

Open-Domain Dialogue Generation: limitation and direction

Michael Shell School of Electrical and
Computer Engineering
Georgia Institute of Technology
Atlanta, Georgia 30332-0250

Email: <http://www.michaelshell.org/contact.html> Homer Simpson Twentieth Century Fox
Springfield, USA

Email: homer@thesimpsons.com James Kirk
and Montgomery Scott Starfleet Academy
San Francisco, California 96678-2391
Telephone: (800) 555-1212
Fax: (888) 555-1212



Abstract—With the development of deep learning, dialogue systems have received more and more attention from researchers, especially open-domain dialogue generation system. The application scenario of it, also chat-bot, is very extensive, from website customer service to dialogue robot such as Siri and Microsoft Xiaoice. In recent years, more and more studies have been done on open-domain dialogue generation. This paper will provide a structured introduction to the common techniques and research area of open-domain dialogue generation and present the limitation of each methods. We first introduce the basic problems and the direction of improvement, which is to relevance and diversity. Next, we will introduce how to achieve this goal from the perspective of representation and framework. Finally, we compare some traditional and novel evaluations to evaluate the performance of generation sentences. The purpose of this paper is to sort out the research results of recent years and make contributions to the research and development of the dialogue system in the future.

1 INTRODUCTION

Open domain dialog generation system has become an important role in research area now. Different from task-oriented dialog system([1], [2], [3], [4], [5], [6], [7], [8], [9]) which aims to deal with special tasks, the problem to be solved by this topic is how to let the computer understand human language and how to make it possible for the computer to conduct open questions and talks with humans([10], [11], [12], [13], [14], [15], [16]).

However, today's open-domain dialogue generation system is still far from the above goal. It mainly has the following four problems: (1)grammar mistake, (2)duplicate word, (3)relevance and (4)diversity. The first two refer to issues in the generated sentence. And these two issues can be promoted by improving the quality of training data. As for the third and fourth problem, we will focus on the next step. Both are based on the correctness of the language

itself, but their own meaning is not the same. Relevance refers to the similarity of response and query while diversity describes the richness and divergence of language. We will follow relevance and diversity in the following sections to sort out existing methods and inspire future directions.

This paper first analysed the research trends in recent years and proposed macro methodology of this paper. Next, we will elaborate on the existing methods according to different perspectives.

In order to promote the performance of dialogue system, researchers solve the dialogue problem from multiple perspectives. Figure 1 lists all these perspectives. Basically speaking, how to effectively model words will be the basis of all models. The following two famous models will be mentioned, word embedding([17], [18], [19], [20]) and bag-of-word. Both models will convert discrete words into vectors which the computer can understand. Based on the representation of words, we also need the computer to understand the meaning of each sentence. The following methods utilizes RNN([21], [22]) and CNN([23]) and their variants to model sentences. Based on the expression of sentences, some researchers believe that the context should be integrated into models in a hierarchical manner([10], [11]). Experiments have also proved the effectiveness of such methods. Furthermore, except extracting textual information, some models incorporate multiple metadata into word-level and sentence-level representations. These metadata include: emotion, domain or topic, speaker role and so on([24], [25], [26], [27], [28], [29]). In addition to adding metadata, researchers combine the sentence-level and dialog-level representations with the attention mechanism([30], [31], [12]) and the gate mechanism([13], [32]). Specifically, the attention mechanism obtains different weighted context representations for the hidden state of the previous time in

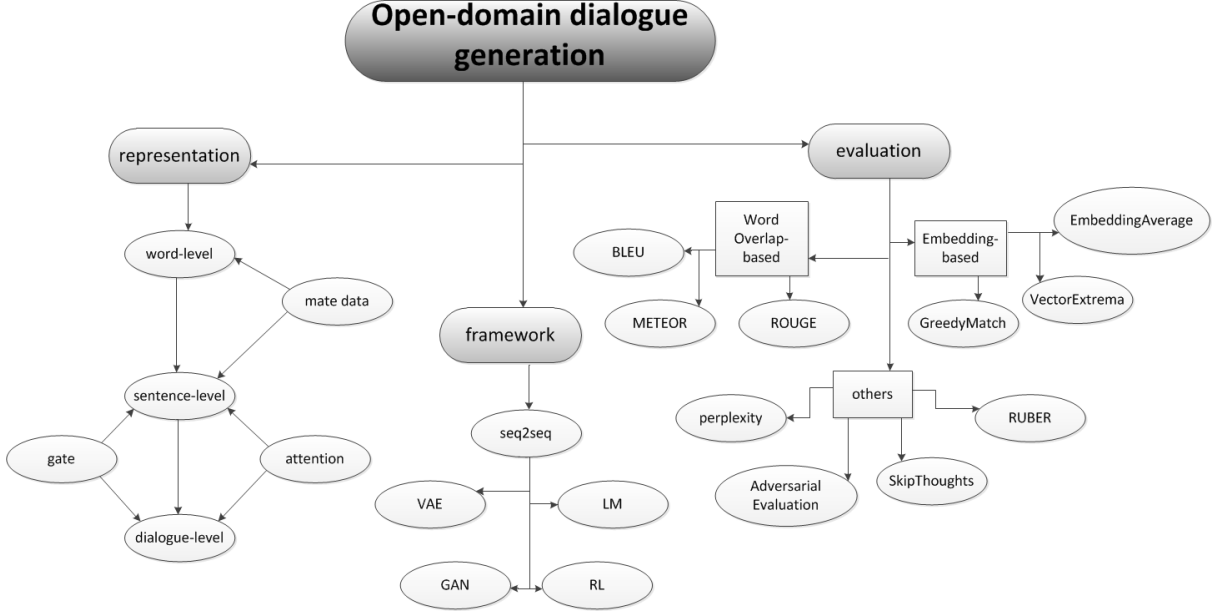


Fig. 1: The architecture of this paper.

each timestep of the generation phase, and then uses this context representation to generate the current word. The fundamental problem that the attention mechanism solves is to allow the network to return to the input sequence instead of encoding all the information into a fixed-length vector. On the other hand, the gate mechanism can be used as a control switch to control the inflow and outflow of information. Just as in LSTM, a cell has three gates: the forget gate, the input gate, and the output gate. They control which parts of the information are forgotten, which parts are input, and which parts are output.

In addition to the perspective of representation, some researchers study the dialogue from the perspective of the framework. The simplest and most basic of all frameworks is the seq2seq model. [33], [34] proposed the most basic seq2seq model, which is almost a baseline for all models. However, the biggest problem with this model is that it always generates a generic reply, such as *I dont know, I am not sure* and so on. To solve the problem of the universal reply in the seq2seq model, different researchers proposed different frameworks. Some researchers have proposed a VAE framework([35], [36]) for learning latent variables behind representation. However, there are still some problems in applying the VAE framework to text generation, which will be mentioned below. The researchers proposed *KL cost annealing* and *Word dropout and historyless decoding* methods to solve these problems, and further applied the VAE framework to the dialogue, and achieved good results([14], [15]). Other researchers proposed the GAN framework([37]). GAN has achieved great success in the image field at first. However, it cannot be applied to text generation directly due to non-differentiable reasons. Researchers have introduced Reinforcement Learning(RL) to bypass the problem of derivation and pass discriminator feedback directly to the generator([16]). There are also some novel methods to use the language model to handle conversations and have achieved good results([38]). Finally, some methods are

proposed to improve the performance of the model. We call these improvement.

All in all, the purpose of researchers is to increase the relevance and diversity of models, which is the main thread of this survey. We will elaborate it from three perspectives. Section 2 will introduce the trend of research in open-domain dialogue generation and the methodology of this review while section 3 and section 4 will describe the models from perspectives of representation and framework respectively. Finally, we will introduce the evaluation metrics in section 5 and give conclusion and future direction in section 6.

2 TREND AND REVIEW METHODOLOGY

The goal of this section aims to explain the trends of open-domain dialogue generation through an analysis of the papers in kinds of authoritative databases up to 2017 and then introduce the methodology of this paper.

2.1 Trend

To illustrate the trend of open-domain dialogue generation in the research field, we listed a number of papers from 2013 to 2017. All these papers have increased the performance of dialogue generation to some extent. Figure 2 and figure 3 display the changes in the amounts of publications with respect to open-domain dialogue generation from 2010 to 2017.

Where the numbers in the x axis means publication years and the numbers in the y axis is publication amounts. These four curves represents the most authoritative four publications or databases on computer science including (1) Association for Computing Machinery(ACM), (2)Springer, (3)Science Direct, (4)Institute of Electrical and Electronics Engineers(IEEE), and(5)arxiv.

From figure 2 and figure 3, it can be seen that open-domain dialogue generation has attracted more and more

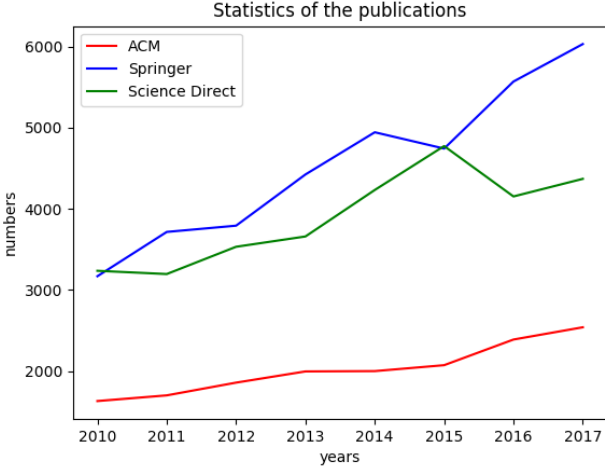


Fig. 2: Statistics of the publications: ACM, Springer and Science Direct.

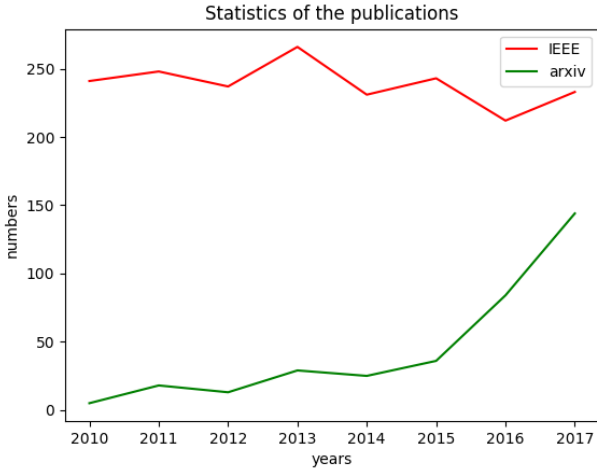


Fig. 3: Statistics of the publications: IEEE and arxiv.

researchers' interest, and the number of papers has been showing an increasing trend, especially in the last three years. Similarly, the applications of open-domain dialogue generation around us has gradually grown from Apple's Siri to Microsofts Cortana to Microsofts Xiaoice. These applications have greatly enriched our lives. Because of this, open-domain dialogue generation is the focus of the next research and we must sort out existing methods and further propose future development directions.

2.2 review methodology

There are many papers([39], [40], [41]) that have already reviewed the dialogue system. However, first of all, they are an overview of the entire dialogue system including task-oriented dialogue system, retrieval-based dialogue system and some hybrid methods, which can't sort out the dialogue generation system in detail, and is difficult to put forward the direction of improvement. Secondly, above all, these reviews classify various papers according to rough criteria such as context, personality, knowledge base and so on.

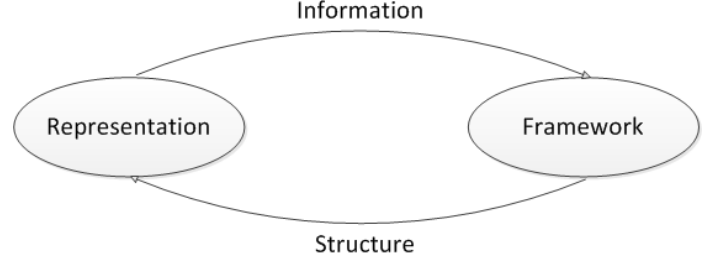


Fig. 4: The architecture of representation and framework.

Instead, this paper detailed classification of various methods: firstly, we decouple representation and framework, the representation is responsible for extracting information while the framework is responsible for completing the end-to-end connection. Of course, the two are not absolutely separate but promote each other. Figure 4 show this process.

Some methods are improved at the presentation level, and some methods are improved at the framework level. For example, [10] proposes HRED to incorporate context information while [42] introduce the *vae* framework to process context information. We will explain in detail below.

3 REPRESENTATION

Figure 1 summarizes the model structure described in this section. From a hierarchical perspective, it is layer by layer to represent a more abstract structure. From discrete words to corresponding word embedding, from multiple word embedding to abstract representations of sentences, then multiple sentences are modelled to represent dialogue. From a horizontal perspective, we can incorporate meta data or join the attention mechanism and the gate mechanism, readers can understand these as a pluggable structure. Finally, effectively using external knowledge and expressing them is also an important way to increase performance. Section 3.1 will talks about hierarchical representation and section 3.2 describes horizontal representation, then we introduce external representation in section 3.3 and give a summary in section 3.4.

3.1 Hierarchical representation

3.1.1 word representation

Many concepts that are difficult for humans to comprehend but are able to formulate are easy to understand for computers. On the contrary, concepts that are easily understood by humans but are not well formulated can be difficultly understand for computers. Words are one of such a concept. [17], [18], [19], [20] proposes word embedding as a representation of words, which achieves the famous equation, *king - man + woman = queen*, although this model has not done well enough in terms of polysemy. Other expressions of words such as the bag-of-word model assumes that for a document, we ignores its word order and grammar, syntax, etc., and considers it only as a collection of several words. The appearance of each word in the document is independent and does not depend on whether other words appear. Based on this model, [43] proposed to treat the context information and the message as a single sentence and use

the bag-of-word model to represent it, which is integrated into every timestep of the RNN and together with input and the hidden state of the previous timestep to generate the hidden state of the current moment. Furthermore, they express the context and message in a bag-of-word model, concatenate together the two representations and perform the same processing in the RNN as in the previous method. This model improves the relevance of the response by introducing the bag-of-word representation of context. In the following we will also see different ways of modelling the context, which can also improve the model's relevance.

3.1.2 sentence representation

Based on word representation, Researchers can use neural networks to abstract sentences. [21], [22] proposed RNN, which broke through the traditional MLP and can better handle sequence information. LSTM([44]) and GRU([45]) can handle long-term dependencies relative to RNN. On the other hand, [23] proposed CNN, which is originally applied primarily to images, TextCNN([46]) apply the structure of CNN to text and realize good results in text categorization tasks.

There are some methods for modelling sentences that have not yet been applied to the dialogue generation field, but given their good performance, it is also a good exploration to apply them to the dialogue generation.

Instead of encoding sentence content directly, RCNN([47]) encodes each words in a sentence into c_l , e_w , and c_r respectively and obtains an abstract representation of the sentence, then a max-pooling layer is selected to convey the most representative vector representation. In International Conference on Learning Representations(ICLR) 2018, [48] first uses an encoder to represent given sentence s and uses another encoder to represent some candidate sentences. Then they utilize a classifier to choose a candidate sentence with the highest co-occurrence rate to given sentence s . Since this model can choose to ignore aspects of the sentence that are irrelevant in constructing a semantic embedding space, it shows a good result. Based on [49] and [50], [51] put forward a universal sentence encoder toolkit with the transformer architecture([49]) and the deep averaging network (DAN)([50]) architecture.

Sentence modelling is very necessary in dialogue generation, however, most models use many other methods such as attention mechanism in addition to sentence modelling. Therefore, we will put these models in the following chapter to describe.

3.1.3 dialogue representation

On basis of various representations of sentences, many researchers have proposed dialog representations for better context modeling. Past research in distributional semantics has suggested the meaningful of language can be inferred from the adjacent context([52], [53]). [54] put forward a simplest representation, consider the context as a single sentence which is used as input in the standard seq2seq encoder, and then use the decoder to generate reply. In order to represent context hierarchically, [10], [11] proposed HRED, a fundamental representation of dialog, which utilizes a RNN to model context information on top of two RNN as encoder and decoder respectively. [55] uses MLP as an encoder and

the remaining modules are basically the same as HRED. Instead of using one RNN to model context, [56] proposed two RNN to represent context information for speaker A and speaker B respectively with cross connection between adjacent turns. Furthermore, they complement external state modeled by another RNN to explicitly abstract dialog context information. For the contextual representation, [57] proposed a novel method. They use two sets of HRED. The first set uses the dialog corpus as the input, and the second set uses the coarse representation of the dialog corpus as an input. This kind of coarse representation is a dominant representation for the nouns and verbs of the original data.

Just as section 3.1.1 saying, [43] and many of the above models greatly improve the model's relevance. Intuitively, when the model is able to model the context, it can to a certain extent understand the essence of the entire conversation, so it can generate highly relevant responses.

3.1.4 summary

Table 1 summarizes the applicable scope, advantages, weakness and representative models of the sentence-level and dialogue-level representations. As can be seen from the table, the two representations each have their own advantages and disadvantages and application scope, the key is to see what characteristics of the data you want to process.

3.2 Horizontal representation

3.2.1 meta data

Researchers have introduced multiple meta data between word-level and sentence-level representation. [58] simply splices word embedding and bag-of-word representation.

emotion [59] combines the traditional word embedding with an emotion representation (VAD) of word, and they splice these after the word embedding. In the loss function, they added the Euclidean distance between the VAD representation of all the words of the message and the VAD representation of all the words of the response, thus ensuring that the emotion of the response is close to the message. In the beam search phase, they use cos similarity to measure different responses to ensure emotional differences, thereby increasing the diversity of emotion. [24] utilize emotion to get the correlation scores.

domain or topic Unlike other models, [60] introduced the domain to assist in the generation of dialog, they first use different domain datasets to train multiple standard seq2seq, and at the same time train a classifier that combines RNN and SVM to estimate the domain of utterance, and finally select the response with the highest score through a Re-ranker module. Unlike [25] and [26], [27] proposes a novel method to introduce topic. They first sample the keyword from the Query using the PMI method, then use the seq2seq model to reverse-regenerate the previous word of the keyword, and finally reverse the generated words and use the other seq2seq model to generate the words after the keyword. [61] uses the profile detector module to derive a speaker keyword based on the output of the encoder, then generates other words based on this keyword backward and forward respectively. By adding *domain or topic*, models can improve the relevance between query or context and response.

Approaches	Applicable Scope	Advantages	Weakness	Representative models
sentence representation	This kind of method focuses on answering the question itself, without extra information.	It avoids redundant information and focuses on the question itself.	It lost a lot of useful information while avoiding redundant information.	[27]
dialogue representation	This type of approach focuses on conversations that require contextual review. The entire conversation process often has certain goals and topics.	It can model conversations and capture a lot of useful information.	The information it captures contains a lot of useless information that dilutes the value of useful information.	[10], [11]

TABLE 1: Summary of hierarchical representation.

speaker role Taking a dialogue speaker role into consideration is also a very good method. [28] utilizes speaker-level representation as one of the inputs to decoder. Furthermore, just as in the TV series Friends, Ross showed different ways of speaking to Monica and Rachel, which means that a person’s way of speaking depends not only on their own but also on the addressee. Therefore, they linearly combine the two vectors that represent the speaker and act as input to the decoder. It is noteworthy that the authors have changed the loss function, added the likelihood of generating message given the response and the length of the generated sentence, which has achieved good results. [62], like the previous model, also introduces the speaker role into the model. However, they combine the speaker role with HRED. Unlike other models, [29] proposed modelling multiple speakers to solve the problem of multi-party conversations. Specifically, they regard all utterances as a single sentence and add a vector representing the speaker at the beginning of each utterance. Furthermore, they proposed a dynamic model, which means that the vector of the speaker will change with time at each timestep. By modelling *speaker role*, the model learns how to respond in a specific role, thus increasing diversity.

3.2.2 attention mechanism

The attention mechanism([30], [31]) has always occupied a very important position in deep learning applications since it was put forward. Researchers also have some improvements to attention itself. [63] uses *Diag-disabled mask*, *Forward mask*, *Backward mask* respectively to perform *masked self-attention mechanism*. Furthermore, [64] perform *Masked block self-attention (mBloSA) mechanism* with *Intra-block self-attention* and *Inter-block self-attention*. Instead of modifying attention mechanism itself, [65] put forward *Key-value separation*, *Key-value-predict separation* and *Concatenation of previous output representations* to segment hidden state at each time step. [66] regards the \tilde{h} obtained by attention mechanism as a memory state, then makes a linear change with it and the current hidden state to get the final h_t .

In the field of dialogue, researchers have deepened computer understanding of representation by introducing attention mechanism into different levels of representations. Based on [54], [67] complement attention mechanism to the standard seq2seq. Unlike [54] and [67], [68] uses the last hidden state of each sentence as a representation of the sentence, then performs sum pooling and concatenation on these representations, and further uses another RNN

to abstract the context information. This RNN takes as inputs all these sentences representation. Finally, they added attention mechanism to give different weights to different sentences. [12] proposed HRAN, a developmental HRED with sentence-level attention and dialogue-level attention, which extended the traditional attention mechanism to the utterance level. [25] performs attention processing on messages and topic words, and puts the results together in the decoder to generate a response. [26] used a different kind of attention method to process message and topic words, and also generated a Topic Summarizer in the encoder to represent topic words more abstractly. The previous attention mechanism is based on message or context. [69] proposed a different approach. First, the decoder generates only one segment of the entire sentence at a time, and then adds this segment to the tail of the message to re-encode and introduce the attention mechanism.

3.2.3 gate mechanism

Researchers have migrated structures in the GRU and LSTM to the field of dialogue, demonstrating the power of the gate. [13] proposed a Internal Memory controlled by a read gate and write gate for reading and writing internal memory respectively. They assume that the internal memory maintains a emotion state and use sigmoid function to process the input and output of the GRU respectively to get read gate and write gate. Unlike the traditional gate mechanism, [32] propose the structure of *encoder-diverter-decoder model* to exactly model response mechanism. Specifically, they define m_i to represent latent mechanism for response generation as first, then, given the context c , the output of encoder, they calculate probability of each latent mechanism m_i conditioned on x . the inputs:

$$p(m_i|x) = \frac{\exp g(m_i, c)}{\sum_{k=1}^M \exp g(m_k, c)} \quad (1)$$

To avoid overfitting g is defined with the maxout activation function([70]):

$$g(m_i, c) = m_i^T W_t t \quad (2)$$

$$t = [\max\{\tilde{t}_{2j-1}, \tilde{t}_{2j}\}]_{j=1,2,\dots,l_c}^T \quad (3)$$

$$\tilde{t} = W_c c \quad (4)$$

where \tilde{t}_j is the j -th element of the vector \tilde{t} . Finally, they feed the concatenation of $[c; m_i]$ to the decoder to model $p(y|mi, x)$ and generate response. Based on this *encoder-diverter-decoder model*, [71] put forward the *encoder-diverter-filter-decoder model*. This model introduce Reinforcement learning (RL) to select mechanism. Specifically, the action contains a termination action a_t and continue action a_c ; the state is the set of selected mechanisms, $S_x^{(k)} \in S$; and the policy is defined as:

$$\pi(a_t|S_x^{(k)}) = \frac{p(m_0|x)}{1 - \sum_{m \in S_x^{(k)}} p(m|x)} \quad (5)$$

$$\pi(a_c|S_x^{(k)}) = 1 - \pi(a_t|S_x^{(k)}) \quad (6)$$

also the reward is defined as:

$$r_K = \log p(y|x; S_x^{(k)}) = \log \frac{\sum_{m \in S_x^{(k)}} p(m|x)p(y|m, x)}{\sum_{m \in S_x^{(k)}} p(m|x)} \quad (7)$$

$$r_k = 0 \quad (k = 1, 2, \dots, K - 1) \quad (8)$$

3.2.4 summary

Table 2 summarizes the applicable scope, advantages, weakness and representative models of the metadata, attention mechanism and gate mechanism. The first three lines do not give Weakness because their weaknesses lie in ignoring other information. As can be seen from the table, metadata is more similar to the introduction of external information, while attention mechanism and gate mechanism tend to reorganize existing information according to its importance.

3.3 external representation

There is still a large amount of data available outside of the dialog corpus. Therefore, how to efficiently extract these data becomes a problem that needs to be solved urgently. [43] uses textcnn to extract relevant external knowledge and integrates the extracted local features into each hidden state in generation. Instead of utilizing CNN to extract external information, [72] processes the content derived from TF-IDF with another RNN called facts encoder. [73] regard commonsense knowledge as triples, $\langle \text{concept } 1, \text{relation}, \text{concept } 2 \rangle$, then they utilize *Tri-LSTM Encoder* to represent these external knowledge and calculate the bilinear results with reply, at the same time, they use *Dual-LSTM Encoder* to encode message and reply. Finally, the compatibility of message and reply is defined as:

$$f(x, y) = (\vec{x} + \vec{o})^T \vec{y} = \vec{x}^T \vec{y} + \left(\sum_i p_i \vec{a}_i \right)^T \vec{y} \quad (9)$$

where x and y is message and reply respectively, \vec{a}_i is representation of each commonsense knowledge and p_i is the attention signal over \vec{a}_i calculated by:

$$p_i = \text{softmax}(\vec{x}^T \vec{a}_i) \quad (10)$$

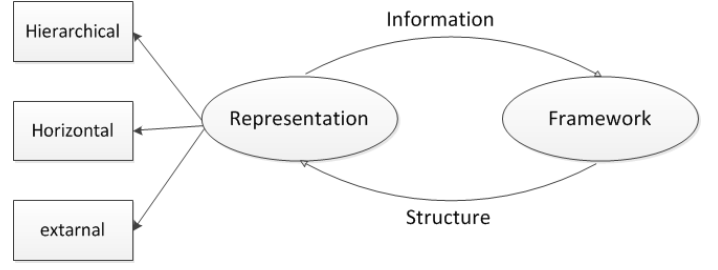


Fig. 5: The architecture of representation.

3.4 summary

Table 3 summarizes the applicable scope, advantages, weakness and representative models of the hierarchical representations, horizontal representations and external representations. This comprehensive comparison will show the meaning and effect of each representation. Figure 5 shows the architecture of this section.

4 FRAMEWORK

As mentioned earlier, the seq2seq model is a basic model, almost all models are based on this model. However, for reasons of solving the *universal reply* problem, the model also needs to be improved from the framework perspective. Section 4.1, 4.2, 4.3 and 4.4 will introduce VAE, GAN, LM and RL respectively and we will describe some improvement methods in section 4.5 then give a summary in section 4.6.

4.1 VAE

[35], [36] proposed VAE, which hypothesizes that prior $p(z)$ is normal distribution. During the training process, on the one hand they reduce the reconstruction error, and on the other hand they make the posterior $q(z|x)$ to approximate the prior $p(z)$ by KL divergence.

[74] first introduce vae to text generation, they put the encoder at one end of the vae, and let the latent variable z learn the representation of the sentence in the encoder, after calculating the mean and variance and using the reparameterization trick, the output is connected to the decoder and let the decoder predict. However, vae for text generation has encountered a problem called *vanishing latent variable problem*, the researchers further proposed two methods: *KL cost annealing* and *Word dropout and historyless decoding*. The former reduces the weight of the KL term in the loss function at the beginning of the training, thus enabling the model to encode as much information in z as possible, as the training progresses, the weights gradually increase until it is 1. The latter is deliberately changing part of the input word to unk during the decoder phase so that the model tries to obtain information through z as much as possible.

Instead of utilizing RNN to represent sentence in VAE framework, [75] uses CNN to replace RNN, convolutes the sentence and passes it to VAE, and deconvolves the output of VAE. They connected the deconvolution operation in the decoder to RNN and ByteNet, respectively, and the outputs was taken as the final predictions. Specifically, [76] uses the VAE framework to generate dialog, they use bidirectional

<i>Approaches</i>	<i>Applicable Scope</i>	<i>Advantages</i>	<i>Weakness</i>	<i>Representative models</i>
metadata-emotion	It applies to the emotionally rich corpus so that it produces a more natural response.	It helps to increase the emotional diversity of the response and express the emotion that groudtruth wants to express.	—	[59], [24]
metadata-domain or topic	It applies to corpora with a clear intention of response or dialogue	It allows computer-generated answers to express the domain or topic to some extent.	—	[25], [26] and [27]
metadata-speaker role	It focuses on the dialogue between two characters or multiple characters.	It can model the speaker roles and to generate responses that matches the speaker roles.	—	[28], [29]
attention machanism	It applies to scenes that require information to be weighted.	It assigns high weights to important information and low weights to minor and unimportant information. This facilitates the selection of important information from redundant information.	It increases the complexity of the models.	[12]
gate machanism	It is suitable when modeling long distance information.	It can maintain a block of information, read from it and write it continuously.	It also increases the complexity of the model, and it cannot be calculated in parallel, so it takes a long time to train.	[13], [32]

TABLE 2: Summary of horizontal representation.

<i>Approaches</i>	<i>Applicable Scope</i>	<i>Advantages</i>	<i>Weakness</i>	<i>Representative models</i>
hierarchical representation	It is suitable for mining the text itself.	It can deeply mine the text itself and extract different levels of information by modeling different levels.	Only part of the information it taps out is valid information, and most of the rest is useless.	[10], [11]
horizontal representation	It is suitable for filtering information.	It can reorganize existing information and improve the effectiveness of information.	It increases the complexity of the model and leads to longer training times.	[12], [13]
external representation	It applies to the representation of increased information.	It introduces information that does not exist in the corpus and improves the diversity of the model.	It adds redundant information and the use of these redundant information is less valuable than the corpus.	[73]

TABLE 3: Summary of representation.

RNN to represent message and reply then takes it as input for VAE.

In order to incorporate contextual information into the model, [14] adds the context representation to the message representation and uses the result as input to VAE. In addition, the model also added lots of meta data, such as conversation floor and dialog act. They also proposed a novel method , bag-of-word loss, to solve the vanishing latent variable problem, which adds the probability of generating the bag-of-word representation of x given latent variable z and context c .

Furthermore, [15] proposes SPHRED representing utterance-level dialog for each speakers and uses this instead of sentence -level representation to be inputs to VAE. If we do not consider interlocutors, all contextual information is represented by a identical structure, [42] offers a solution.

Similar to HRED, they use another RNN to model the context in addition to the encoder and decoder. Each context hidden state are used as one of the inputs to the VAE, and the VAE output is passed sequentially to each timestep of the decoder.

4.2 GAN

At the same time as VAE, Goodfellow([37]) proposed another framework, GAN. GAN has a discriminator and a generator. The generator generates the target data. The discriminator judges the real data and the generated data at the same time and feeds back the corresponding signal to the generator to increase the generating ability of the generator.

[77] proposed seqGAN, which first introduce GAN to generate sequences of discrete tokens.

Specifically, they use RNN as a generator and textcnn as a discriminator. Because the discrete data is not derivable, the signal of the discriminator cannot be passed to the generator. They propose to use the strategy of reinforcement learning to pass the discriminator signal as a reward to the generator. Therefore, the objective of the generator model(policy) $G_\theta(y_t|Y_{1:t-1})$ is to generate a sequence from the start state s_0 to maximize its expected end reward:

$$J(\theta) = E[R_T|s_0, \theta] = \sum_{y_1 \in Y} G_\theta(y_1|s_0) \cdot Q_{D_\phi}^{G_\theta}(s_0, y_1) \quad (11)$$

where

$$\{Y_{1:T}^1, \dots, Y_{1:T}^N\} = MC^{G_\beta}(Y_{1:T}; N) \quad (12)$$

and

$$Q_{D_\phi}^{G_\theta}(s = Y_{1:t-1}, a = y_t) =$$

$$\begin{cases} \frac{1}{N} \sum_{n=1}^N D_\phi(Y_{1:T}^n), & Y_{1:T}^n \in MC^{G_\beta}(Y_{1:T}; N) \quad \text{for } t < T \\ D_\phi(Y_{1:t}) & \text{for } t = T \end{cases} \quad (13)$$

Similarly, Li Jiwei ([16]) uses the GAN framework and the same strategy as seqGAN to enhance the diversity of dialogue generation. [78] proposes a novel method for generating text using GAN called MaskGAN. On the one hand, the generator of this model takes as inputs a masked sequence with some token replaced with underscore and then uses encoder-decoder architecture to generate the filled-in sequence. On the other hand, the discriminator has an identical architecture to the generator except that the output is a scalar probability at each time point, rather than a distribution over the vocabulary size. Similarly, they also use RL to transmit signals.

4.3 LM

The language model is used to model a sentence and the probability of a sentence can be calculated. [17] uses neural networks to model sentences. Furthermore, [79] proposes the use of RNNs to model sentences so that they can hold as much effective information as possible. [80] further proves that the language model is easier to model dialogue than the seq2seq model and they add attention mechanism to the Recurrent Neural Network Language Model(RNNLM). ([38]) used RNNLM to model the conversation and added the topic obtained by LDA to the RNN's hidden state.

4.4 RL

Reinforcement learning (RL) is an area of machine learning inspired by behaviourist psychology, concerned with how software agents ought to take actions in an environment so as to maximize some notion of cumulative reward. Task in dialogue is formulated as a sequential decision problem: current decision (or action) of word generation affects following decisions, which can be naturally addressed by policy gradient method ([81]). In upstream applications, [82] proposed that structured sentences can be obtained by

RL to classify texts. In their model, they first propose the Information Distilled LSTM (ID-LSTM) model which decides whether to retain or delete each time he reads a word. This model is very similar to attention mechanism, extracting the important content of the sentence. Furthermore, they put forward the Hierarchically Structured LSTM (HS-LSTM) model. Instead of deciding to retain or delete, they judge which word is inside or end in this model. So intuitively speaking, they segmented the sentence by structure and get better results. In the filed of dialogue, the key point is how to set the right reward. [83] proposed three types of rewards. The first is *Ease of answering*, which rewards those easily answerable responses to prevent the other person from answering the generic response in the next round since they cannot answer. Specifically, they hopes that the probability of generation of universal reply s is as small as possible given a response a .

$$r_1 = -\frac{1}{N_S} \sum_{s \in S} \frac{1}{N_s} \log p_{seq2seq}(s|a) \quad (14)$$

where N_S denotes the cardinality of N_S and N_s denotes the number of tokens in the dull response s . Next, *Information Flow* rewards answers with more information, which achieved their goal by punishing similar answers between consecutive turns.

$$r_2 = -\log \cos(h_{p_i}, h_{p_{i+1}}) = -\log \frac{h_{p_i} \cdot h_{p_{i+1}}}{\|h_{p_i}\| \|h_{p_{i+1}}\|} \quad (15)$$

where h_{p_i} and $h_{p_{i+1}}$ denote representations obtained from the encoder for two consecutive turns p_i and p_{i+1} . Finally, *Semantic Coherence* measure the adequacy of responses to avoid situations in which the generated replies are highly rewarded but are ungrammatical or not coherent.

$$r_1 = \frac{1}{N_a} \log p_{seq2seq}(a|q_i, p_i) + \frac{1}{N_{q_i}} \log p_{seq2seq}^{backward}(q_i|a) \quad (16)$$

where $p_{seq2seq}(a|q_i, p_i)$ denotes the probability of generating response a given the previous dialogue utterances $[p_i, q_i]$. $p_{seq2seq}^{backward}(q_i|a)$ denotes the backward probability of generating the previous dialogue utterance q_i based on response a . The final reward for action a is a weighted sum of the rewards discussed above. One of this paper's authors, Jiwei Li, once said: The most important thing to apply RL to dialogue generation is the setting of reward and the reward must be able to measure the quality of the dialogue. However, this is the most difficult in dialogue, therefore we think of using adversarial training to generate rewards in [16].

4.5 improvement

Different frameworks can use some methods to improve the quality of generated sentences. This section describes these common methods.

4.5.1 beam search

The traditional Breadth First Strategy(BFS) can find the optimal path, but in the case of very large search space, the memory usage is exponentially increasing and it is easy to

cause memory overflow. Therefore, a beam search algorithm is proposed. Beam search attempts to optimize the search space on a basis of BFS to reduce memory consumption, which is similar to prune. [84] proposed to connect Beam Search algorithm after the seq2seq framework. They first set a beam width, then select the top beam width number highest score in the generated candidates at each time and continue to iterate until the end of the generation.

[85] is such a method. In their study of how to improve the beam search’s diversity, they subtract penalty from each sub branches in the corresponding order of different branches and the penalty increases with the order. So they can highlight the highest score in different branches and not just a few higher scores in the same branch.

4.5.2 loss function

The loss function directly defines the training objective of the model, so it is also a very important research direction to modify the loss function so that it meets the training results we want. [86] gives two changed objective function.

$$\hat{T} = \operatorname{argmax}_T \{ \log p(T|S) - \lambda \log p(T) \} \quad (17)$$

and

$$\hat{T} = \operatorname{argmax}_T \{ (1 - \lambda) \log p(T|S) + \lambda \log p(S|T) \} \quad (18)$$

however, there are some problems in the two loss functions. For this, the author gives the corresponding solutions in practical consideration. The former may result in ungrammatical responses. After observing, the author found that the first few words in the generated response are often very important, so they replaced $p(T)$ with the following formula.

$$U(T) = \prod_{i=1}^{L_t} p(t_i | t_1, t_2, \dots, t_I) \cdot g(i) \quad (19)$$

then they make the first $\gamma g(i)$ values 1 and the rest 0. The problem of the latter is intractable. To solve this problem, the author only uses the second item as a standard to sort the content generated by the standard loss.

4.5.3 dynamic vocabulary

In the previous generation model, when generating a word, the model needed to traverse the entire dictionary. If the dictionary is very large, then each traversal is very time consuming. Therefore, based on the traditional sequence-to-sequence with attention model, [87] put forward the DVS2S model whose *Dynamic Vocabulary Construction* component significantly improved the time consumption problem caused by the large vocabulary. Firstly, they divide words into function words which guarantee grammatical correctness and fluency of response and content words which express semantics of responses. For function words, $\beta_{i,j}$, the probability of being left, is always equal to 1, and for content words, $\beta_{i,j}$ is calculated as:

$$\beta_{I(c)} = \sigma(W_c^T h_t + b_c) \quad (20)$$

after removing some unnecessary words, the remaining words is utilized by the decoder as dictionary.

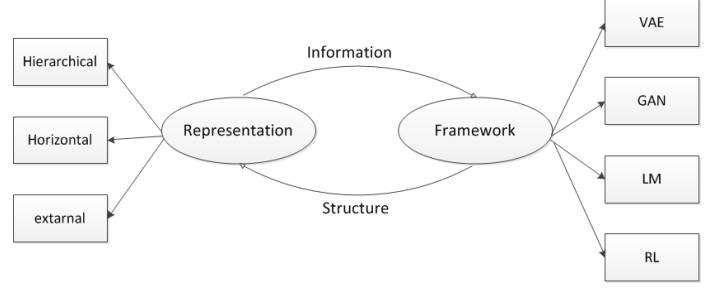


Fig. 6: The architecture of representation and framework in detail.

4.5.4 classified words

Similar to [87], [88] utilize the standard encoder-decoder model and change the projection layer of the decoder. The author finds that different types of words are always placed in the same vocabulary with the same weight, however, when human beings communicate with each others, it is not like this. Therefore they put forward two types of methods: STD and HTD.

Specifically, for each hidden state of decoder, STD first estimate a distribution for these types then generate different vocabulary distribution for these types and multiply these corresponding types distribution, finally, it mix all these vocabulary to get the final results. Instead of assuming that each word is a distribution over the word types implicitly, HTD explicitly divide the vocabulary into different types. This method utilizes *Gumbel-Softmax* method to obtain a differentiable surrogate to the *argmax* function then use it to multiplies the corresponding vocabulary to get the final results.

4.6 summary

Table 4 summarizes the advantages, weakness and representative models of the hierarchical representations, horizontal representations and external representations. This comprehensive comparison will show the meaning and effect of each representation. Figure 6 shows the architecture of the representation and framework. Note that these components are not simply stacked together, but they come together to work together.

5 EVALUATION

The previous chapters described various generation models. However, humans always need to evaluate the models given by computers. At present, it is still most reasonable for humans to evaluate these models themselves, that is, human evaluation. But this method is too expensive and too inefficient. Therefore, in order to evaluate the quality of different models automatically, researchers proposed a variety of metrics.

5.1 Word Overlap-based

The evaluation metrics in this section are basically learned from the field of machine translation. Such metrics do not fully reflect the goodness of a reply but can, to some extent, reflect the similarity to the original reply.

<i>Approaches</i>	<i>Advantages</i>	<i>Weakness</i>	<i>Representative models</i>
VAE	It helps to improve the consistency and diversity of the generated text by modeling latent variables.	This framework itself is not very compatible with the encoder-decoder architecture.	[14], [15]
GAN	It uses a discriminator to define and measure "good generated response" that is difficult to define clearly.	Its training process is unstable and the discriminator plays a much lower role than expected so that Teacher-Forcing is constantly introduced during the training process.	[16]
LM	It is able to model the language most accurately.	It is completely incompatible with the current mainstream, the encoder-decoder architecture..	[80]
RL	It can pass signals from one side to the other in two phases that isn't differentiable.	Its training process is unstable, and the values of signals passed are less than expected.	[83]

TABLE 4: Summary of framework.

5.1.1 BLEU

BLEU ([89]) calculates the joint occurrence rate of n-grams in the ground truth and the proposed responses. It first computes an n-gram precision for the whole dataset:

$$P_n(R, \hat{R}) = \frac{\sum_i \sum_k \min(h(k, r_i), h(k, \hat{r}_i))}{\sum_i \sum_k h(k, r_i)} \quad (21)$$

where k donates all possible n-grams of length n, $h(k, r_i)$ means the number of n-grams k in r_i . In order to bypass the drawbacks of utilizing a precision score, which means that it favours shorter sentences, the authors add a brevity penalty. BLEU-N where N is the maximum length of the n-grams considered is donated as:

$$BLEU - N := b(R, \hat{R}) \exp(\sum_{n=1}^N \beta_n \log P_n(R, \hat{R})) \quad (22)$$

β_n , a weighting, is usually uniform, and $b(\cdot)$ is the brevity penalty. Note that BLEU is usually calculated at the corpus-level, and has been shown to correlate with human judgement in the translation domain when there are multiple ground truth candidates available.

5.1.2 METEOR

The METEOR metric ([90]) was introduced to address several weaknesses in BLEU. It creates an explicit alignment between the candidate and target responses. The alignment is based on exact token matching, followed by WordNet synonyms, stemmed tokens, and then paraphrases. Given a set of alignments, the METEOR score is the harmonic mean of precision and recall between the proposed and ground truth sentence.

5.1.3 ROUGE

ROUGE ([91]) is a set of evaluation metrics used for automatic summarization. We consider ROUGE-L, which is a Fmeasure based on the Longest Common Subsequence (LCS) between a candidate and target sentence. The LCS is

a set of words which occur in two sentences in the same order; however, unlike n-grams the words do not have to be contiguous, i.e. there can be other words in between the words of the LCS.

5.2 Embedding-based

If we consider the above metrics as *Word Overlap-based Metrics*, then [92] put forward some novel metrics called *Embedding-based Metrics*.

5.2.1 Embedding Average

The first calculates the mean of the word embeddings of each token in a sentence called *Embedding Average* in the following:

$$\bar{e}_r = \frac{\sum_{w \in r} e_w}{|\sum_{w' \in r} e_{w'}|} \quad (23)$$

where $|\cdot|$ represents the modulus of the vector and we can compute the cosine similarity between their respective sentence level embeddings $EA := \cos(\bar{e}_r, \bar{e}_{\hat{r}})$.

5.2.2 Vector Extrema

Next, *Vector Extrema* take the most extreme value amongst all word vectors in the sentence for each dimension of the word vectors.

$$e_{rd} = \begin{cases} \max_{w \in r} e_{wd} & \text{if } e_{wd} > |\min_{w' \in r} e_{w'd}| \\ \min_{w' \in r} e_{w'd} & \text{otherwise} \end{cases}$$

where d indexes the dimensions of a vector; e_{wd} is the dth dimensions of e_w (ws embedding) and $|\cdot|$ calculated the absolute value. By taking the extrema along each dimension, they are therefore more likely to ignore common words.

5.2.3 Greedy Matching

Finally, they also proposed *Greedy Matching* which averaged across all words in two sentence r and \hat{r} :

$$G(r, \hat{r}) = \frac{\sum_{w \in r; \max_{\hat{w} \in \hat{r}} \cos(e_w, e_{\hat{w}})} |r|}{|r|} \quad (24)$$

$$GM(r, \hat{r}) = \frac{G(r, \hat{r}) + G(\hat{r}, r)}{2} \quad (25)$$

this greedy approach favours responses with key words that are semantically similar to those in the ground truth response.

5.3 others

5.3.1 perplexity

In information theory, perplexity is a measurement of how well a probability distribution or probability model predicts a sample. It may be used to compare probability models. A low perplexity indicates the probability distribution is good at predicting the sample. In NLP, perplexity is used to estimate the probability of a sentence appearing. Specifically:

$$\begin{aligned} PP(S) &= P(w_1 w_2 \dots w_N)^{-\frac{1}{N}} \\ &= \sqrt[N]{\frac{1}{p(w_1 w_2 \dots w_N)}} \\ &= \sqrt[N]{\prod_{i=1}^N \frac{1}{p(w_i | w_1 w_2 \dots w_{i-1})}} \quad (\text{chain rule}) \quad (26) \\ &= \sqrt[N]{\prod_{i=1}^N \frac{1}{p(w_i | w_{i-1})}} \quad (\text{for bigram}) \end{aligned}$$

5.3.2 Adversarial Evaluation

Instead of using *Embedding-based Metrics*, [93] utilizes a generator and a discriminator as same as a GAN, specifically, the generator is a sequence-to-sequence model to generate response and the discriminator is an RNN with only an encoder followed by a binary classifier. During adversarial training, the discriminator gradually improve its ability and can be used as metric.

5.3.3 RUBER

Different from the previous metrics, [94] first calculate the referenced score:

$$s_R(r, \hat{r}) = \cos(v_r, v_{\hat{r}}) = \frac{v_r^T v_{\hat{r}}}{\|v_r\| \cdot \|v_{\hat{r}}\|} \quad (27)$$

where $v = [v_{max}; v_{min}]$ and

$$v_{max}[i] = \max\{\omega_1[i], \omega_2[i], \dots, \omega_n[i]\} \quad (28)$$

the calculation of $v_{min}[i]$ is similar to $v_{max}[i]$. Via the referenced score, this metric evaluates the similarity between *Generated Reply* \hat{r} and *Groundtruth Reply* r . Furthermore, they utilize two different Bi-GRU to model *Query* and *Reply* respectively and then concatenate this two representation and their bilinear results to get the final score between *Query*

and *Reply*. Finally, they normalize each score to the range (0, 1) via:

$$\tilde{s} = \frac{s - \min(s')}{\max(s') - \min(s')} \quad (29)$$

and combine \tilde{s}_R and \tilde{s}_U as ultimate RUBER metric.

5.3.4 skip-thoughts

[95] uses a recurrent network to encode a given sentence into an embedding and then decoder it to predict the previous and next sentence. This method is trained in an unsupervised fashion on the BookCorpus dataset([96]). The embeddings produced by the encoder have a robust performance on semantic relatedness tasks.

5.4 Human evaluation

Many papers use human evaluation in addition to the above methods. However, this most precise metric is too expensive to utilize it widely. [97] put forward to use human evaluation via crowdsourcing. Crowdsourcing is a sourcing model in which individuals or organizations obtain goods and services.

5.5 summary

[98] implements an open source library that contains most of the above criteria and facilitates the open-domain dialogue generation system researchers.

Table 5 summarizes the advantages, weakness and representative models of the hierarchical representations, horizontal representations and external representations. The current situation is still mixed using several methods.

6 CONCLUSION AND FUTURE DIRECTION

6.1 Conclusion

With the sequence-to-sequence model applied to the open-domain dialogue system, chat-bot do have a certain degree of understanding, but the problem of the universal answer such as *I don't know* and *I'm not sure* haunts all researchers. There are two main directions for solving this problem: relevance and diversity, researchers have proposed various methods to achieve it. Only by combing the previous methods can we better understand the problems and challenges faced by dialogue generation so as to better improve the dialogue system. Therefore, this paper will classify these methods according to the representation and framework and then give the traditional and advanced evaluation for the open-domain dialogue system.

From the respective of representation, we can find out a hierarchical structure and a horizontal architecture. In the hierarchical structure, researchers express the words as word embedding and obtain sentence representation and dialogue representation further. By modelling the utterances of dialogue, dialogue system can understand the context information to generate more meaningful and relevance reply. And in the horizontal architecture, attention mechanism and gate mechanism are placed on various level of hierarchical representation, also meta data such as dialogue emotion, domain or topic and speaker role are introduced to these

<i>Approaches</i>	<i>Advantages</i>	<i>Weakness</i>	<i>Representative metrics</i>
Word Overlap-based	This evaluation standard is learned from the field of machine translation and can faithfully reflect the difference between computer-generated responses and groundtruth.	Different people give different answers to the same question, so it cannot give correct scores for these responses which are correctly generated but have different angles and ways from groundtruth.	BLEU, ROUGE
Embedding-based	It evaluates the sentence generated based on the word embedding, to a certain extent, according to the sentence meaning rather than the similarity with groundtruth.	It is not yet perfect and does not take into account the consistency of the sentences.	Greedy Matching
Human evaluation	It is the standard that best meets the rules of human dialogue.	It is too expensive and it is difficult to expand on a large scale.	—

TABLE 5: Summary of evaluation.

level of hierarchical representation to improve diversity of dialogue generation. Besides, external representation for example commonsense knowledge and so on can also make the response more diverse.

Sequence-to-sequence model is the basic framework for the open-domain dialogue generation, however in order to avoid the universal response problem, VAE, GAN, LM and RL are introduced to dialogue system. At present, VAE and RL are mainstream. All this advanced framework don not discard the sequence-to-sequence architecture, they only modify it by introducing other methods. Also, some additional methods that perfect these framework are introduced in the *improvement* section.

Finally, we introduce some evaluation metrics to evaluate these methods above such as BLEU, ROUGE and so on.

6.2 future direction

Open-domain dialogue generation system can be extended in several ways. From a methodological perspective, this includes (1)digging deeper into the text itself, discovering features that have not previously been valued, and using models to model it, (2)optimizing model structure to improve modelling performance. From a technical perspective, this may include (1)combining multiple representations, (2)exploring the new frameworks. The variety of models is only tools. The purpose of using these tools should be the goal that researchers must strive for.

6.2.1 text

The various texts themselves contain rich information. The more effective information is extracted, the higher the quality of response generated. According to the development of open-domain dialogue generation, there are two general directions that can be advanced.

Information representation. With the development of the World Wide Web(WWW), the Internet is flooded with a large number of dialogue corpus and other related texts.

We need to express this information hierarchically and consider whether we can propose a new presentation layer to increase the abstraction of the presentation. For example, from [27] to [10] and [11], [10] and [11] raise the dialogue-level representation, furthermore, [28] proposes speaker-level representation. These representations all increase the level of abstraction of information and increase the relevance of responses. At the same time, we can also study how to accurately position efficient external knowledge to supplement the training corpora and propose more meta-data that can improve the response.

Information recombination. Although there is a lot of information, how to use it efficiently is still a problem. There are many kinds of attention mechanisms have been studied, such as [99], [100], [101] and [49]. Whether or not the attention mechanism behind a variety of models can effectively improve the outcome of dialogues is a topic worthy of study. Whether or not it is possible to develop a more efficient gate mechanism is also a direction.

6.2.2 structure

From [33] and [34] to [35] and [37], the framework itself has also made tremendous progress. Especially for GAN, there are many other forms of exploration of the framework structure. For example, the GAN of multiple discriminators([102], [103], [104] and [105]) are worth exploring.

ACKNOWLEDGMENT

We would like to thank the colleague of Engineering Research Center of Ministry of Education for Advanced Computer Application Technology of Beihang University(BUAA).

REFERENCES

- [1] X. Li, Y. Chen, L. Li, J. Gao, and A. Çelikyilmaz, "End-to-end task-completion neural dialogue systems," in *Proceedings of the Eighth International Joint Conference on Natural Language Processing, IJCNLP 2017, Taipei, Taiwan, November 27 - December 1, 2017 - Volume 1: Long Papers*, 2017, pp. 733–743. [Online]. Available: <https://aclanthology.info/papers/I17-1074/i17-1074>

- [2] L. E. Asri, H. Schulz, S. Sharma, J. Zumer, J. Harris, E. Fine, R. Mehrotra, and K. Suleman, "Frames: a corpus for adding memory to goal-oriented dialogue systems," in *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue, Saarbrücken, Germany, August 15-17, 2017*, 2017, pp. 207–219. [Online]. Available: <https://aclanthology.info/papers/W17-5526/w17-5526>
- [3] B. Peng, X. Li, L. Li, J. Gao, A. Çelikyilmaz, S. Lee, and K. Wong, "Composite task-completion dialogue policy learning via hierarchical deep reinforcement learning," in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, 2017, pp. 2231–2240. [Online]. Available: <https://aclanthology.info/papers/D17-1237/d17-1237>
- [4] J. D. Williams, K. Asadi, and G. Zweig, "Hybrid code networks: practical and efficient end-to-end dialog control with supervised and reinforcement learning," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, 2017, pp. 665–677. [Online]. Available: <https://doi.org/10.18653/v1/P17-1062>
- [5] M. Eric, L. Krishnan, F. Charette, and C. D. Manning, "Key-value retrieval networks for task-oriented dialogue," in *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue, Saarbrücken, Germany, August 15-17, 2017*, 2017, pp. 37–49. [Online]. Available: <https://aclanthology.info/papers/W17-5506/w17-5506>
- [6] A. Bordes and J. Weston, "Learning end-to-end goal-oriented dialog," *CoRR*, vol. abs/1605.07683, 2016. [Online]. Available: <http://arxiv.org/abs/1605.07683>
- [7] P. Su, P. Budzianowski, S. Ultes, M. Gasic, and S. J. Young, "Sample-efficient actor-critic reinforcement learning with supervised data for dialogue management," in *Proceedings of the 18th Annual SIGdial Meeting on Discourse and Dialogue, Saarbrücken, Germany, August 15-17, 2017*, 2017, pp. 147–157. [Online]. Available: <https://aclanthology.info/papers/W17-5518/w17-5518>
- [8] T. Zhao and M. Eskénazi, "Towards end-to-end learning for dialog state tracking and management using deep reinforcement learning," in *Proceedings of the SIGDIAL 2016 Conference, The 17th Annual Meeting of the Special Interest Group on Discourse and Dialogue, 13-15 September 2016, Los Angeles, CA, USA, 2016*, pp. 1–10. [Online]. Available: <http://aclweb.org/anthology/W/W16/W16-3601.pdf>
- [9] B. Dhingra, L. Li, X. Li, J. Gao, Y. Chen, F. Ahmed, and L. Deng, "Towards end-to-end reinforcement learning of dialogue agents for information access," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, 2017, pp. 484–495. [Online]. Available: <https://doi.org/10.18653/v1/P17-1045>
- [10] A. Sordoni, Y. Bengio, H. Vahabi, C. Lioma, J. G. Simonsen, and J. Nie, "A hierarchical recurrent encoder-decoder for generative context-aware query suggestion," in *Proceedings of the 24th ACM International Conference on Information and Knowledge Management, CIKM 2015, Melbourne, VIC, Australia, October 19 - 23, 2015*, 2015, pp. 553–562. [Online]. Available: <http://doi.acm.org/10.1145/2806416.2806493>
- [11] I. V. Serban, A. Sordoni, Y. Bengio, A. C. Courville, and J. Pineau, "Building end-to-end dialogue systems using generative hierarchical neural network models," in *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA, 2016*, pp. 3776–3784. [Online]. Available: <http://www.aaai.org/ocs/index.php/AAAI/AAAI16/paper/view/11957>
- [12] C. Xing, W. Wu, Y. Wu, M. Zhou, Y. Huang, and W. Ma, "Hierarchical recurrent attention network for response generation," *CoRR*, vol. abs/1701.07149, 2017. [Online]. Available: <http://arxiv.org/abs/1701.07149>
- [13] H. Zhou, M. Huang, T. Zhang, X. Zhu, and B. Liu, "Emotional chatting machine: Emotional conversation generation with internal and external memory," *CoRR*, vol. abs/1704.01074, 2017. [Online]. Available: <http://arxiv.org/abs/1704.01074>
- [14] T. Zhao, R. Zhao, and M. Eskénazi, "Learning discourse-level diversity for neural dialog models using conditional variational autoencoders," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 1: Long Papers*, 2017, pp. 654–664. [Online]. Available: <https://doi.org/10.18653/v1/P17-1061>
- [15] X. Shen, H. Su, Y. Li, W. Li, S. Niu, Y. Zhao, A. Aizawa, and G. Long, "A conditional variational framework for dialog generation," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 2: Short Papers*, 2017, pp. 504–509. [Online]. Available: <https://doi.org/10.18653/v1/P17-2080>
- [16] J. Li, W. Monroe, T. Shi, S. Jean, A. Ritter, and D. Jurafsky, "Adversarial learning for neural dialogue generation," in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, 2017, pp. 2157–2169. [Online]. Available: <https://aclanthology.info/papers/D17-1230/d17-1230>
- [17] Y. Bengio, R. Ducharme, and P. Vincent, "A neural probabilistic language model," in *Advances in Neural Information Processing Systems 13, Papers from Neural Information Processing Systems (NIPS) 2000, Denver, CO, USA, 2000*, pp. 932–938. [Online]. Available: <http://papers.nips.cc/paper/1839-a-neural-probabilistic-language-model>
- [18] T. Mikolov, K. Chen, G. Corrado, and J. Dean, "Efficient estimation of word representations in vector space," *CoRR*, vol. abs/1301.3781, 2013. [Online]. Available: <http://arxiv.org/abs/1301.3781>
- [19] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Advances in Neural Information Processing Systems 26: 27th Annual Conference on Neural Information Processing Systems 2013. Proceedings of a meeting held December 5-8, 2013, Lake Tahoe, Nevada, United States., 2013*, pp. 3111–3119. [Online]. Available: <http://papers.nips.cc/paper/5021-distributed-representations-of-words-and-phrases-and-their-compositionality>
- [20] J. Pennington, R. Socher, and C. D. Manning, "Glove: Global vectors for word representation," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL, 2014*, pp. 1532–1543. [Online]. Available: <http://aclweb.org/anthology/D/D14/D14-1162.pdf>
- [21] Y. Bengio, "Deep learning of representations: Looking forward," in *Statistical Language and Speech Processing - First International Conference, SLSP 2013, Tarragona, Spain, July 29-31, 2013. Proceedings, 2013*, pp. 1–37. [Online]. Available: <https://doi.org/10.1007/978-3-642-39593-2-1>
- [22] G. E. Hinton and J. L. McClelland, "Learning representations by recirculation," in *Neural Information Processing Systems, Denver, Colorado, USA, 1987*, 1987, pp. 358–366. [Online]. Available: <http://papers.nips.cc/paper/78-learning-representations-by-recirculation>
- [23] Y. LeCun, B. E. Boser, J. S. Denker, D. Henderson, R. E. Howard, W. E. Hubbard, and L. D. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural Computation*, vol. 1, no. 4, pp. 541–551, 1989. [Online]. Available: <https://doi.org/10.1162/neco.1989.1.4.541>
- [24] A. C. Samson, S. D. Kreibig, B. Soderstrom, A. A. Wade, and J. J. Gross, "Eliciting positive, negative and mixed emotional states: A film library for affective scientists," *Cognition & Emotion*, vol. 30, no. 5, pp. 827–856, 2016.
- [25] C. Xing, W. Wu, Y. Wu, J. Liu, Y. Huang, M. Zhou, and W. Ma, "Topic aware neural response generation," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA, 2017*, pp. 3351–3357. [Online]. Available: <http://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/14563>
- [26] —, "Topic augmented neural response generation with a joint attention mechanism," *CoRR*, vol. abs/1606.08340, 2016. [Online]. Available: <http://arxiv.org/abs/1606.08340>
- [27] L. Mou, Y. Song, R. Yan, G. Li, L. Zhang, and Z. Jin, "Sequence to backward and forward sequences: A content-introducing approach to generative short-text conversation," in *COLING 2016, 26th International Conference on Computational Linguistics, Proceedings of the Conference: Technical Papers, December 11-16, 2016, Osaka, Japan, 2016*, pp. 3349–3358. [Online]. Available: <http://aclweb.org/anthology/C/C16/C16-1316.pdf>
- [28] J. Li, M. Galley, C. Brockett, G. P. Spithourakis, J. Gao, and W. B. Dolan, "A persona-based neural conversation model,"

- in *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*, 2016. [Online]. Available: <http://aclweb.org/anthology/P/P16/P16-1094.pdf>
- [29] H. Ouchi and Y. Tsuboi, "Addressee and response selection for multi-party conversation," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, 2016, pp. 2133–2143. [Online]. Available: <http://aclweb.org/anthology/D/D16/D16-1231.pdf>
- [30] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," *CoRR*, vol. abs/1409.0473, 2014. [Online]. Available: <http://arxiv.org/abs/1409.0473>
- [31] T. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," in *Proceedings of the 2015 Conference on Empirical Methods in Natural Language Processing, EMNLP 2015, Lisbon, Portugal, September 17-21, 2015*, 2015, pp. 1412–1421. [Online]. Available: <http://aclweb.org/anthology/D/D15/D15-1166.pdf>
- [32] G. Zhou, P. Luo, R. Cao, F. Lin, B. Chen, and Q. He, "Mechanism-aware neural machine for dialogue response generation," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*, 2017, pp. 3400–3407. [Online]. Available: <http://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/14471>
- [33] K. Cho, B. van Merriënboer, Ç. Gülçehre, D. Bahdanau, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL, 2014*, pp. 1724–1734. [Online]. Available: <http://aclweb.org/anthology/D/D14/D14-1179.pdf>
- [34] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Advances in Neural Information Processing Systems 27: Annual Conference on Neural Information Processing Systems 2014, December 8-13 2014, Montreal, Quebec, Canada*, 2014, pp. 3104–3112. [Online]. Available: <http://papers.nips.cc/paper/5346-sequence-to-sequence-learning-with-neural-networks>
- [35] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *CoRR*, vol. abs/1312.6114, 2013. [Online]. Available: <http://arxiv.org/abs/1312.6114>
- [36] D. J. Rezende, S. Mohamed, and D. Wierstra, "Stochastic backpropagation and approximate inference in deep generative models," in *Proceedings of the 31th International Conference on Machine Learning, ICML 2014, Beijing, China, 21-26 June 2014*, 2014, pp. 1278–1286. [Online]. Available: <http://jmlr.org/proceedings/papers/v32/rezende14.html>
- [37] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. C. Courville, and Y. Bengio, "Generative adversarial networks," *CoRR*, vol. abs/1406.2661, 2014. [Online]. Available: <http://arxiv.org/abs/1406.2661>
- [38] Y. Luan, Y. Ji, and M. Ostendorf, "LSTM based conversation models," *CoRR*, vol. abs/1603.09457, 2016. [Online]. Available: <http://arxiv.org/abs/1603.09457>
- [39] H. Chen, X. Liu, D. Yin, and J. Tang, "A survey on dialogue systems: Recent advances and new frontiers," *SIGKDD Explorations*, vol. 19, no. 2, pp. 25–35, 2017. [Online]. Available: <http://doi.acm.org/10.1145/3166054.3166058>
- [40] S. Mallios and N. G. Bourbakis, "A survey on human machine dialogue systems," in *7th International Conference on Information, Intelligence, Systems & Applications, IISA 2016, Chalkidiki, Greece, July 13-15, 2016*, 2016, pp. 1–7. [Online]. Available: <https://doi.org/10.1109/IISA.2016.7785371>
- [41] G. E. Churcher, E. S. Atwell, and C. Souter, "Dialogue management systems: a survey and overview," 1997.
- [42] I. V. Serban, A. Sordoni, R. Lowe, L. Charlin, J. Pineau, A. C. Courville, and Y. Bengio, "A hierarchical latent variable encoder-decoder model for generating dialogues," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*, 2017, pp. 3295–3301. [Online]. Available: <http://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/14567>
- [43] A. Sordoni, M. Galley, M. Auli, C. Brockett, Y. Ji, M. Mitchell, J. Nie, J. Gao, and B. Dolan, "A neural network approach to context-sensitive generation of conversational responses," in *NAACL HLT 2015, The 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Denver, Colorado, USA, May 31 - June 5, 2015*, 2015, pp. 196–205. [Online]. Available: <http://aclweb.org/anthology/N/N15/N15-1020.pdf>
- [44] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Computation*, vol. 9, no. 8, pp. 1735–1780, 1997. [Online]. Available: <https://doi.org/10.1162/neco.1997.9.8.1735>
- [45] J. Chung, Ç. Gülçehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *CoRR*, vol. abs/1412.3555, 2014. [Online]. Available: <http://arxiv.org/abs/1412.3555>
- [46] Y. Kim, "Convolutional neural networks for sentence classification," in *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing, EMNLP 2014, October 25-29, 2014, Doha, Qatar, A meeting of SIGDAT, a Special Interest Group of the ACL, 2014*, pp. 1746–1751. [Online]. Available: <http://aclweb.org/anthology/D/D14/D14-1181.pdf>
- [47] S. Lai, L. Xu, K. Liu, and J. Zhao, "Recurrent convolutional neural networks for text classification," in *Proceedings of the Twenty-Ninth AAAI Conference on Artificial Intelligence, January 25-30, 2015, Austin, Texas, USA*, 2015, pp. 2267–2273. [Online]. Available: <http://www.aaai.org/ocs/index.php/AAAI/AAAI15/paper/view/9745>
- [48] L. Logeswaran and H. Lee, "An efficient framework for learning sentence representations," in *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=rJvJXZb0W>
- [49] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, 2017, pp. 6000–6010. [Online]. Available: <http://papers.nips.cc/paper/7181-attention-is-all-you-need>
- [50] M. Iyyer, V. Manjunatha, J. L. Boyd-Graber, and H. D. III, "Deep unordered composition rivals syntactic methods for text classification," in *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL 2015, July 26-31, 2015, Beijing, China, Volume 1: Long Papers*, 2015, pp. 1681–1691. [Online]. Available: <http://aclweb.org/anthology/P/P15/P15-1162.pdf>
- [51] D. Cer, Y. Yang, S. Kong, N. Hua, N. Limtiaco, R. S. John, N. Constant, M. Guajardo-Cespedes, S. Yuan, C. Tar, Y. Sung, B. Strope, and R. Kurzweil, "Universal sentence encoder," *CoRR*, vol. abs/1803.11175, 2018. [Online]. Available: <http://arxiv.org/abs/1803.11175>
- [52] Z. S. Harris, "Distributional structure," 1954.
- [53] F. Hill, K. Cho, and A. Korhonen, "Learning distributed representations of sentences from unlabelled data," in *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, 2016, pp. 1367–1377. [Online]. Available: <http://aclweb.org/anthology/N/N16/N16-1162.pdf>
- [54] O. Vinyals and Q. V. Le, "A neural conversational model," *CoRR*, vol. abs/1506.05869, 2015. [Online]. Available: <http://arxiv.org/abs/1506.05869>
- [55] K. Yao, B. Peng, G. Zweig, and K. Wong, "An attentional neural conversation model with improved specificity," *CoRR*, vol. abs/1606.01292, 2016. [Online]. Available: <http://arxiv.org/abs/1606.01292>
- [56] B. Liu and I. Lane, "Dialog context language modeling with recurrent neural networks," in *2017 IEEE International Conference on Acoustics, Speech and Signal Processing, ICASSP 2017, New Orleans, LA, USA, March 5-9, 2017*, pp. 5715–5719. [Online]. Available: <https://doi.org/10.1109/ICASSP.2017.7953251>
- [57] I. V. Serban, T. Klinger, G. Tesauro, K. Talamadupula, B. Zhou, Y. Bengio, and A. C. Courville, "Multiresolution recurrent neural networks: An application to dialogue response generation," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*, 2017, pp. 3288–3294. [Online]. Available: <http://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/14571>

- [58] G. Forgues, J. Pineau, J.-M. Larchevêque, and R. Tremblay, "Bootstrapping dialog systems with word embeddings," 2014.
- [59] N. Asghar, P. Poupart, J. Hoey, X. Jiang, and L. Mou, "Affective neural response generation," in *Advances in Information Retrieval - 40th European Conference on IR Research, ECIR 2018, Grenoble, France, March 26-29, 2018, Proceedings*, 2018, pp. 154–166. [Online]. Available: <https://doi.org/10.1007/978-3-319-76941-7-12>
- [60] S. Choudhary, P. Srivastava, L. H. Ungar, and J. Sedoc, "Domain aware neural dialog system," *CoRR*, vol. abs/1708.00897, 2017. [Online]. Available: <http://arxiv.org/abs/1708.00897>
- [61] Q. Qian, M. Huang, H. Zhao, J. Xu, and X. Zhu, "Assigning personality/identity to a chatting machine for coherent conversation generation," *CoRR*, vol. abs/1706.02861, 2017. [Online]. Available: <http://arxiv.org/abs/1706.02861>
- [62] S. Kottur, X. Wang, and V. Carvalho, "Exploring personalized neural conversational models," in *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, 2017, pp. 3728–3734. [Online]. Available: <https://doi.org/10.24963/ijcai.2017/521>
- [63] T. Shen, T. Zhou, G. Long, J. Jiang, S. Pan, and C. Zhang, "Disan: Directional self-attention network for rnn/cnn-free language understanding," *CoRR*, vol. abs/1709.04696, 2017. [Online]. Available: <http://arxiv.org/abs/1709.04696>
- [64] T. Shen, T. Zhou, G. Long, J. Jiang, and C. Zhang, "Bi-directional block self-attention for fast and memory-efficient sequence modeling," in *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=H1cWzoxA->
- [65] M. Daniluk, T. Rocktäschel, J. Welbl, and S. Riedel, "Frustratingly short attention spans in neural language modeling," *CoRR*, vol. abs/1702.04521, 2017. [Online]. Available: <http://arxiv.org/abs/1702.04521>
- [66] D. Yogatama, Y. Miao, G. Melis, W. Ling, A. Kuncoro, C. Dyer, and P. Blunsom, "Memory architectures in recurrent neural network language models," in *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=SkFqf0IAZ>
- [67] L. Shang, Z. Lu, and H. Li, "Neural responding machine for short-text conversation," in *Proceedings of the 53rd Annual Meeting of the Association for Computational Linguistics and the 7th International Joint Conference on Natural Language Processing of the Asian Federation of Natural Language Processing, ACL 2015, July 26-31, 2015, Beijing, China, Volume 1: Long Papers*, 2015, pp. 1577–1586. [Online]. Available: <http://aclweb.org/anthology/P/P15/P15-1152.pdf>
- [68] Z. Tian, R. Yan, L. Mou, Y. Song, Y. Feng, and D. Zhao, "How to make context more useful? an empirical study on context-aware neural conversational models," in *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics, ACL 2017, Vancouver, Canada, July 30 - August 4, Volume 2: Short Papers*, 2017, pp. 231–236. [Online]. Available: <https://doi.org/10.18653/v1/P17-2036>
- [69] Y. Shao, S. Gouws, D. Britz, A. Goldie, B. Strope, and R. Kurzweil, "Generating high-quality and informative conversation responses with sequence-to-sequence models," in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, 2017, pp. 2210–2219. [Online]. Available: <https://aclanthology.info/papers/D17-1235/d17-1235>
- [70] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. C. Courville, and Y. Bengio, "Maxout networks," *CoRR*, vol. abs/1302.4389, 2013. [Online]. Available: <http://arxiv.org/abs/1302.4389>
- [71] G. Zhou, P. Luo, Y. Xiao, F. Lin, B. Chen, and Q. He, "Elastic responding machine for dialog generation with dynamically mechanism selecting," in *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, New Orleans, Louisiana, USA, February 2-7, 2018*, 2018. [Online]. Available: <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16318>
- [72] M. Ghazvininejad, C. Brockett, M. Chang, B. Dolan, J. Gao, W. Yih, and M. Galley, "A knowledge-grounded neural conversation model," *CoRR*, vol. abs/1702.01932, 2017. [Online]. Available: <http://arxiv.org/abs/1702.01932>
- [73] T. Young, E. Cambria, I. Chaturvedi, H. Zhou, S. Biswas, and M. Huang, "Augmenting end-to-end dialogue systems with commonsense knowledge," in *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, New Orleans, Louisiana, USA, February 2-7, 2018*, 2018. [Online]. Available: <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16573>
- [74] S. R. Bowman, L. Vilnis, O. Vinyals, A. M. Dai, R. Józefowicz, and S. Bengio, "Generating sentences from a continuous space," in *Proceedings of the 20th SIGNLL Conference on Computational Natural Language Learning, CoNLL 2016, Berlin, Germany, August 11-12, 2016*, 2016, pp. 10–21. [Online]. Available: <http://aclweb.org/anthology/K/K16/K16-1002.pdf>
- [75] S. Semeniuta, A. Severyn, and E. Barth, "A hybrid convolutional variational autoencoder for text generation," in *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP 2017, Copenhagen, Denmark, September 9-11, 2017*, 2017, pp. 627–637. [Online]. Available: <https://aclanthology.info/papers/D17-1066/d17-1066>
- [76] S. Clark and K. Cao, "Latent variable dialogue models and their diversity," in *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2017, Valencia, Spain, April 3-7, 2017, Volume 2: Short Papers*, 2017, pp. 182–187. [Online]. Available: <https://aclanthology.info/papers/E17-2029/e17-2029>
- [77] L. Yu, W. Zhang, J. Wang, and Y. Yu, "Seqgan: Sequence generative adversarial nets with policy gradient," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*, 2017, pp. 2852–2858. [Online]. Available: <http://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/14344>
- [78] W. Fedus, I. Goodfellow, and A. M. Dai, "MaskGAN: Better text generation via filling in the," in *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=ByOExmWAB>
- [79] T. Mikolov, M. Karafiát, L. Burget, J. Cernocký, and S. Khudanpur, "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*, 2010, pp. 1045–1048. [Online]. Available: <http://www.isca-speech.org/archive/interspeech-2010/i10-1045.html>
- [80] H. Mei, M. Bansal, and M. R. Walter, "Coherent dialogue with attention-based language models," in *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*, 2017, pp. 3252–3258. [Online]. Available: <http://aaai.org/ocs/index.php/AAAI/AAAI17/paper/view/14164>
- [81] R. S. Sutton, D. A. McAllester, S. P. Singh, and Y. Mansour, "Policy gradient methods for reinforcement learning with function approximation," in *Advances in Neural Information Processing Systems 12, [NIPS Conference, Denver, Colorado, USA, November 29 - December 4, 1999]*, 1999, pp. 1057–1063. [Online]. Available: <http://papers.nips.cc/paper/1713-policy-gradient-methods-for-reinforcement-learning-with-function-approximation>
- [82] T. Zhang, M. Huang, and L. Zhao, "Learning structured representation for text classification via reinforcement learning," in *AAAI-18 AAAI Conference on Artificial Intelligence*, 2018.
- [83] J. Li, W. Monroe, A. Ritter, D. Jurafsky, M. Galley, and J. Gao, "Deep reinforcement learning for dialogue generation," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, 2016, pp. 1192–1202. [Online]. Available: <http://aclweb.org/anthology/D/D16/D16-1127.pdf>
- [84] S. Wiseman and A. M. Rush, "Sequence-to-sequence learning as beam-search optimization," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, 2016, pp. 1296–1306. [Online]. Available: <http://aclweb.org/anthology/D/D16/D16-1137.pdf>
- [85] J. Li, W. Monroe, and D. Jurafsky, "A simple, fast diverse decoding algorithm for neural generation," *CoRR*, vol. abs/1611.08562, 2016. [Online]. Available: <http://arxiv.org/abs/1611.08562>
- [86] J. Li, M. Galley, C. Brockett, J. Gao, and B. Dolan, "A diversity-promoting objective function for neural conversation models," in *NAACL HLT 2016, The 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, San Diego California, USA, June 12-17, 2016*, 2016, pp. 110–119. [Online]. Available: <http://aclweb.org/anthology/N/N16/N16-1014.pdf>
- [87] Y. Wu, W. Wu, D. Yang, C. Xu, and Z. Li, "Neural

- response generation with dynamic vocabularies," in *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, New Orleans, Louisiana, USA, February 2-7, 2018*, 2018. [Online]. Available: <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16136>
- [88] Y. Wang, C. Liu, M. Huang, and L. Nie, "Learning to ask questions in open-domain conversational systems with typed decoders," 2018.
- [89] K. Papineni, S. Roukos, T. Ward, and W. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics, July 6-12, 2002, Philadelphia, PA, USA., 2002*, pp. 311–318. [Online]. Available: <http://www.aclweb.org/anthology/P02-1040.pdf>
- [90] S. Banerjee and A. Lavie, "METEOR: an automatic metric for MT evaluation with improved correlation with human judgments," in *Proceedings of the Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization@ACL 2005, Ann Arbor, Michigan, June 29, 2005*, 2005, pp. 65–72. [Online]. Available: <https://aclanthology.info/papers/W05-0909/w05-0909>
- [91] C.-Y. Lin, "Rouge: A package for automatic evaluation of summaries," in *Text Summarization Branches Out: Proceedings of the ACL-04 Workshop*, 2004, pp. 74–81.
- [92] C. Liu, R. Lowe, I. Serban, M. Noseworthy, L. Charlin, and J. Pineau, "How NOT to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation," in *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, 2016, pp. 2122–2132. [Online]. Available: <http://aclweb.org/anthology/D/D16/D16-1230.pdf>
- [93] A. Kannan and O. Vinyals, "Adversarial evaluation of dialogue models," *CoRR*, vol. abs/1701.08198, 2017. [Online]. Available: <http://arxiv.org/abs/1701.08198>
- [94] C. Tao, L. Mou, D. Zhao, and R. Yan, "RUBER: an unsupervised method for automatic evaluation of open-domain dialog systems," in *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence, New Orleans, Louisiana, USA, February 2-7, 2018*, 2018. [Online]. Available: <https://www.aaai.org/ocs/index.php/AAAI/AAAI18/paper/view/16179>
- [95] R. Kiros, Y. Zhu, R. Salakhutdinov, R. S. Zemel, R. Urtasun, A. Torralba, and S. Fidler, "Skip-thought vectors," in *Advances in Neural Information Processing Systems 28: Annual Conference on Neural Information Processing Systems 2015, December 7-12, 2015, Montreal, Quebec, Canada, 2015*, pp. 3294–3302. [Online]. Available: <http://papers.nips.cc/paper/5950-skip-thought-vectors>
- [96] Y. Zhu, R. Kiros, R. S. Zemel, R. Salakhutdinov, R. Urtasun, A. Torralba, and S. Fidler, "Aligning books and movies: Towards story-like visual explanations by watching movies and reading books," in *2015 IEEE International Conference on Computer Vision, ICCV 2015, Santiago, Chile, December 7-13, 2015*, 2015, pp. 19–27. [Online]. Available: <https://doi.org/10.1109/ICCV.2015.11>
- [97] Y. Zhao and Q. Zhu, "Evaluation on crowdsourcing research: Current status and future direction," *Information Systems Frontiers*, vol. 16, no. 3, pp. 417–434, 2014. [Online]. Available: <https://doi.org/10.1007/s10796-012-9350-4>
- [98] S. Sharma, L. E. Asri, H. Schulz, and J. Zumer, "Relevance of unsupervised metrics in task-oriented dialogue for evaluating natural language generation," *CoRR*, vol. abs/1706.09799, 2017. [Online]. Available: <http://arxiv.org/abs/1706.09799>
- [99] L. Kaiser and S. Bengio, "Can active memory replace attention?" in *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain, 2016*, pp. 3774–3782. [Online]. Available: <http://papers.nips.cc/paper/6295-can-active-memory-replace-attention>
- [100] N. Kalchbrenner, L. Espeholt, K. Simonyan, A. van den Oord, A. Graves, and K. Kavukcuoglu, "Neural machine translation in linear time," *CoRR*, vol. abs/1610.10099, 2016. [Online]. Available: <http://arxiv.org/abs/1610.10099>
- [101] J. Gehring, M. Auli, D. Grangier, D. Yarats, and Y. N. Dauphin, "Convolutional sequence to sequence learning," in *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, 2017, pp. 1243–1252. [Online]. Available: <http://proceedings.mlr.press/v70/gehring17a.html>
- [102] M. Chen and L. Denoyer, "Multi-view generative adversarial networks," in *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2017, Skopje, Macedonia, September 18-22, 2017, Proceedings, Part II*, 2017, pp. 13675–188. [Online]. Available: <https://doi.org/10.1007/978-3-319-71246-8-11>
- [103] J. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired image-to-image translation using cycle-consistent adversarial networks," in *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, 2017, pp. 2242–2251. [Online]. Available: <https://doi.org/10.1109/ICCV.2017.244>
- [104] Z. Yi, H. R. Zhang, P. Tan, and M. Gong, "Dualgan: Unsupervised dual learning for image-to-image translation," in *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, 2017, pp. 2868–2876. [Online]. Available: <https://doi.org/10.1109/ICCV.2017.310>
- [105] T. Kim, M. Cha, H. Kim, J. K. Lee, and J. Kim, "Learning to discover cross-domain relations with generative adversarial networks," in *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, 2017, pp. 1857–1865. [Online]. Available: <http://proceedings.mlr.press/v70/kim17a.html>