater complexity.

We also find that the FULL MODEL is substantially more creative in generating new messages beyond those seen in its training data. In total, FULL MODEL sends 875 unique message strings, of which 525 (60%) do not appear in the training data. In contrast, HIERARCHICAL sends fewer unique message strings (751), and just 18% of these are not copied from the training data.

9.2. Self-consistency of Messages

Models can sometimes output inconsistent messages, such as *I really need the hat. I can give you the hat and one ball.* We searched for messages that mentioned the same item type multiple times, and then manually evaluated whether it was consistent. We found that the FULL MODEL model was more prone to this error than HIERARCHICAL (23 times vs. 11), though this fact may be a consequence of its greater creativity, and the problem only occurred in roughly 1% of messages.

9.3. Consistency with Input

We also investigate whether messages are consistent with the context—for example, models may emit messages such as *I'd like the hat and books; you keep the 3 balls* when there are not 3 balls available. We use simple pattern matching for several such errors, and found that the FULL MODEL performed slightly better (15 errors vs. 19).

9.4. Consistency with Dialogue Context

Lewis et al. (2017) describe cases where an agent indicates it is simply re-stating an agreement, when it is actually proposing a new deal (e.g. *you get 2 hats / Okay deal, so I get 3 hats*). Interestingly, we found this behaviour only happened with the models using rollouts. While this tactic is effective against our models, it would be frustrating for humans, and future work should address this issue.

9.5. Repetitiveness

Previous work noted that reinforcement learning models were prone to an extortion tactic of simply repeating the same demand until acceptance. We measured how often agents repeated the same message in a dialogue, compar-

Input		
DIVERSE ROLLOUTS	1xbook value=9 1xhat value=1 4xball value=0	
BASELINE ROLLOUTS	1xbook value=0 1xhat value=6 4xball value=1	
DIVERSE ROLLOUTS	I will take the book and hat and you can have the balls.	
BASELINE ROLLOUTS	I need the hat and two balls	
DIVERSE ROLLOUTS	The balls are worthless, I need the hat and the book.	
BASELINE ROLLOUTS	I need the hat or no deal	
DIVERSE ROLLOUTS	Then no deal.	
BASELINE ROLLOUTS	What about the balls?	
DIVERSE ROLLOUTS	You can have the hat but I need the book.	
BASELINE ROLLOUTS	How about I get the balls and 1 hat?	
DIVERSE ROLLOUTS	Ok	
Output		Reward
DIVERSE ROLLOUTS	1x book	9/10
BASELINE ROLLOUTS	1x hat 4x ball	10/10

Figure 3: Example dialogue between two models using different type of rollouts. The DIVERSE ROLLOUTS model makes several attempts to win the *hat*, before finally compromizing.

ing the ALL-RL+SV model based on previous work, with our PRED-RL model. We found that our model was substantially less repetitive: only 1% of dialogues contained a repetition of the same message, compared to 12% for the baseline.

10. Related Work

Traditional goal-orientated dialogue models have first generated symbolic *intents*, capturing the meaning of the message, and then generated text to match the intent (e.g. Williams & Young (2007), Keizer et al. (2017)). Our approach can be seen as a latent model for generating intents.

Our model is most closely related to other recent latent variable hierarchical dialogue models from Serban et al. (2016b), Wen et al. (2017) and Cao & Clark (2017). A key difference is that both these approaches optimize latent representations z to maximize the likelihood of generating the next message—whereas our model pretrain's z to maximize the likelihood of the continuation of the dialogue, to better capture the semantics of the message rather than its surface form. We have shown that our approach leads to higher performance on a strategic dialogue task.

Other work has explored generating sentence embeddings for open domain text—for example, based on maximizing the likelihood of surrounding sentences (Kiros et al., 2015), supervised entailment data (Conneau et al., 2017), and auto-encoders (Bowman et al., 2015).

11. Conclusion

We have introduced a novel approach to creating sentence representations, within the context of an end-to-end strategic dialogue system, and have shown that our hierarchical approach improves text generation and planning. We identified a number of challenges faced by previous work, and show empirically that our model improves on these aspects. Future work should apply our model to other dialogue settings, such as cooperative strategic dialogue games (He et al., 2017), or multi-sentence generation tasks, such as long document language modelling (Merity et al., 2016).

References

Bowman, Samuel R., Vilnis, Luke, Vinyals, Oriol, Dai, Andrew M., Józefowicz, Rafal, and Bengio, Samy. Generating sentences from a continuous space. *CoRR*, abs/1511.06349, 2015.

Cao, Kris and Clark, Stephen. Latent variable dialogue models and their diversity. *CoRR*, abs/1702.05962, 2017. URL http://arxiv.org/abs/1702.05962.

Conneau, Alexis, Kiela, Douwe, Schwenk, Holger, Barrault, Loic, and Bordes, Antoine. Supervised learning of universal sentence representations from natural language inference data. *arXiv* preprint arXiv:1705.02364, 2017.

Dempster, Arthur P, Laird, Nan M, and Rubin, Donald B. Maximum Likelihood from Incomplete Data via the EM Algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, pp. 1–38, 1977.

He, H., Balakrishnan, A., Eric, M., and Liang, P. Learning Symmetric Collaborative Dialogue Agents with Dynamic Knowledge Graph Embeddings. In *Association for Computational Linguistics (ACL)*, 2017.

Keizer, Simon, Guhe, Markus, Cuayhuitl, Heriberto, Efstathiou, Ioannis, Engelbrecht, Klaus-Peter, Dobre, Mihai, Lascarides, Alexandra, and Lemon, Oliver. Evaluating Persuasion Strategies and Deep Reinforcement Learning methods for Negotiation Dialogue agents. In Proceedings of the European Chapter of the Association for Computational Linguistics (EACL 2017), 2 2017.

Kiros, Ryan, Zhu, Yukun, Salakhutdinov, Ruslan, Zemel, Richard S., Torralba, Antonio, Urtasun, Raquel, and Fidler, Sanja. Skip-thought vectors. In *Proceedings of the 28th International Conference on Neural Information Processing Systems - Volume 2*, NIPS'15, pp. 3294–3302, Cambridge, MA, USA, 2015. MIT Press.

Lewis, Mike, Yarats, Denis, Dauphin, Yann N, Parikh, Devi, and Batra, Dhruv. Deal or no deal? end-to-end learning for negotiation dialogues. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language*, pp. 2433–2443. Association for Computational Linguistics, September 2017.

Li, Jiwei, Galley, Michel, Brockett, Chris, Gao, Jianfeng, and Dolan, Bill. A Diversity-promoting Objective Func-

- tion for Neural Conversation Models. arXiv preprint arXiv:1510.03055, 2015.
- Merity, Stephen, Xiong, Caiming, Bradbury, James, and Socher, Richard. Pointer sentinel mixture models. *arXiv* preprint arXiv:1609.07843, 2016.
- Serban, Iulian Vlad, Sordoni, Alessandro, Bengio, Yoshua, Courville, Aaron C, and Pineau, Joelle. Building end-to-end dialogue systems using generative hierarchical neural network models. In *AAAI*, pp. 3776–3784, 2016a.
- Serban, Iulian Vlad, Sordoni, Alessandro, Lowe, Ryan, Charlin, Laurent, Pineau, Joelle, Courville, Aaron C., and Bengio, Yoshua. A hierarchical latent variable encoder-decoder model for generating dialogues. *CoRR*, abs/1605.06069, 2016b. URL http://arxiv.org/abs/1605.06069.
- Tieleman, Tijmen and Hinton, Geoffrey. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural networks for machine learning*, 4(2):26–31, 2012.
- Wen, Tsung-Hsien, Miao, Yishu, Blunsom, Phil, and Young, Steve J. Latent intention dialogue models. *CoRR*, abs/1705.10229, 2017. URL http://arxiv.org/abs/1705.10229.
- Williams, Jason D and Young, Steve. Partially Observable Markov Decision Processes for Spoken Dialog Systems. *Computer Speech & Language*, 21(2):393–422, 2007.