

Hierarchical Recurrent Attention Network for Response Generation

Chen Xing^{12*}, Wei Wu³, Yu Wu⁴, Ming Zhou³, Yalou Huang¹², Wei-Ying Ma³

¹College of Computer and Control Engineering, Nankai University, Tianjin, China

²College of Software, Nankai University, Tianjin, China

³ Microsoft Research, Beijing, China

⁴State Key Lab of Software Development Environment, Beihang University, Beijing, China

{v-chxing,wuwei,v-wuyu,mingzhou,wyma}@microsoft.com yluhuan@nankai.edu.cn

Abstract

We study multi-turn response generation in chatbots where a response is generated according to a conversation context. Existing work has modeled the hierarchy of the context, but does not pay enough attention to the fact that words and utterances in the context are differentially important. As a result, they may lose important information in context and generate irrelevant responses. We propose a hierarchical recurrent attention network (HRAN) to model both aspects in a unified framework. In HRAN, a hierarchical attention mechanism attends to important parts within and among utterances with word level attention and utterance level attention respectively. With the word level attention, hidden vectors of a word level encoder are synthesized as utterance vectors and fed to an utterance level encoder to construct hidden representations of the context. The hidden vectors of the context are then processed by the utterance level attention and formed as context vectors for decoding the response. Empirical studies on both automatic evaluation and human judgment show that HRAN can significantly outperform state-of-the-art models for multi-turn response generation.

1 Introduction

Conversational agents include task-oriented dialog systems which are built in vertical domains for specific tasks (Young et al., 2013; Boden, 2006; Wallace, 2009; Young et al., 2010), and non-task-oriented chatbots which aim to realize natural and human-like conversations with people

*The work was done when the first author was an intern in Microsoft Research Asia.

Context	
u₁ (Speaker A):	征男友, 160cm的妹子真的找不到男友吗 I want a boyfriend. Why can't a 160cm girl find a boyfriend?
u₂ (Speaker B):	你找不到一定不是因为160 It's definitely not because you are 160cm.
u₃ (Speaker A):	我知道脸也是硬伤嘛 Well I know I'm not good-looking
u₄ (Speaker B):	是你非要175以上 No, it's because you always hit on someone higher than 175cm.
Response Candidates	
身高不是硬性要求	✓
No, I don't care much about height.	
你是男的还是女的啊	✗
Are you a man or a woman?	

Figure 1: An example of multi-turn conversation

regarding to a wide range of issues in open domains (Jafarpour et al., 2010). A common practice to build a chatbot is to learn a response generation model within an encoder-decoder framework from large scale message-response pairs (Shang et al., 2015; Vinyals and Le, 2015). Such models ignore conversation history when responding, which is contradictory to the nature of real conversation between humans. To resolve the problem, researchers have taken conversation history into consideration and proposed response generation for multi-turn conversation (Sordoni et al., 2015; Serban et al., 2015; Serban et al., 2016b; Serban et al., 2016c).

In this work, we study multi-turn response generation for open domain conversation in chatbots in which we try to learn a response generation model from responses and their contexts. A context refers to a message and several utterances in its previous turns. In practice, when a message comes, the model takes the context as input and generate a response as the next turn. Multi-turn conversation requires a model to generate a response relevant to the whole context. The complexity of the task lies in two aspects: 1) a conversation context is in a hierarchical structure (words form an utterance, and utterances form the context) and has two levels of sequential relationships among both words and utterances within the struc-

ture; 2) not all parts of the context are equally important to response generation. Words are differentially informative and important, and so are the utterances. State-of-the-art methods such as HRED (Serban et al., 2016a) and VHRED (Serban et al., 2016c) focus on modeling the hierarchy of the context, whereas there is little exploration on how to select important parts from the context, although it is often a crucial step for generating a proper response. Without this step, existing models may lose important information in context and generate irrelevant responses¹. Figure 1 gives an example from our data to illustrate the problem. The context is a conversation between two speakers about height and boyfriend, therefore, to respond to the context, words like “girl”, “boyfriend” and numbers indicating height such as “160” and “175” are more important than “not good-looking”. Moreover, u_1 and u_4 convey main semantics of the context, and therefore are more important than the others for generating a proper response. Without modeling the word and utterance importance, the state-of-the-art model VHRED (Serban et al., 2016c) misses important points and gives a response “are you a man or a woman” which is OK if there were only u_3 left, but nonsense given the whole context. After paying attention to the important words and utterances, we can have a reasonable response like “No, I don’t care much about height” (the response is generated by our model, as will be seen in experiments).

We aim to model the hierarchy and the important parts of contexts in a unified framework. Inspired by the success of the attention mechanism in single-turn response generation (Shang et al., 2015), we propose a hierarchical recurrent attention network (HRAN) for multi-turn response generation in which we introduce a hierarchical attention mechanism to dynamically highlight important parts of word sequences and the utterance sequence when generating a response. Specifically, HRAN is built in a hierarchical structure. At the bottom of HRAN, a word level recurrent neural network (RNN) encodes each utterance into a sequence of hidden vectors. In generation of each word in the response, a word level attention mech-

anism assigns a weight to each vector in the hidden sequence of an utterance and forms an utterance vector by a linear combination of the vectors. Important hidden vectors correspond to important parts in the utterance regarding to the generation of the word, and contribute more to the formation of the utterance vector. The utterance vectors are then fed to an utterance level RNN which constructs hidden representations of the context. Different from classic attention mechanism, the word level attention mechanism in HRAN is dependent on both the decoder and the utterance level RNN. Thus, both the current generated part of the response and the content of context can help select important parts in utterances. At the third layer, an utterance attention mechanism attends to important utterances in the utterance sequence and summarizes the sequence as a context vector. Finally, at the top of HRAN, a decoder takes the context vector as input and generates the word in the response. HRAN mirrors the data structure in multi-turn response generation by growing from words to utterances and then from utterances to the output. It extends the architecture of current hierarchical response generation models by a hierarchical attention mechanism which not only results in better generation quality, but also provides insight into which parts in an utterance and which utterances in context contribute to response generation.

We conduct an empirical study on large scale open domain conversation data and compare our model with state-of-the-art models using both automatic evaluation and side-by-side human comparison. The results show that on both metrics our model can significantly outperform existing models for multi-turn response generation. We release our source code and data at <https://github.com/LynetteXing1991/HRAN>.

The contributions of the paper include (1) proposal of attending to important parts in contexts in multi-turn response generation; (2) proposal of a hierarchical recurrent attention network which models hierarchy of contexts, word importance, and utterance importance in a unified framework; (3) empirical verification of the effectiveness of the model by both automatic evaluation and human judgment.

2 Related Work

Most existing effort on response generation is paid to single-turn conversation. Starting from the ba-

¹Note that one can simply concatenate all utterances and employs the classic sequence-to-sequence with attention to model word importance in generation. This method, however, loses utterance relationships and results in bad generation quality, as will be seen in experiments.

sic sequence to sequence model (Sutskever et al., 2014), various models (Shang et al., 2015; Vinyals and Le, 2015; Li et al., 2015; Xing et al., 2016; Li et al., 2016; ?) have been proposed under an encoder-decoder framework to improve generation quality from different perspectives such as relevance, diversity, and personality. Recently, multi-turn response generation has drawn attention from academia. For example, Sordani et al. (2015) proposed DCGM where context information is encoded with a multi-layer perceptron (MLP). Serban et al. (2016a) proposed HRED which models contexts in a hierarchical encoder-decoder framework. Under the architecture of HRED, more variants including VHRED (Serban et al., 2016c) and MrRNN (Serban et al., 2016b) are proposed in order to introduce latent and explicit variables into the generation process. In this work, we also study multi-turn response generation. Different from the existing models which do not model word and utterance importance in generation, our hierarchical recurrent attention network simultaneously models the hierarchy of contexts and the importance of words and utterances in a unified framework.

Attention mechanism is first proposed for machine translation (Bahdanau et al., 2014; Cho et al., 2015), and is quickly applied to single-turn response generation afterwards (Shang et al., 2015; Vinyals and Le, 2015). Recently, Yang et al. (2016) proposed a hierarchical attention network for document classification in which two levels of attention mechanisms are used to model the contributions of words and sentences in classification decision. Seo et al. (2016) proposed a hierarchical attention network to precisely attending objects of different scales and shapes in images. Inspired by these work, we extend the attention mechanism for single-turn response generation to a hierarchical attention mechanism for multi-turn response generation. To the best of our knowledge, we are the first who apply the hierarchical attention technique to response generation in chatbots.

3 Problem Formalization

Suppose that we have a data set $\mathcal{D} = \{(\mathbf{U}_i, \mathbf{Y}_i)\}_{i=1}^N$. $\forall i$, $(\mathbf{U}_i, \mathbf{Y}_i)$ consists of a response $\mathbf{Y}_i = (y_{i,1}, \dots, y_{i,T_i})$ and its context $\mathbf{U}_i = (u_{i,1}, \dots, u_{i,m_i})$ with $y_{i,j}$ the j -th word, u_{i,m_i} the message, and $(u_{i,1}, \dots, u_{i,m_i-1})$ the utterances in previous turns. In this work, we require $m_i \geq 2$ and thus each context has at

least one utterance as conversation history. $\forall j$, $u_{i,j} = (w_{i,j,1}, \dots, w_{i,j,T_{i,j}})$ where $w_{i,j,k}$ is the k -th word. We aim to estimate a generation probability $p(y_1, \dots, y_T | \mathbf{U})$ from \mathcal{D} , and thus given a new conversation context \mathbf{U} , we can generate a response $\mathbf{Y} = (y_1, \dots, y_T)$ according to $p(y_1, \dots, y_T | \mathbf{U})$.

In the following, we will elaborate how to construct $p(y_1, \dots, y_T | \mathbf{U})$ and how to learn it.

4 Hierarchical Recurrent Attention Network

We propose a hierarchical recurrent attention network (HRAN) to model the generation probability $p(y_1, \dots, y_T | \mathbf{U})$. Figure 2 gives the architecture of HRAN. Roughly speaking, before generation, HRAN employs a word level encoder to encode information of every utterance in context as hidden vectors. Then, when generating every word, a hierarchical attention mechanism attends to important parts within and among utterances with word level attention and utterance level attention respectively. With the two levels of attention, HRAN works in a bottom-up way: hidden vectors of utterances are processed by the word level attention and uploaded to an utterance level encoder to form hidden vectors of the context. Hidden vectors of the context are further processed by the utterance level attention as a context vector and uploaded to the decoder to generate the word.

In the following, we will describe details and the learning objective of HRAN.

4.1 Word Level Encoder

Given $\mathbf{U} = (u_1, \dots, u_m)$, we employ a bidirectional recurrent neural network with gated recurrent units (BiGRU) (Bahdanau et al., 2014) to encode each $u_i, i \in \{1, \dots, m\}$ as hidden vectors $(\mathbf{h}_{i,1}, \dots, \mathbf{h}_{i,T_i})$. Formally, suppose that $u_i = (w_{i,1}, \dots, w_{i,T_i})$, then $\forall k \in \{1, \dots, T_i\}$, $\mathbf{h}_{i,k}$ is given by

$$\mathbf{h}_{i,k} = \text{concat}(\overrightarrow{\mathbf{h}}_{i,k}, \overleftarrow{\mathbf{h}}_{i,k}), \quad (1)$$

where $\text{concat}(\cdot, \cdot)$ is an operation defined as concatenating the two arguments together, $\overrightarrow{\mathbf{h}}_{i,k}$ is the k -th hidden state of a forward GRU (Cho et al., 2014), and $\overleftarrow{\mathbf{h}}_{i,k}$ is the k -th hidden state of a backward GRU. The forward GRU reads u_i in its order

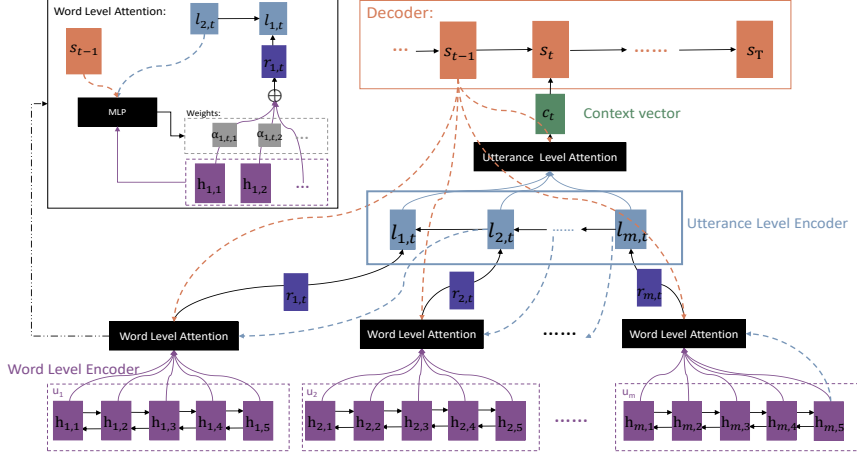


Figure 2: Hierarchical Recurrent Attention Network

(i.e., from $w_{i,1}$ to w_{i,T_i}), and calculates $\vec{\mathbf{h}}_{i,k}$ as

$$\begin{aligned}
 \mathbf{z}_k &= \sigma(\mathbf{W}_z \mathbf{e}_{i,k} + \mathbf{V}_z \vec{\mathbf{h}}_{i,k-1}) \\
 \mathbf{r}_k &= \sigma(\mathbf{W}_r \mathbf{e}_{i,k} + \mathbf{V}_r \vec{\mathbf{h}}_{i,k-1}) \\
 \mathbf{s}_k &= \tanh(\mathbf{W}_s \mathbf{e}_{i,k} + \mathbf{V}_s (\vec{\mathbf{h}}_{i,k-1} \circ \mathbf{r}_k)) \\
 \vec{\mathbf{h}}_{i,k} &= (1 - \mathbf{z}_k) \circ \mathbf{s}_k + \mathbf{z}_k \circ \vec{\mathbf{h}}_{i,k-1},
 \end{aligned} \tag{2}$$

where $\vec{\mathbf{h}}_{i,0}$ is initialized with a isotropic Gaussian distribution, $\mathbf{e}_{i,k}$ is the embedding of $w_{i,k}$, \mathbf{z}_k and \mathbf{r}_k are an update gate and a reset gate respectively, $\sigma(\cdot)$ is a sigmoid function, and $\mathbf{W}_z, \mathbf{W}_r, \mathbf{W}_s, \mathbf{V}_z, \mathbf{V}_r, \mathbf{V}_s$ are parameters. The backward GRU reads u_i in its reverse order (i.e., from w_{i,T_i} to $w_{i,1}$) and generates $\{\vec{\mathbf{h}}_{i,k}\}_{k=1}^{T_i}$ with a parameterization similar to the forward GRU.

4.2 Hierarchical Attention and Utterance Encoder

Suppose that the decoder has generated $t - 1$ words, at step t , word level attention calculates a weight vector $(\alpha_{i,t,1}, \dots, \alpha_{i,t,T_i})$ (details are described later) for $\{\mathbf{h}_{i,j}\}_{j=1}^{T_i}$ and represents utterance u_i as a vector $\mathbf{r}_{i,t}$. $\forall i \in \{1, \dots, m\}$, $\mathbf{r}_{i,t}$ is defined by

$$\mathbf{r}_{i,t} = \sum_{j=1}^{T_i} \alpha_{i,t,j} \mathbf{h}_{i,j}. \tag{3}$$

$\{\mathbf{r}_{i,t}\}_{i=1}^m$ are then utilized as input of an utterance level encoder and transformed to $(\mathbf{l}_{1,t}, \dots, \mathbf{l}_{m,t})$ as hidden vectors of the context. After that, utterance level attention assigns a weight $\beta_{i,t}$ to $\mathbf{l}_{i,t}$ (details are described later) and forms a context vector \mathbf{c}_t as

$$\mathbf{c}_t = \sum_{i=1}^m \beta_{i,t} \mathbf{l}_{i,t}. \tag{4}$$

In both Equation (3) and Equation (4), the more important a hidden vector is, the larger weight it

will have, and the more contributions it will make to the high level vector (i.e., the utterance vector and the context vector). This is how the two levels of attention attends to the important parts of utterances and the important utterances in generation.

More specifically, the utterance level encoder is a backward GRU which processes $\{\mathbf{r}_{i,t}\}_{i=1}^m$ from the message $\mathbf{r}_{m,t}$ to the earliest history $\mathbf{r}_{1,t}$. Similar to Equation (2), $\forall i \in \{m, \dots, 1\}$, $\mathbf{l}_{i,t}$ is calculated as

$$\begin{aligned}
 \mathbf{z}'_i &= \sigma(\mathbf{W}_{zl} \mathbf{r}_{i,t} + \mathbf{V}_{zl} \mathbf{l}_{i+1,t}) \\
 \mathbf{r}'_i &= \sigma(\mathbf{W}_{rl} \mathbf{r}_{i,t} + \mathbf{V}_{rl} \mathbf{l}_{i+1,t}) \\
 \mathbf{s}'_i &= \tanh(\mathbf{W}_{sl} \mathbf{r}_{i,t} + \mathbf{V}_{sl} (\mathbf{l}_{i+1,t} \circ \mathbf{r}'_i)) \\
 \mathbf{l}_{i,t} &= (1 - \mathbf{z}'_i) \circ \mathbf{s}'_i + \mathbf{z}'_i \circ \mathbf{l}_{i+1,t},
 \end{aligned} \tag{5}$$

where $\mathbf{l}_{m+1,t}$ is initialized with a isotropic Gaussian distribution, \mathbf{z}'_i and \mathbf{r}'_i are the update gate and the reset gate of the utterance level GRU respectively, and $\mathbf{W}_{zl}, \mathbf{V}_{zl}, \mathbf{W}_{rl}, \mathbf{V}_{rl}, \mathbf{W}_{sl}, \mathbf{V}_{sl}$ are parameters.

Different from the classic attention mechanism, word level attention in HRAN depends on both the hidden states of the decoder and the hidden states of the utterance level encoder. It works in a reverse order by first weighting $\{\mathbf{h}_{m,j}\}_{j=1}^{T_m}$ and then moving towards $\{\mathbf{h}_{1,j}\}_{j=1}^{T_1}$ along the utterance sequence. $\forall i \in \{m, \dots, 1\}, j \in \{1, \dots, T_i\}$, weight $\alpha_{i,t,j}$ is calculated as

$$\begin{aligned}
 e_{i,t,j} &= \eta(\mathbf{s}_{t-1}, \mathbf{l}_{i+1,t}, \mathbf{h}_{i,j}); \\
 \alpha_{i,t,j} &= \frac{\exp(e_{i,t,j})}{\sum_{k=1}^{T_i} \exp(e_{i,t,k})},
 \end{aligned} \tag{6}$$

where $\mathbf{l}_{m+1,t}$ is initialized with a isotropic Gaussian distribution, \mathbf{s}_{t-1} is the $(t - 1)$ -th hidden state of the decoder, and $\eta(\cdot)$ is a multi-layer perceptron (MLP) with tanh as an activation function.

Note that the word level attention and the utterance level encoding are dependent with each other and alternatively conducted (first attention then encoding). The motivation we establish the dependency between $\alpha_{i,t,j}$ and $\mathbf{l}_{i+1,t}$ is that content from the context (i.e., $\mathbf{l}_{i+1,t}$) could help identify important information in utterances, especially when \mathbf{s}_{t-1} is not informative enough (e.g., the generated part of the response are almost function words). We require the utterance encoder and the word level attention to work reversely, because we think that compared to history, conversation that happened after an utterance in the context is more likely to be capable of identifying important information in the utterance for generating a proper response to the context.

With $\{\mathbf{l}_{i,t}\}_{i=1}^m$, the utterance level attention calculates a weight $\beta_{i,t}$ for $\mathbf{l}_{i,t}$ as

$$\begin{aligned} e'_{i,t} &= \eta(\mathbf{s}_{t-1}, \mathbf{l}_{i,t}); \\ \beta_{i,t} &= \frac{\exp(e'_{i,t})}{\sum_{i=1}^m \exp(e'_{i,t})}. \end{aligned} \quad (7)$$

4.3 Decoding the Response

The decoder of HRAN is a RNN language model (Mikolov et al., 2010) conditioned on the context vectors $\{\mathbf{c}_t\}_{t=1}^T$ given by Equation (4). Formally, the probability distribution $p(y_1, \dots, y_T | \mathbf{U})$ is defined as

$$p(y_1, \dots, y_T | \mathbf{U}) = p(y_1 | \mathbf{c}_1) \prod_{t=2}^T p(y_t | \mathbf{c}_t, y_1, \dots, y_{t-1}). \quad (8)$$

where $p(y_t | \mathbf{c}_t, y_1, \dots, y_{t-1})$ is given by

$$\begin{aligned} \mathbf{s}_t &= f(\mathbf{e}_{y_{t-1}}, \mathbf{s}_{t-1}, \mathbf{c}_t) \\ p(y_t | \mathbf{c}_t, y_1, \dots, y_{t-1}) &= \mathbb{I}_{y_t} \cdot \text{softmax}(\mathbf{s}_t, \mathbf{e}_{y_{t-1}}), \end{aligned} \quad (9)$$

where \mathbf{s}_t is the hidden state of the decoder at step t , $\mathbf{e}_{y_{t-1}}$ is the embedding of y_{t-1} , f is a GRU, \mathbb{I}_{y_t} is the one-hot vector for y_t , and $\text{softmax}(\mathbf{s}_t, \mathbf{e}_{y_{t-1}})$ is a V -dimensional vector with V the response vocabulary size and each element the generation probability of a word. In practice, we employ the beam search (Tillmann and Ney, 2003) technique to generate the n -best responses.

Let us denote Θ as the parameter set of HRAN, then we estimate Θ from $\mathcal{D} = \{(\mathbf{U}_i, \mathbf{Y}_i)\}_{i=1}^N$ by minimizing the following objective function:

$$\hat{\Theta} = \arg \min_{\Theta} - \sum_{i=1}^N \log(p(y_{i,1}, \dots, y_{i,T_i} | \mathbf{U}_i)) \quad (10)$$

5 Experiments

We compared HRAN with state-of-the-art methods by both automatic evaluation and side-by-side human judgment.

5.1 Data Set

We built a data set from Douban Group² which is a popular Chinese social networking service (SNS) allowing users to discuss a wide range of topics in groups through posting and commenting. In Douban Group, regarding to a post under a specific topic, two persons can converse with each other by one posting a comment and the other quoting it and posting another comment. We crawled 20 million conversations between two persons with the average number of turns as 6.32. The data covers many different topics and can be viewed as a simulation of open domain conversations in a chatbot. In each conversation, we treated the last turn as response, and the remaining turns as context. As preprocessing, we first employed Stanford Chinese word segmenter³ to tokenize each utterance in the data. Then we removed the conversations whose response appearing more than 50 times in the whole data to prevent them from dominating learning. We also removed the conversations shorter than 3 turns and the conversations with an utterance longer than 50 words. After the preprocessing, there are 1,656,652 conversations left. From them, we randomly sampled 1 million conversations as training data, 10,000 conversations as validation data, and 1,000 conversations as test data, and made sure that there is no overlap among them. In the test data, the contexts were used to generate responses and their responses were used as ground truth to calculate perplexity of generation models. We kept the 40,000 most frequent words in the contexts of the training data to construct a context vocabulary. The vocabulary covers 98.8% of words appearing in the contexts of the training data. Similarly, we constructed a response vocabulary that contains the 40,000 most frequent words in the responses of the training data which covers 99.0% words appearing in the responses. Words outside the two vocabularies were treated as ‘‘UNK’’. The data will be publicly available.

²<https://www.douban.com/group/explore>

³<http://nlp.stanford.edu/software/segmenter.shtml>

Model	Validation Perplexity	Test Perplexity
S2SA	43.679	44.508
HRED	46.279	47.467
VHRED	44.548	45.484
HRAN	40.257	41.138

Table 1: Perplexity results

5.2 Baselines

We compared HRAN with the following models:

S2SA: we concatenated all utterances in a context as a long sequence and treated the sequence and the response as a message-response pair. By this means, we transformed the problem of multi-turn response generation to a problem of single-turn response generation and employed the standard sequence to sequence with attention (Shang et al., 2015) as a baseline.

HRED: the hierarchical encoder-decoder model proposed by (Serban et al., 2016a).

VHRED: a modification of HRED (Serban et al., 2016c) where latent variables are introduced in to generation. In all models, we set the dimensionality of hidden states of encoders and decoders as 1000, and the dimensionality of word embedding as 620. All models were initialized with isotropic Gaussian distributions $\mathcal{X} \sim \mathcal{N}(0, 0.01)$ and trained with an AdaDelta algorithm (Zeiler, 2012) on a NVIDIA Tesla K40 GPU. The batch size is 128. We set the initial learning rate as 1.0 and reduced it by half if the perplexity on validation began to increase. We implemented the models with an open source deep learning tool Blocks⁴.

5.3 Evaluation Metrics

How to evaluate a response generation model is still an open problem but not the focus of the paper. We followed the existing work and employed the following metrics:

Perplexity: following (Vinyals and Le, 2015), we employed perplexity as an evaluation metric. Perplexity is defined by Equation (11). It measures how well a model predicts human responses. Lower perplexity generally indicates better generation performance. In our experiments, perplexity on validation was used to determine when to stop training. If the perplexity stops decreasing and the difference is smaller than 2.0 five times in validation, we think that the algorithm has reached convergence and terminate training. We tested the generation ability of different models by perplex-

Models	Win	Loss	Tie	Kappa
HRAN v.s. S2SA	27.3	20.6	52.1	0.37
HRAN v.s. HRED	27.2	21.2	51.6	0.35
HRAN v.s. VHRED	25.2	20.4	54.4	0.34

Table 2: Human annotation results (in %)

ity on the test data.

$$PPL = exp \left\{ -\frac{1}{N} \sum_{i=1}^N \log(p(\mathbf{Y}_i | \mathbf{U}_i)) \right\}. \quad (11)$$

Side-by-side human annotation: we also compared HRAN with every baseline model by side-by-side human comparison. Specifically, we recruited three native speakers with rich Douban Group experience as human annotators. To each annotator, we showed a context of a test example with two generated responses, one from HRAN and the other one from a baseline model. Both responses are the top one results in beam search. The two responses were presented in random order. We then asked the annotator to judge which one is better. The criteria is, response A is better than response B if (1) A is relevant, logically consistent to the context, and fluent, while B is either irrelevant or logically contradictory to the context, or it is disfluent (e.g., with grammatical errors or UNKs); or (2) both A and B are relevant, consistent, and fluent, but A is more informative and interesting than B (e.g., B is a universal reply like “I see”). If the annotator cannot tell which one is better, he/she was asked to label a “tie”. Each annotator individually judged 1000 test examples for each HRAN-baseline pair, and in total, each one judged 3000 examples (for three pairs). Agreements among the annotators were calculated using Fleiss’ kappa (Fleiss and Cohen, 1973).

Note that we do not choose BLEU (Papineni et al., 2002) as an evaluation metric, because (1) Liu et al. (Liu et al., 2016) have proven that BLEU is not a proper metric for evaluating conversation models as there is weak correlation between BLEU and human judgment; (2) different from the single-turn case, in multi-turn conversation, one context usually has one copy in the whole data. Thus, without any human effort like what Sordoni et al. (Sordoni et al., 2015) did in their work, each context only has a single reference in test. This makes BLEU even unreliable as a measurement of generation quality in open domain conversation due to the diversity of responses.

5.4 Evaluation Results

Table 1 gives the results on perplexity. HRAN achieves the lowest perplexity on both validation

⁴<https://github.com/mila-udem/blocks>

ID	Context	Response
1	求唇部死皮怎么去⇒用化妆水一遍遍的擦嘴唇就能擦下来⇒好的回去试试⇒我是能擦下来 How to exfoliate my lips?⇒Wipe your lips with lotion again and again. ⇒OK, let me try it. ⇒It works for me.	HRAN: 你用的什么牌子的化妆水 Which brand of lotion are you using? VHRED: 嗯我知道了 OK I see.
2	我不能去, 你可以找丽丽陪你吃饭⇒她住松江, 离市区太远了⇒有好吃的多远都要去啊⇒你为什么不能来呢 I can't have dinner with you, but you can ask Lily instead.⇒She lives in Songjiang, too far away from downtown.⇒I think she'd like to go because of the delicious food.⇒ Why can't you come?	HRAN: 吃多了上火了 Because I've eaten too much and gotten inflamed. VHRED: 因为我是男的 Because I'm a man.
3	征男友, 160cm的妹子真的找不到男友吗⇒你找不到一定不是因为160⇒我知道脸也是硬伤嘛⇒是你非要175以上 I want a boyfriend. Why can't a 160cm girl find a boyfriend? ⇒It's definitely not because you are 160cm.⇒ Well I know I'm not good-looking. ⇒No, it's because you always hit on someone higher than 175cm.	HRAN: 身高不是硬性要求 No, I don't care much about height. VHRED: 你是男的还是女的啊 Are you a man or a woman?
4	啦啦啦销售助理在哪里呀在哪里⇒学历不达标⇒对做销售有兴趣么⇒没做过不过我不挑工作 I'm looking for a sales assistant.⇒I might not be qualified because of my education.⇒Are you interested in sales? ⇒I have no experience before but I'm open to any jobs.	HRAN: 那你做过设计么 Then have you done any designing jobs before? VHRED: 好吧好吧我也觉得 OK, I think so.

Figure 3: Case study (utterances between two persons in contexts are split by “⇒”)

and test. We conducted t-test on test perplexity and the result shows that the improvement of HRAN over all baseline models is statistically significant (p -value < 0.01).

Table 2 shows the human annotation results. The ratios were calculated by combining the annotations from the three judges together. We can see that HRAN outperforms all baseline models and all comparisons have relatively high kappa scores which indicates that the annotators reached relatively high agreements in judgment. Compared with S2SA, HRED, and VHRED, HRAN achieves preference gains (win-loss) 6.7%, 6%, 4.8% respectively. Sign test results show that the improvement is statistically significant (p -value < 0.01 for HRAN v.s. S2SA and HRAN v.s. HRED, and p -value < 0.05 for HRAN v.s. VHRED). Among the three baseline models, S2SA is the worst one, because it loses relationships among utterances in the context. VHRED is the best baseline model, which is consistent with the existing literatures (Serban et al., 2016c). We checked the cases on which VHRED loses to HRAN and found that on 56% cases, VHRED generated irrelevant responses while responses given by HRAN are relevant, logically consistent, and fluent.

5.5 Discussions

Case study: Figure 3 lists some cases from the test set to compare HRAN with the best baseline VHRED. We can see that HRAN not only can answer the last turn in the context (i.e., the mes-

Model	Win	Loss	Tie	PPL
No UD Att	22.3%	24.8%	52.9%	41.54
No Word Att	20.4%	25.0%	50.6%	43.24
No Utterance Att	21.1%	23.7%	55.2%	47.35

Table 3: Model ablation results

sage) properly by understanding the context (e.g., case 2), but also be capable of starting a new topic according to the conversation history to keep the conversation going (e.g., case 1). In case 2, HRAN understands that the message is actually asking “why can’t you come to have dinner with me?” and generates a proper response that gives a plausible reason. In case 1, HRAN properly brings up a new topic by asking the “brand” of the user’s “lotion” when the current topic “how to exfoliate my skin” has come to an end. The new topic is based on the content of the context and thus can naturally extend the conversation in the case.

Visualization of attention: to further illustrate why HRAN can generate high quality responses, we visualized the hierarchical attention for the cases in Figure 3 in Figure 4. In every sub-figure, each line is an utterance with blue color indicating word importance. The leftmost column of each sub-figure uses red color to indicate utterance importance. Darker color means more important words or utterances. The importance of a word or an utterance was calculated by the average weight of the word or the utterance assigned by attention in generating the response given at the bottom of each sub-figure. It reflects an overall contribution of the word or the utterance to generate the response. Above each line, we gave the transla-

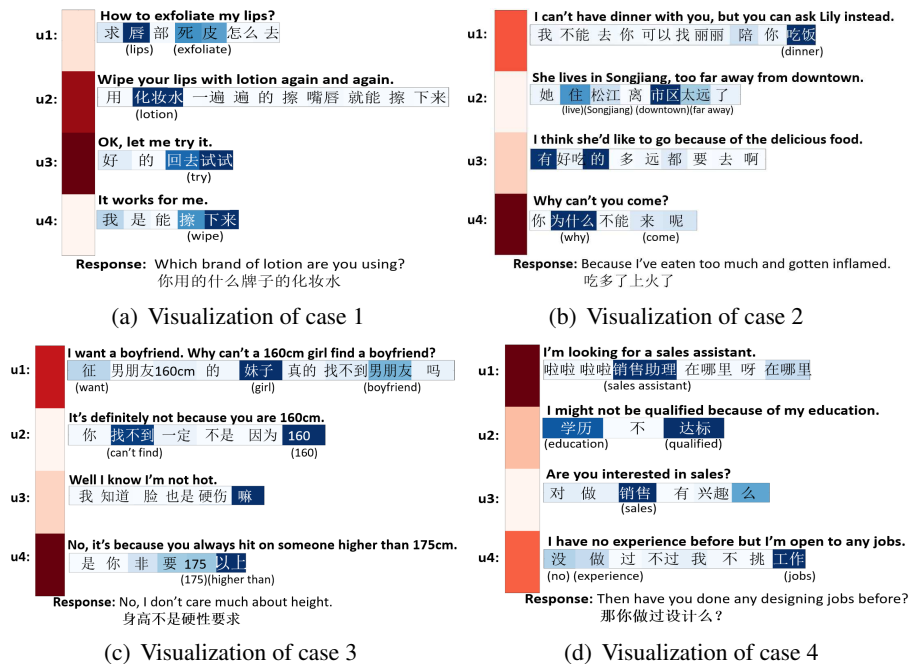


Figure 4: Attention visualization (the importance of a word or an utterance is calculated as their average weights when generating the whole response)

tion of the utterance, and below it, we translated important words. Note that word-to-word translation may cause confusion sometimes, therefore, we left some words (most of them are function words) untranslated. We can see that the hierarchical attention mechanism in HRAN can attend to both important words and important utterances in contexts. For example, in Figure 4(c), words including “girl” and “boyfriend” and numbers including “160” and “175” are highlighted, and u_1 and u_4 are more important than others. The result matches our intuition in introduction. In Figure 4(b), HRAN assigned larger weights to u_1 , u_4 and words like “dinner” and “why”. This explains why the model can understand that the message is actually asking “why can’t you come to have dinner with me?”. The figures provide us insights on how HRAN understands contexts in generation.

Model ablation: we then examine the effect of different components of HRAN by removing them one by one. We first removed l_{i+1} from $\eta(s_{t-1}, l_{i+1,t}, h_{i,j})$ in Equation (6) (i.e., removing utterance dependency from word level attention) and denoted the model as “No UD Att”, then we removed word level attention and utterance level attention separately, and denoted the models as “No Word Att” and “No Utterance Att” respectively. We conducted side-by-side human comparison on these models with the full HRAN on the test data and also calculated their test perplexity (PPL). Table 3 gives the results. We can see that

all the components are useful because removing any of them will cause performance drop. Among them, word level attention is the most important one as HRAN achieved the most preference gain (4.6%) to No Word Att on human comparison.

Error analysis: we finally investigate how to improve HRAN in the future by analyzing the cases on which HRAN loses to VHRED. The errors can be summarized as: 51.81% logic contradiction, 26.95% universal reply, 7.77% irrelevant response, and 13.47% others. Most bad cases come from universal replies and responses that are logically contradictory to contexts. This is easy to understand as HRAN does not explicitly model the two issues. The result also indicates that (1) although contexts provide more information than single messages, multi-turn response generation still has the “safe response” problem as the single-turn case; (2) although attending to important words and utterances in generation can lead to informative and logically consistent responses for many cases like those in Figure 3, it is still not enough for fully understanding contexts due to the complex nature of conversations. The irrelevant responses might be caused by wrong attention in generation. Although the analysis might not cover all bad cases (e.g., HRAN and VHRED may both give bad responses), it sheds light on our future directions: (1) improving response diversity, e.g., by introducing extra content into generation like Xing et al. (Xing et al., 2016) and Mou et al. (Mou et al.,

2016) did for single-turn conversation; (2) modeling logics in contexts; (3) improving attention.

6 Conclusion

We propose a hierarchical recurrent attention network (HRAN) for multi-turn response generation in chatbots. Empirical studies on large scale conversation data show that HRAN can significantly outperform state-of-the-art models.

References

- Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. 2014. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*.
- Margaret Ann Boden. 2006. *Mind as machine: A history of cognitive science*. Clarendon Press.
- Kyunghyun Cho, Bart van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. 2014. On the properties of neural machine translation: Encoder–decoder approaches. *Syntax, Semantics and Structure in Statistical Translation*, page 103.
- Kyunghyun Cho, Aaron Courville, and Yoshua Bengio. 2015. Describing multimedia content using attention-based encoder-decoder networks. *Multimedia, IEEE Transactions on*, 17(11):1875–1886.
- Joseph L Fleiss and Jacob Cohen. 1973. The equivalence of weighted kappa and the intraclass correlation coefficient as measures of reliability. *Educational and psychological measurement*.
- Sina Jafarpour, Christopher JC Burges, and Alan Ritter. 2010. Filter, rank, and transfer the knowledge: Learning to chat. *Advances in Ranking*, 10.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2015. A diversity-promoting objective function for neural conversation models. *arXiv preprint arXiv:1510.03055*.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016. A persona-based neural conversation model. *arXiv preprint arXiv:1603.06155*.
- 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*.
- Tomas Mikolov, Martin Karafiát, Lukas Burget, Jan Cernocký, and Sanjeev Khudanpur. 2010. Recurrent neural network based language model. In *INTERSPEECH 2010, 11th Annual Conference of the International Speech Communication Association, Makuhari, Chiba, Japan, September 26-30, 2010*, pages 1045–1048.
- Lili Mou, Yiping Song, Rui Yan, Ge Li, Lu Zhang, and Zhi Jin. 2016. Sequence to backward and forward sequences: A content-introducing approach to generative short-text conversation. *arXiv preprint arXiv:1607.00970*.
- 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. Association for Computational Linguistics.
- Paul Hongsuck Seo, Zhe Lin, Scott Cohen, Xiaohui Shen, and Bohyung Han. 2016. Hierarchical attention networks. *arXiv preprint arXiv:1606.02393*.
- Iulian V Serban, Alessandro Sordoni, Yoshua Bengio, Aaron Courville, and Joelle Pineau. 2015. Building end-to-end dialogue systems using generative hierarchical neural network models. *arXiv preprint arXiv:1507.04808*.
- Iulian V Serban, Alessandro Sordoni, Yoshua Bengio, Aaron Courville, and Joelle Pineau. 2016a. Building end-to-end dialogue systems using generative hierarchical neural network models. In *Proceedings of the 30th AAAI Conference on Artificial Intelligence (AAAI-16)*.
- Iulian Vlad Serban, Tim Klinger, Gerald Tesauro, Kartik Talamadupula, Bowen Zhou, Yoshua Bengio, and Aaron Courville. 2016b. Multiresolution recurrent neural networks: An application to dialogue response generation. *arXiv preprint arXiv:1606.00776*.
- Iulian Vlad Serban, Alessandro Sordoni, Ryan Lowe, Laurent Charlin, Joelle Pineau, Aaron Courville, and Yoshua Bengio. 2016c. A hierarchical latent variable encoder-decoder model for generating dialogues. *arXiv preprint arXiv:1605.06069*.
- Lifeng Shang, Zhengdong Lu, and Hang Li. 2015. Neural responding machine for short-text conversation. *arXiv preprint arXiv:1503.02364*.
- Alessandro Sordoni, Michel Galley, Michael Auli, Chris Brockett, Yangfeng Ji, Margaret Mitchell, Jian-Yun Nie, Jianfeng Gao, and Bill Dolan. 2015. A neural network approach to context-sensitive generation of conversational responses. *arXiv preprint arXiv:1506.06714*.
- 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.
- Christoph Tillmann and Hermann Ney. 2003. Word reordering and a dynamic programming beam search algorithm for statistical machine translation. *Computational linguistics*, 29(1):97–133.
- Oriol Vinyals and Quoc Le. 2015. A neural conversational model. *arXiv preprint arXiv:1506.05869*.

- Richard S Wallace. 2009. *The anatomy of ALICE*. Springer.
- Chen Xing, Wei Wu, Yu Wu, Jie Liu, Yalou Huang, Ming Zhou, and Wei-Ying Ma. 2016. Topic aware neural response generation. *arXiv preprint arXiv:1606.08340*.
- Zichao Yang, Diyi Yang, Chris Dyer, Xiaodong He, Alex Smola, and Eduard Hovy. 2016. Hierarchical attention networks for document classification. In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*.
- 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.
- Stephanie Young, Milica Gasic, Blaise Thomson, and John D Williams. 2013. Pomdp-based statistical spoken dialog systems: A review. *Proceedings of the IEEE*, 101(5):1160–1179.
- Matthew D Zeiler. 2012. Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*.