

Date of publication xxxx 00, 0000, date of current version xxxx 00, 0000.

Digital Object Identifier 10.1109/ACCESS.2017.DOI

Open Domain Response Generation Guided by Retrieved Conversations

CHANG SHU¹, ZIJIAN ZHANG^{2,3}, YOUXIN CHEN², JING XIAO², JEY HAN LAU⁴, QIAN ZHANG¹ AND ZHENG LU¹

¹University of Nottingham Ningbo China, Ningbo, China (e-mail: scxcs1, qian.zhang, zheng.lu@nottingham.edu.cn)

²Ping An Technology Co., Ltd, Shanghai, China (e-mail: chenyouxin149, xiaojing661@pingan.com.cn)

³Meituan-Dianping Group, Shanghai, China (e-mail: zijian.zh96@gmail.com)

⁴University of Melbourne, Melbourne, Australia (e-mail: laujh@unimelb.edu.au)

Corresponding author: Zheng Lu (e-mail: zheng.lu@nottingham.edu.cn).

This work is supported by Ningbo Science and Technology Bureau under Service Industry S&T Programme with project code 2019F1028 and Major Projects Fund with project code 2021Z089.

ABSTRACT Open domain response generation is the task of creating a response given a user query in any topics/domain. Limited by context and reference information, responses generated by current systems are often “bland” or generic. In this paper, we combine a response generation model with a retrieval system that searches for relevant utterances and responses. The generation model has two main components: a keyword extraction module and a two-stage transformer. The keyword extraction module aims to extract two types of keywords in an unsupervised fashion from the retrieved results: (1) keywords in the *query* not found in the *retrieved utterances* (DiffKey), and (2) overlapping keywords among the *retrieved responses* (SimKey). Given these keywords, the two-stage transformer first decides where to insert the keywords in the response, and the second generates the full response given the location of the keywords. The keyword extraction module and the two-stage transformer are connected in a single network, and so our system is trained end-to-end. Experimental results on Cornell Movie-Dialog corpus, Douban and Weibo demonstrate that our model outperforms state-of-the-art systems in terms of ROUGE, relevance scores and human evaluation.

INDEX TERMS Dialogue generation, Hybrid retrieval-generation, Deep learning

I. INTRODUCTION

OPEN domain response generation aims to develop conversational agents that can interact and communicate in a variety of topics [1]–[4], and it differs from task-oriented dialogue systems which are designed to work towards a specific goal in a particular domain (e.g. finding a restaurant). There are a host of dialogue agents nowadays, and they often combine open domain and task-oriented dialogue generation, e.g. Microsoft XiaoIce, Apple Siri and Google Assistant.

Our work focuses on open domain response generation, and specifically in the single-turn setting (i.e. the conversation lasts only one turn, and so consists of a query and reply). There are generally two approaches for this task: retrieval and generation methods. Retrieval approaches search for answers from an existing corpus of dialogues to use them as response. Responses created by retrieval methods tend to be partially relevant and often do not directly address the queries, as the corpus is unlikely to have full coverage for all queries. Generation methods, on the other hand, are able to

create fitting and natural responses but they tend to be short and generic [2]. Combining both approaches would allow us to generate responses that are more diverse, interesting and relevant. Although there are a number of studies that explore combining both approaches, recent studies train the retrieval and generation component separately, with each component requiring different training data [5]–[8]. Existing studies also tend to use hidden representations as additional signals — e.g. that of skeleton words [5] or abstract [9] — while our method uses keywords *directly*, and so is more interpretable as we can analyse specifically what keywords lead to a particular generated response.

In this paper, we introduce a novel end-to-end system that combines retrieval and generation methods for open domain response generation. Given a query (q), our model first searches for relevant utterances (u') and responses (r') in a corpus of existing dialogues, and extracts two sets of keywords: DiffKey and SimKey. Using Table 1 as an example, DiffKey corresponds to words in the query (q) that are not

TABLE 1. An example of query (q) and its retrieved utterances (u') and responses (r'). DIFFKEY is highlighted, while SIMKEY is bolded.

Query (q)	What's your suggestion about holiday ? How about going to the seaside ?
Utterance1 (u'_1)	What's your suggestion for tomorrow? How about outdoor sports?
Utterance2 (u'_2)	Tomorrow is the weekend. I give your husband a suggestion to have a rest.
Retrieved	
DIFFKEY	holiday / going / to / seaside
Response1 (r_1)	Leave me alone. I just want to have a rest at home and do some housework.
Response2 (r_2)	You are right. He just at home for one day last month and didn't have enough rest.
SIMKEY	just/have/rest/at/home
Response (r)	I just want to rest at home on holiday , not go to the seaside .

found in the retrieved utterances (u'_1 and u'_2), e.g. *holiday* and *seaside*; while SimKey are overlapping words in the retrieved responses (r'_1 and r'_2), e.g. *rest* and *home*. Intuitively, DiffKey are keywords that are not captured by existing dialogues, and they are extracted so that the generated response would include them to improve its relevance to the query. SimKey can be interpreted as guiding keywords — the fact that they are frequently mentioned in the retrieved responses indicate that they are useful keywords to be incorporated in the generated response. To capture word similarity beyond their surface forms, our system leverages transformer’s attention mechanism to extract these keywords.

Given these keywords, we use a two-stage transformer to generate the final response. The first transformer takes the keywords as (unordered) input and decides where to insert them in the final response, creating a sentence where the predicted positions contain the keywords and other positions are masked tokens (e.g. “[mask] just [mask] [mask] rest at home [mask] holiday [mask] [mask] go to [mask] seaside”). The second transformer works like a text infilling model [10], where it takes the masked sentence as input and “fill in the blanks” to generate the final response. The keyword extraction module and the two-stage transformer are connected in a single network, and as such the full model is trained end-to-end, requiring only a dialogue corpus as training data.

We conduct experiments on English and Chinese dialogue datasets and demonstrate that our system outperforms benchmark systems consistently across three datasets based on ROUGE [11], relevance, diversity [2] scores and human evaluation, creating a new state-of-the-art for open domain response generation. We also perform an ablation study to measure the impact of DiffKey, SimKey and two-stage transformer.

To summarise, our contributions are given as follows: (1)

we propose a new method to extract keywords based on multi-source alignment to guide generation; (2) we design a two-stage transformer that uses two BERT models to do generation; (3) we verify the empirical effectiveness of our approach by comparing to benchmark systems and find state-of-the-art performance in open domain dialogue generation.

II. RELATED WORK

Response generation can be broadly categorised into retrieval-based, generation-based and hybrid methods, which we review below.

Retrieval-based methods. Given an utterance, retrieval-based methods rely on matching algorithms to find the most relevant utterance in the conversation history to use its response as the output. The key is in developing matching algorithms that can measure textual relevance between two utterances [12]. Nowadays, retrieval models typically use semantic retrieval rather than keyword retrieval, thanks to the advent of semantic matching [13]. For hybrid or large-scale models, the latter is faster and more efficient [14]. Early studies mainly focus on response selection for single-turn conversations [15]. More recently, multi-turn retrieval-based conversation methods [16] or retrieval of documents given single question [17] are also explored.

Generation-based methods. By and large, generation methods use the sequence-to-sequence framework [18]–[20] for response generation. Attention [20] and copy [21] mechanisms have been widely used to improve the performance of the original sequence-to-sequence framework. Shen et al. [22] propose a hierarchical self-attention mechanism and distant supervision to find related information globally when decoding. As generated responses tend to be generic, several methods are proposed to improve the diversity of the generated responses, e.g. by incorporating topic information [23] or using latent variable models and gate mechanism [24]. Li et al. [25] train an encoder-decoder model in a bidirectional method by adding a backward reasoning step to avoid generic and dull responses. Also some popular pre-train language models, such as UniLM [26] and GPT-3 [27] are fine-tuned as an encoder-decoder architecture for response generation. Li et al. [28] experiment with reinforcement learning to further improve generation quality, and Liu et al. [29] incorporate adversarial learning to reduce gender bias in its response generation. Ling et al. [30] explore a response generation model that utilizes contextual topics to find relevant transition words. Xu et al. [31] incorporate a keyword decoder to generate keywords based on the dialogue history and feed these keywords to the response generator. Other studies experiment with integrating extra information such as tables [32] and knowledge graph [33]. Xu et al. [34] inject knowledge into the pre-trained language models to diversify and enrich generated dialogue responses. Conditional variational auto-encoder proposed by Zhao et al. [35] is proposed to similarly improve the diversity during decoding. [36]–[38] attempt to generate responses that are empathetic.

Hybrid methods. Song et al. [39] propose combining both generation and retrieval methods for generating responses. Pandey et al. [40] retrieve similar conversations and weight them to guide generation. Miao, Cao and et al. [41], [42] develop retrieve-then-edit techniques for text generation which can improve the quality of the generated response. Cai, Tian, Kazemnejad and et al. [6], [7], [43] treat the retrieval and generation as disjointed components and train them separately, but this means additional data is needed. Multi-task learning that jointly optimize retrieval and generation steps are also explored [44]. Unlike other studies that largely focus on improving the generation component, Wu et al. [45] propose improving the performance of the retrieval component through entity alignment. Xu et al. [46] found retrieval-augmented methods that have the ability to summarize and recall previous conversations helpful. Xia et al. [47] use two decoders to generate a raw sequence and revise the draft from scratch. Li et al. [48] utilise an edit keywords to guide sentences generation.

Compared to previous studies, our proposed method is unique in how it extracts keywords from the retrieved conversations (most studies only use the retrieved conversations as additional input without keyword extraction [6], [9], [49]). Furthermore, our model is different to a traditional encoder-decoder architecture, where we only use two BERT [50] models to generate the final response. The closest work to ours is Wu et al. [51], but it only uses the top-1 retrieved results and so do not consider overlapping keywords among the responses. In contrast, our system retrieves top- K ($K > 1$) results from the corpus, and use multi-source alignment to extract two different types of keywords.

III. MODEL ARCHITECTURE

A. MODEL OVERVIEW

The overall architecture of our system is presented in Figure 1, which consists of a retrieval model and a generation model. Given a query q , the retrieval model first retrieves top- N utterance-response pairs (u'_i, r'_i) from corpus \mathcal{D} . The generation model then extracts the keywords (DiffKey and SimKey) using the semantic alignment keyword extraction (SAKE) module, and the extracted keywords are fed to the two-stage transformer to generate the final response \hat{r} .

B. RETRIEVAL MODEL

We use Lucene¹ to index and find top- K ($K = 2$) best utterances in corpus \mathcal{D} based on Jaccard similarity:²

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

where A and B are the bag-of-words of utterances.

¹<https://lucene.apache.org/>

²We consider only utterances that have Jaccard similarity between 0.5–0.9.

C. SEMANTIC ALIGNMENT KEYWORD EXTRACTION (SAKE)

In SAKE, the goal is to extract DiffKey and SimKey, given the query q and retrieved utterance-response pairs (u'_i, r'_i) . We first use a 1-layer transformer to encode the text:

$$e = \mathbf{TF}(S, P) \quad (1)$$

where S represents either q , u'_i , or r'_i ; and P is the corresponding positional embeddings. $e \in \mathbb{R}^{L \times D}$ is the contextualised word embeddings (where L denotes the length of sentence and D is the embedding dimension).

To align two sentences, we adapt the alignment approach by Tsai et al. [52], which was proposed to align sequences of different modalities (e.g. text with audio). We first provide a generic description of the alignment method, and come back to explain how to extract DiffKey and SimKey based on q , u'_i , and r'_i .

Single-source Alignment. Given two sentences, α and β , our goal is to align words in α to the words in β via query/key/value attention.

After encoding the two sentences with a transformer (Equation 1), we have two embeddings e_α and e_β . They have the same embedding dimensions D but different sentence lengths. We then project the two embeddings to query, key and value vectors by three learnable weights matrixes. Thus, both of them have the same projected dimension H . The semantic alignment output Y is computed as follows:

$$Y = \mathbf{SA}(e_\alpha, e_\beta) = \text{softmax} \left(\frac{Q_\alpha K_\beta^T}{\sqrt{d_k}} \right) V_\beta$$

Multi-source Alignment. Here we extend the alignment of one sentence to N sentences, i.e. α is one sentence but β is now a group of sentences $\beta = \{\beta_1, \dots, \beta_N\}$, by aligning a pair of sentences iteratively and summing up their outputs:

$$Y^{[N]} = \sum_i^N \mathbf{SA}(e_\alpha, e_{\beta_i})$$

Figure 2 presents an illustration of the multi-source alignment method. Note that the key and value projection matrix (W_K and W_V) are shared by all $\{\beta_1, \dots, \beta_N\}$.

Recall that DiffKey are keywords in query q that are not found in the retrieved utterance u'_i . In this case, $\alpha = q$, and $\beta = \{u'_i\}_{i=1}^N$. For SimKey, they are the overlapping words between the top-1 retrieved response r'_1 and other retrieved responses $\{r'_i\}_{i=2}^N$, and so $\alpha = r'_1$ and $\beta = \{r'_i\}_{i=2}^N$. Formally:

$$Y_{\text{DiffKey}}^{[N]} = \sum_{i=1}^N \mathbf{SA}(e_q, e_{u'_i})$$

$$Y_{\text{SimKey}}^{[N]} = \sum_{i=2}^N \mathbf{SA}(e_{r'_1}, e_{r'_i}) \quad (2)$$

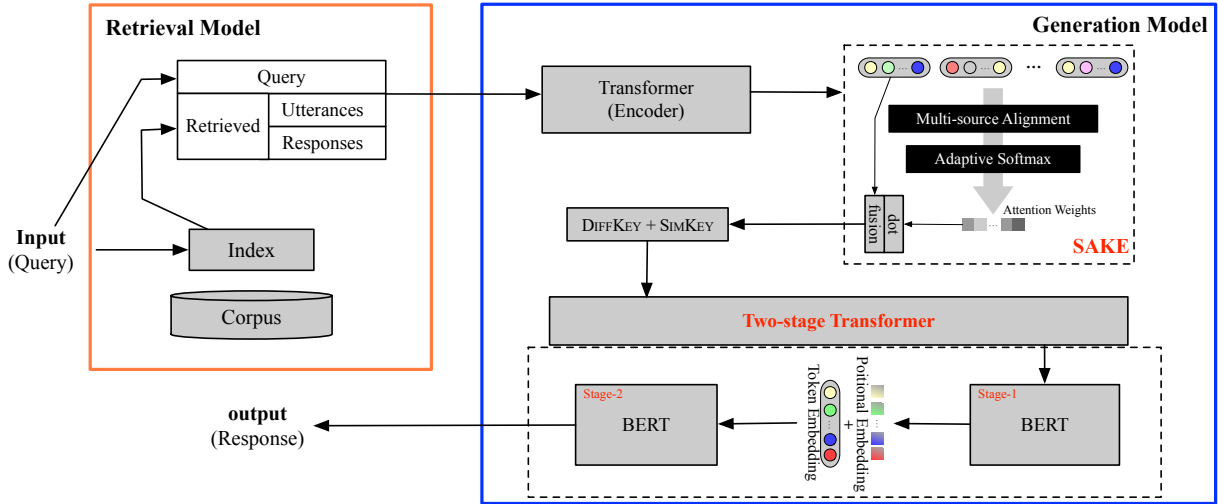


FIGURE 1. The architecture of proposed retrieval-generation model. The red box denotes the retrieval model and blue box is the generation model, which consists of a SAKE module to produce DIFFKEY and SIMKEY and a two-stage transformer to generation response.

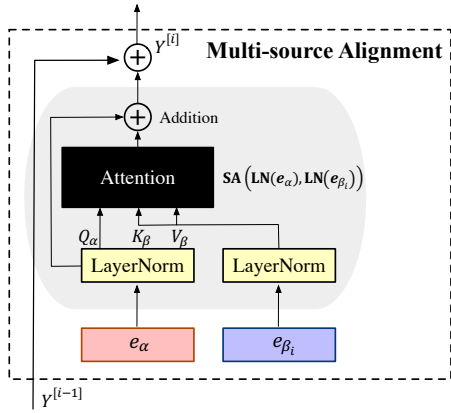


FIGURE 2. Multi-source alignment α (single sentence) and β (N sentences). β_i denotes the i^{th} sentence of retrieved utterance and $Y^{[i]}$ is the i^{th} aligned matrix.

To calculate the attention weights for each words in the query (q) for DiffKey or in the first response (r'_1) for SimKey, we compute $M = softmax(Y^{[N]} \cdot W_S + b_S)$, where $Y^{[N]}$ represents either $Y_{DiffKey}^{[N]}$ and $Y_{SimKey}^{[N]}$.

After obtaining the attention weights, we use them to weight the word embeddings as a soft approach to ‘extract’ the keywords.³ Using Table 1 as an example for DiffKey, we would weight all the word embeddings in the query text (q) “What’s your suggestion about holiday? How about going to the seaside?”, and words that receive low attention weights (such as *what*) would be effectively masked out. Note that at test time, we use an *argmin* and *argmax* operator to extract ηL words from q for DiffKey or r'_1 for SimKey respectively, where η is a scaling hyper-parameter that controls how many keywords to extract based on the original length of q or r'_1 .

³Strictly speaking they are subword embeddings, but for ease of interpretation we use the term “word embeddings” here.

D. TWO-STAGE TRANSFORMER

The DiffKey and SimKey produced by SAKE are keywords without ordering or positional information. To use them as input to guide the response generation,⁴ we first use a BERT model [50] to predict their positions in the response, and then use another BERT to generate the final response. Note that this second BERT is *not* a fill-in-the-blanks model, as the final response is constructed by taking the highest probability word at every position. Also, during training we update only the second BERT (first BERT parameters are kept static).

Stage-1 Transformer. To imbue the keywords with positional information, we feed them to BERT to predict their positions:

$$\begin{aligned}
 g &= \mathbf{BERT}_1([\text{DiffKey}; \text{SimKey}]) \\
 p_i &= \text{softmax}(W_1 g_i + b_1) \\
 q_i &= \sum_j p_{i,j} P_j
 \end{aligned} \tag{3}$$

where P_j is the predefined static positional embedding [53] for position j . The maximum length of predicted position is 150, which is a hyperparameter that can be defined. Intuitively, for a word in DiffKey or SimKey, p_i represents its probability distribution over different positions, and q_i is weighted positional embedding.

Stage-2 Transformer. The second transformer is also a BERT, and similarly takes DiffKey and SimKey as input and its goal is to generate the final response. Here, we add the weighted positional embeddings (q) from the stage-1

⁴To clarify, during training DIFFKEY consists of the whole query (q) and SIMKEY the first response (r'_1) (noting that their embeddings are weighted by SAKE), but at test time DIFFKEY and SIMKEY contain only a subset of words (selected by the *argmin* and *argmax* operators respectively).

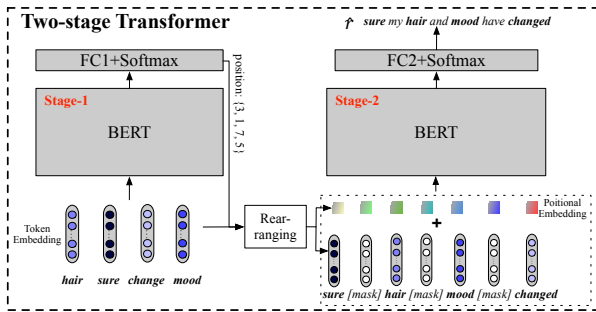


FIGURE 3. Input and output of the two-stage transformer at test time. The left subfigure is the stage-1 transformer to imbue the set of DIFFKEY and SIMKEY with position information. And the stage-2 transformer (right) aims to generate a target response given the keywords.

transformer to the input:

$$h = \mathbf{BERT}_2([\text{DiffKey}; \text{SimKey}] + Q([\text{DiffKey}; \text{SimKey}]))$$

$$\hat{r}_i = \text{softmax}(W_2 h_i + b_2)$$

where Q is a function that applies the weighted positional embeddings (Equation 3) to each input word, and \hat{r}_i is the output word probability distribution, and the whole model is optimised end-to-end based on cross-entropy loss: $\mathcal{L} = -\sum_i \log P(\hat{r}_i)$.

At test time, instead of computing the weighted positional embeddings (Equation 3), we use argmax to select the best position for each keyword, and introduce an additional step to re-arrange the keywords before feeding them to the stage-2 transformer; see Figure 3 for an illustration. We also truncate the generated response after $\langle \text{EOS} \rangle$ is produced at test time (i.e. all words to the right of $\langle \text{EOS} \rangle$ are discarded).

Intuitively, our stage-2 transformer can be interpreted as a text infilling model, where it takes a masked sentence (that contains only important keywords) and learns how to “fill in the blanks” to create the response. As such, the generation process does not involve any decoding algorithms.

IV. EXPERIMENTS

A. DATASETS AND EVALUATION METRICS

We use two Chinese datasets (Douban [51] and Weibo⁵) and an English dataset (Cornell Movie-Dialog corpus⁶ [54]) to evaluate our response generation system. All three datasets consist of human conversations in the form of utterance and response pairs. For Douban, there are 19,623,374 original pairs. After removing pairs with high proportion of symbols (e.g. punctuations and emoticons) and very long sentences (> 50 words), we retain 11,321,313 pairs. Weibo and Cornell each has 4,281,692 and 430,579 pairs after undergoing the same preprocessing. We release the source code of our experiments and the Weibo dataset to facilitate replication and research.⁷ To evaluate our model, we use the following metrics:

⁵We collect the data by scraping microblog posts from <https://weibo.com/>

⁶https://www.cs.cornell.edu/cristian/Cornell_Movie-Dialogs_Corpus.html

⁷ANONYMISED.COM

ROUGE: This metric measures the similarity between the generated response and ground truth response by evaluating their n -gram overlap.

Relevance: This also measures similarity using cosine similarity of word embeddings instead of evaluating textual relevance. To aggregate the word embeddings of a sentence, we follow Liu et al. [55] by taking the mean embeddings (“Average”) and max-pooled embeddings (i.e. maximum value over words for each dimension; “Max”) before computing the cosine similarity. We also compute another variant where we do not pool the word embeddings but greedily find the best matching words in the text pairs (“Greedy”).

Diversity: This measures the repetitiveness of the generated response, and is computed based on the ratios of distinct unigrams (Dist-1) and bigrams (Dist-2). This metric does not use the ground truth response.

Human: Thirty annotators are invited to judge the generated responses of different systems on two aspects on a 4-point ordinal scale: fluency (F) and relevance (R). Details of the crowdsourcing experiments are provided in the Appendix.

B. BASELINES/BENCHMARKS

We compare our method against the following baseline/benchmark systems (which covers sequence-to-sequence, retrieval and hybrid methods):

S2S+Attn: Recurrent network-based sequence-to-sequence with attention model.

CVAE: Conditional variational auto-encoder proposed by to improve the diversity of generated responses.

KW+S2S: A generation-based model that uses a keyword encoder-decoder to generate keywords given the dialogue history, which are then concatenated with the dialogue history to generate the response. KW+S2S is trained end-to-end and the ground truth keywords are extracted using TF-IDF.

UniLM, BERT, GPT-3: These are pre-trained language models fine-tuned for response generation.⁸ We only have English results for GPT-3 as it does not support Chinese.

Retrieval: Baseline retrieval model that searches for the most relevant utterance (Lucene) and returns its response as the result.

Rtv+Rank: Retrieval method that searches for top-20 utterances based on **Retrieval**; two LSTM models are trained to encode pairs of utterances to select the most relevant response [56].

Edit: Hybrid method that retrieves the most relevant utterance and computes two edit vectors to represent novel words in the query and the retrieved utterance to guide response generation. Note that this method retrieves only 1 relevant utterance, and as such does not capture similarity among relevant responses like our model.

Reranker: Hybrid method that has 2 components: (1) a generator that takes encoded representations of conversation context and retrieved responses as input to generate a response;

⁸Note that for BERT we generate the full sentence by selecting the highest probability word in each position in one step and do not do left-to-right decoding (as it does not have a decoder).

TABLE 2. Results on the datasets of Cornell (top), Weibo (middle) and Douban (bottom). **Boldfont** indicates optimal performance for a metric in a dataset. For types, "Rtv" = retrieval method, "Gen"= generation method, and "Pretrained" = whether it uses pre-trained models.

Types			Models	ROUGE			Relevance			Diversity		Human	
Pretrained	Rtv	Gen		R-1	R-2	R-L	Average	Max	Greedy	Dist-1	Dist-2	F	R
		✓	S2S+Attn	37.82	17.87	33.73	0.314	0.157	0.327	0.049	0.088	2.81	3.04
		✓	CVAE	41.89	20.86	39.49	0.339	0.182	0.357	0.076	0.145	3.14	3.13
		✓	KW+S2S	47.14	24.05	42.86	0.378	0.214	0.387	0.133	0.242	3.51	3.52
✓		✓	BERT	42.81	21.49	39.92	0.364	0.211	0.366	0.104	0.172	3.43	3.57
✓		✓	UniLM	44.24	23.07	40.27	0.373	0.205	0.371	0.121	0.189	3.47	3.38
✓		✓	GPT-3	47.11	23.92	41.77	0.378	0.209	0.389	0.137	0.211	3.58	3.56
	✓		Retrieval	30.81	13.87	27.33	0.252	0.131	0.269	0.103	0.249	3.81	3.20
	✓		Rtv+Rank	34.57	18.23	32.44	0.327	0.157	0.336	0.129	0.212	3.84	3.25
	✓	✓	Edit	45.81	21.99	43.01	0.369	0.198	0.376	0.112	0.207	3.44	3.51
	✓	✓	Reranker	45.16	21.04	42.71	0.357	0.194	0.380	0.094	0.182	3.47	3.51
	✓	✓	MemDistill	46.76	22.59	43.67	0.374	0.204	0.386	0.107	0.212	3.48	3.49
	✓	✓	SkelGen _{LSTM}	46.69	21.13	43.09	0.369	0.211	0.375	0.116	0.231	3.43	3.51
	✓	✓	SkelGen _{GPT-2}	48.16	23.09	43.22	0.377	0.216	0.389	0.108	0.221	3.53	3.59
✓	✓	✓	RAG	46.81	22.19	42.78	0.375	0.211	0.382	0.104	0.192	3.52	3.57
✓	✓	✓	Ours	50.25	25.02	45.22	0.393	0.225	0.405	0.124	0.233	3.60	3.62

Types			Models	ROUGE			Relevance			Diversity		Human	
Pretrained	Rtv	Gen		R-1	R-2	R-L	Average	Max	Greedy	Dist-1	Dist-2	F	R
		✓	S2S+Attn	36.77	20.14	35.07	0.346	0.179	0.358	0.026	0.084	3.12	3.08
		✓	CVAE	44.15	23.12	41.39	0.361	0.187	0.374	0.086	0.142	3.30	3.25
		✓	KW+S2S	50.32	25.23	48.01	0.390	0.241	0.392	0.163	0.222	3.45	3.49
✓		✓	BERT	48.64	27.71	48.73	0.391	0.262	0.397	0.147	0.285	3.58	3.57
✓		✓	UniLM	50.33	30.19	49.81	0.403	0.267	0.408	0.142	0.342	3.51	3.56
	✓		Retrieval	32.41	15.21	28.03	0.302	0.153	0.322	0.111	0.472	3.82	3.41
	✓		Rtv+Rank	38.29	18.17	35.14	0.361	0.174	0.378	0.167	0.494	3.87	3.40
	✓	✓	Edit	50.82	26.71	48.38	0.394	0.253	0.401	0.152	0.158	3.44	3.56
	✓	✓	Reranker	50.64	26.60	47.81	0.386	0.236	0.397	0.125	0.161	3.41	3.52
	✓	✓	MemDistill	50.82	26.71	48.38	0.394	0.253	0.401	0.152	0.158	3.47	3.55
	✓	✓	SkelGen _{LSTM}	51.14	26.79	49.03	0.401	0.261	0.407	0.158	0.181	3.45	3.56
	✓	✓	SkelGen _{GPT-2}	53.28	30.11	51.02	0.414	0.289	0.411	0.142	0.201	3.54	3.60
✓	✓	✓	RAG	52.02	27.69	50.31	0.409	0.276	0.413	0.152	0.271	3.50	3.57
✓	✓	✓	Ours	55.13	31.21	52.79	0.430	0.305	0.434	0.218	0.350	3.59	3.66

Types			Models	ROUGE			Relevance			Diversity		Human	
Pretrained	Rtv	Gen		R-1	R-2	R-L	Average	Max	Greedy	Dist-1	Dist-2	F	R
		✓	S2S+Attn	33.74	18.16	30.62	0.331	0.172	0.345	0.061	0.081	3.08	3.22
		✓	CVAE	42.32	21.83	38.78	0.362	0.189	0.368	0.076	0.201	3.31	3.29
		✓	KW+S2S	47.92	29.11	45.98	0.374	0.208	0.378	0.142	0.281	3.52	3.48
✓		✓	BERT	45.11	27.78	46.03	0.371	0.216	0.379	0.135	0.206	3.46	3.51
✓		✓	UniLM	48.76	31.89	47.09	0.382	0.223	0.394	0.202	0.364	3.51	3.49
	✓		Retrieval	30.19	14.88	28.36	0.274	0.142	0.289	0.131	0.466	3.75	3.38
	✓		Rtv+Rank	36.49	17.67	35.44	0.354	0.181	0.361	0.137	0.431	3.84	3.41
	✓	✓	Edit	48.27	29.81	46.53	0.393	0.216	0.385	0.134	0.189	3.48	3.50
	✓	✓	Reranker	48.21	29.91	47.37	0.378	0.206	0.383	0.128	0.237	3.46	3.46
	✓	✓	MemDistill	48.82	30.21	47.71	0.381	0.212	0.387	0.131	0.231	3.52	3.57
	✓	✓	SkelGen _{LSTM}	49.18	30.31	48.72	0.384	0.220	0.393	0.102	0.251	3.54	3.56
	✓	✓	SkelGen _{GPT-2}	52.15	32.03	52.79	0.397	0.270	0.421	0.110	0.212	3.56	3.59
✓	✓	✓	RAG	50.28	30.18	49.62	0.387	0.238	0.402	0.137	0.302	3.54	3.58
✓	✓	✓	Ours	54.64	33.72	53.68	0.406	0.275	0.432	0.198	0.314	3.61	3.65

and (2) a neural reranker that selects the best response among generated and retrieved responses.

MemDistill: Hybrid method that first clusters training query-response pairs and stores them in memory, and trains a generator to retrieve the most relevant query-response cluster from the memory to guide its generation. The method is unique in that it uses query-response cluster as a guide (rather than individual responses like our and other benchmark

systems).

SkelGen: Hybrid transformer-based method that reranks a set of retrieved responses to select the best response as input for the generator to create a response. The reranker is trained separately (using ground truth query-response pairs) to the generator, and the framework does not extract any keywords (the best response only serves as additional sequence to generator). Depending on whether the generator

TABLE 3. Ablation results where we measure the impact of DIFFKEY, SIMKEY and stage-1 transformer ("Ours" presents the average scores on main experiments). Symbol "-X" denotes that module X is removed.

Models	ROUGE	Relevance	Diversity	Human	
				F	R
Ours	45.07	0.357	0.243	3.59	3.61
-SIMKEY	37.12	0.289	0.136	-	-
-DIFFKEY	38.54	0.301	0.148	-	-
-Stage-1	43.21	0.347	0.191	-	-

uses LSTM or GPT-2, the methods are called SkelGen_{LSTM} and SkelGen_{GPT-2}.

RAG: End-to-end hybrid model that uses BERT as the neural retriever and BART as the generator. RAG is designed as a general purpose retrieval-augmented generation system, and so uses Wikipedia as the knowledge source [57].

C. EXPERIMENTAL SETTINGS

We set word embedding dimension to 512, transformer hidden state dimension to 1024, and dropout rate to 0.3. We use a vocabulary size of 30,004 (30,000 words and 4 special symbols). For SAKE, the number of retrieved results $K = 2$, the projected dimension $H = 5$ and $\eta = 0.2$. We use a batch size of 512 and train for 30 epochs for all baselines and our model, and halve the learning rate when development performance worsens. We use the base model for BERT, and the uncased variant for English. All baseline/benchmark models use their default recommended hyper-parameter configuration.

D. RESULTS

a: Overall Experiments.

Table 2 presents the full results, where the top, middle and bottom sub-tables are Cornell, Weibo and Douban results, respectively. Generally, we see that the hybrid systems are better models compared to pure generation and retrieval systems. Our model shows strong performance: it substantially outperforms most baselines and benchmark systems in ROUGE and relevance scores across all 3 datasets, creating a new state-of-the-art.

In terms of human evaluation, the generated responses of our model are also more fluent and relevant than all generation and hybrid systems, although they are admittedly less fluent compared to retrieval systems (as their output are human-written responses). For diversity, we see a similar trend where retrieval systems tend to have an upper hand, although when compared to non-retrieval systems, our model outperforms all these systems by a comfortable margin.

b: Ablation Study.

To study the influence of the individual components (e.g. the impact of the number retrieved results K , SAKE and two-stage transformer) in our system, we perform several ablation studies based on Douban (test set). All studies present the

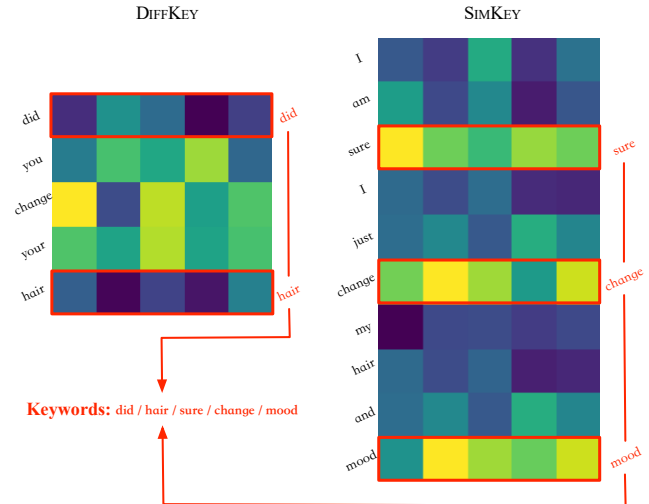


FIGURE 4. Multi-source alignment output $Y^{[N]}$ produced by SAKE (Equation 2) for extracting DIFFKEY from the query (left) and SIMKEY from the first retrieved response (right). Darker colour indicates lower magnitude/strength.

average scores of the different variants of ROUGE, relevance and diversity.

We assess the effectiveness of our keyword extraction module by removing either SimKey and DiffKey and present the results in Table 3. It appears that removing either keywords degrades the response substantially across all metrics, indicating the importance of both keywords. That said, SimKey seems to be marginally more effective than DiffKey in guiding the response generation.

We also test the impact of ordering the keywords by creating a variant where we remove the stage-1 transformer and feed DiffKey and SimKey to stage-2 transformer without ordering information (i.e. positional embeddings are not added to the input). Results are in the last row of Table 3. We see a dip in performance across all metrics, suggesting it is beneficial to decompose the generation task into a two-step process where we predict the order of the keywords before using them to drive the response generation.

c: Qualitative Analysis.

We present the generated responses and the retrieved conversations by our system for two queries from the Cornell Movie-Dialog corpus (top) and Douban (bottom) Table 4. We can see the retrieved utterance and response pairs provide additional context for the query, and the generated responses are largely driven by the extracted keywords (SimKey and DiffKey).

To qualitatively understand the output of the SAKE through multi-source alignment, we present the alignment output $Y_{DiffKey}^{[N]}$ and $Y_{SimKey}^{[N]}$ (produced by Equation 2) in Figure 4. The query and top-2 retrieved utterance-response pairs are presented at the top of Table 4. Here we can see that words such as *did* and *hair* are selected as DiffKey from the query due to their low alignment strength with the retrieved utterances (see Table 4), while *sure*, *change* and *mood* are

TABLE 4. Generated responses and retrieved conversations for two utterances from the Cornell Movie-Dialog corpus (top) and Douban (bottom). Keywords indicate the DIFFKEY and SIMKEY extraction. (via *argmin* and *argmax*). The performance in English and Chinese single-turn dialogue is intuitively displayed.

	Query	Did you change your hair
Retrieved	Utterance1	What happened about your hair
	Utterance2	Are you sure you won't change your mind
	Response1	I am sure I just change my hair and mood
	Response2	Sure, my mind change with mood
	Keywords	did / hair / sure / change / mood
	Generated Response	Sure my hair and mood have changed
	Query	我明天想出去晒晒 I want to go out tomorrow
Retrieved	Utterance1	我想出去晒晒太阳 I want going out to sunbathe
	Utterance2	我想明天出门因为在家太久 I'm going out tomorrow because I've been at home too long
	Response1	明天是个好天气 It's a fine day tomorrow
	Response2	天气预报说明天是多云没有太阳 The weather forecast says it will be cloudy tomorrow and there will be no sun
	Keywords	明天 / 出去 / 好天气 tomorrow / go out / fine weather
	Generated Response	明天会是个适合出去的好天气 Tomorrow is a fine day to go out

extracted as SimKey from the first retrieved response as these words are also mentioned in the second retrieved response.

Seeing function words such as *did* and *sure* are being selected as keywords (which seem counter-intuitive), we did another experiment where we use a stopword list to filter these words in SAKE. We found that the results worsen, and hypothesise that these words could be more important than they appear as we are working with response generation for casual conversation/dialogue.

V. DISCUSSION

Open domain response generation is an arguably difficult task. While studies find that retrieval-based methods are unable to produce the right response even with massive conversation corpora, generation-based methods also have their drawbacks where they struggle to produce creative outputs. To address this, hybrid methods that combine both have been explored and with some success.

Wu et al. [51] found that the retrieved results provide a good starting point for generation because it is grammatical and informative. They extract words from these results and show that these keywords help improve relevance, diversity and originality of the generated responses, and our work is inspired by them. Our experiments found a similar finding, showing that this is a promising direction. That said, our study is different to theirs in that we focus more on the keyword extraction process where we explore extracting two different

types of keywords (DiffKey and SimKey). Ablation results demonstrate that both DiffKey and SimKey have contributed to improving the generation process. Also, our work is unique in how we introduce a two-stage transformer to predict the order of the keywords and generate the response by “filling in the blanks”, which does not involve any decoding.

As single-turn dialogue systems becoming more mature, the next step is to develop multi-turn dialogue systems. Multi-turn conversation is challenging, as the model would need to consider beyond the last utterance as context to generate responses. Our keywords extraction module can be potentially extended here by framing it as a time series problem using the dialogue history. With the advent of very large pre-train language models such as GPT-3, MT-NLG [58] and PaLM [59], it'd also be interesting to explore using them to replace BERT. We suspect this can potentially bring in further improvement, as these models have demonstrated outstanding performance across a number of NLP tasks.

VI. CONCLUSION

We introduce an end-to-end response generation model that extracts keywords from retrieved conversations to guide the response generation. Our system combines the benefits of retrieval and generation methods, and utilises modern pre-trained language models and their attention mechanism for keyword extraction and response generation. We evaluate our system on 3 datasets over two languages (English and Chinese), and demonstrate that it outperforms benchmark systems in ROUGE, relevance scores and human evaluation, creating a new state-of-the-art. In future work, we will further consider a sequence set of keywords extracted from the contextual information to guide the response generation in multi-turn dialogue task. Some large language models are also considered to replace the generation model.

APPENDIX A HUMAN EVALUATION

We use the same methodology to collect human annotations for all three datasets. For each dataset, we randomly sample 200 generated dialogues (original query+generated response) and divide them into four batches (50 dialogues each batch). Sixteen native speakers (Chinese or English depending on the dataset) were invited to rate the generated responses on a 4-point scale;⁹ Table 5 presents an example. The judges are broken into four groups, and each batch of dialogues is annotated by two groups of judges. For each dialogue, we have 2 ratings for each aspect (fluency or relevance) and we take the average as the final rating. Within a batch, if the ratings differ substantially between the two groups of judges, a third group of judges will be invited to annotate the batch. The judges do not have access to the ground-truth response, and see only the query and system-generated responses. Each worker is paid USD \$0.15 for annotating a query. For fluency evaluation, the 4-point scale is described as follows:

⁹We use the Tencent online document platform for conducting the crowdsourcing experiments: <https://docs.qq.com/>

- 1: *hard to read;*
- 2: *not quite fluent and has several grammatical errors;*
- 3: *fluent response with few errors*
- 4: *fluent response without errors.*

For relevancy:

- 1: *totally irrelevant;*
- 2: *marginally relevant;*
- 3: *somewhat relevant but not directly related to the query*
- 4: *relevant.*

TABLE 5. An example of scoring criteria.

Original Query	Generated Response	
	Fluency	Relevance
1	location I course not.	I like apple best.
2	I of course for it.	Shanghai is the most international city in China.
3	No preference for I.	Shanghai is good for me.
4	Of course, I have no preference for location.	Of course, I have no preference for location.

原始问句	生成回答	
	流畅度	相关性
1	这好你是。	你今天看起来真漂亮。
2	我这好想好啊。	今天天气真好!
3	好啊, 我这天想。	我今天想吃面。
4	好, 我正好想出去吃。	没问题我们出去吃。

REFERENCES

- [1] Alessandro Sordani, Michel Galley, Michael Auli, Chris Brockett, Yangfeng Ji, Margaret Mitchell, Jian-Yun Nie, Jianfeng Gao, and Bill Dolan, "A neural network approach to context-sensitive generation of conversational responses," *arXiv preprint arXiv:1506.06714*, 2015.
- [2] Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan, "A diversity-promoting objective function for neural conversation models," *arXiv preprint arXiv:1510.03055*, 2015.
- [3] Oriol Vinyals and Quoc Le, "A neural conversational model," *arXiv preprint arXiv:1506.05869*, 2015.
- [4] Julian V Serban, Alessandro Sordani, Yoshua Bengio, and Joelle Pineau, "Building end-to-end dialogue systems using generative hierarchical neural network models," *arXiv preprint arXiv:1507.04808*, 2015.
- [5] Deng Cai, Yan Wang, Victoria Bi, Zhaopeng Tu, Xiaojiang Liu, Wai Lam, and Shuming Shi, "Skeleton-to-response: Dialogue generation guided by retrieval memory," *arXiv preprint arXiv:1809.05296*, 2018.
- [6] Deng Cai, Yan Wang, Wei Bi, Zhaopeng Tu, Xiaojiang Liu, and Shuming Shi, "Retrieval-guided dialogue response generation via a matching-to-generation framework," in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, 2019, pp. 1866–1875.
- [7] Zhiliang Tian, Wei Bi, Dongkyu Lee, Lanqing Xue, Yiping Song, Xiaojiang Liu, and Nevin L Zhang, "Response-anticipated memory for on-demand knowledge integration in response generation," *arXiv preprint arXiv:2005.06128*, 2020.
- [8] Shen Gao, Xiuying Chen, Zhaochun Ren, Dongyan Zhao, and Rui Yan, "Meaningful answer generation of e-commerce question-answering," *ACM Transactions on Information Systems (TOIS)*, vol. 39, no. 2, pp. 1–26, 2021.
- [9] Zhiliang Tian, Wei Bi, Xiaopeng Li, and Nevin L Zhang, "Learning to abstract for memory-augmented conversational response generation," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, 2019, pp. 3816–3825.
- [10] Chris Donahue, Mina Lee, and Percy Liang, "Enabling language models to fill in the blanks," *arXiv preprint arXiv:2005.05339*, 2020.
- [11] Chin-Yew Lin, "Rouge: A package for automatic evaluation of summaries," in *Text summarization branches out*, 2004, pp. 74–81.
- [12] Baotian Hu, Zhengdong Lu, Hang Li, and Qingcai Chen, "Convolutional neural network architectures for matching natural language sentences," *Advances in neural information processing systems*, vol. 27, 2014.
- [13] Md Atabuzzaman, Md Shajalal, M Elius Ahmed, Masud Ibn Afjal, and Masaki Aono, "Leveraging grammatical roles for measuring semantic similarity between texts," *IEEE Access*, vol. 9, pp. 62972–62983, 2021.
- [14] Zhixue Jiang, Chengying Chi, and Yunyun Zhan, "Research on medical question answering system based on knowledge graph," *IEEE Access*, vol. 9, pp. 21094–21101, 2021.
- [15] Hao Wang, Zhengdong Lu, Hang Li, and Enhong Chen, "A dataset for research on short-text conversations," in *Proceedings of the 2013 conference on empirical methods in natural language processing*, 2013, pp. 935–945.
- [16] Zhuosheng Zhang, Jiangtong Li, Pengfei Zhu, Hai Zhao, and Gongshen Liu, "Modeling multi-turn conversation with deep utterance aggregation," *arXiv preprint arXiv:1806.09102*, 2018.
- [17] Donghyun Choi, Myeongcheol Shin, Eunngyung Kim, and Dong Ryeol Shin, "Adaptive batch scheduling for open-domain question answering," *IEEE Access*, vol. 9, pp. 112097–112103, 2021.
- [18] Ilya Sutskever, Oriol Vinyals, and Quoc V Le, "Sequence to sequence learning with neural networks," *Advances in neural information processing systems*, vol. 27, 2014.
- [19] Sepp Hochreiter and Jürgen Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.
- [20] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio, "Neural machine translation by jointly learning to align and translate," in *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun, Eds., 2015.
- [21] Jiatao Gu, Zhengdong Lu, Hang Li, and Victor OK Li, "Incorporating copying mechanism in sequence-to-sequence learning," *arXiv preprint arXiv:1603.06393*, 2016.
- [22] Lei Shen, Haolan Zhan, Xin Shen, and Yang Feng, "Learning to select context in a hierarchical and global perspective for open-domain dialogue generation," in *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2021, pp. 7438–7442.
- [23] Yu Wu, Wei Wu, Chen Xing, Can Xu, Zhoujun Li, and Ming Zhou, "A sequential matching framework for multi-turn response selection in retrieval-based chatbots," *Computational Linguistics*, vol. 45, no. 1, pp. 163–197, 2019.
- [24] Xiao Sun and Bingbing Ding, "Neural network with hierarchical attention mechanism for contextual topic dialogue generation," *IEEE Access*, 2022.
- [25] Ziming Li, Julia Kiseleva, and Maarten de Rijke, "Improving response quality with backward reasoning in open-domain dialogue systems," in *Proceedings of the 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, 2021, pp. 1940–1944.
- [26] Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon, "Unified language model pre-training for natural language understanding and generation," *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [27] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al., "Language models are few-shot learners," *Advances in neural information processing systems*, vol. 33, pp. 1877–1901, 2020.

- [28] Jiwei Li, Will Monroe, Alan Ritter, Michel Galley, Jianfeng Gao, and Dan Jurafsky, "Deep reinforcement learning for dialogue generation," *arXiv preprint arXiv:1606.01541*, 2016.
- [29] Haochen Liu, Wentao Wang, Yiqi Wang, Hui Liu, Zitao Liu, and Jiliang Tang, "Mitigating gender bias for neural dialogue generation with adversarial learning," *arXiv preprint arXiv:2009.13028*, 2020.
- [30] Yanxiang Ling, Fei Cai, Xuejun Hu, Jun Liu, Wanyu Chen, and Honghui Chen, "Context-controlled topic-aware neural response generation for open-domain dialog systems," *Information Processing & Management*, vol. 58, no. 1, pp. 102392, 2021.
- [31] Heng-Da Xu, Xian-Ling Mao, Zewen Chi, Fanshu Sun, Jingjing Zhu, and Heyan Huang, "Generating informative dialogue responses with keywords-guided networks," in *CCF International Conference on Natural Language Processing and Chinese Computing*. Springer, 2021, pp. 179–192.
- [32] E Haihong, Zecheng Zhan, and Meina Song, "Table-to-dialog: Building dialog assistants to chat with people on behalf of you," *IEEE Access*, vol. 8, pp. 102313–102320, 2020.
- [33] Jianfeng Yu, Yan Yang, Chengcai Chen, Liang He, and Zhou Yu, "Dafa: Dialog system domain adaptation with a filter and an amplifier," *IEEE Access*, vol. 8, pp. 45041–45049, 2020.
- [34] Yan Xu, Etsuko Ishii, Samuel Cahyawijaya, Zihan Liu, Genta Indra Winata, Andrea Madotto, Dan Su, and Pascale Fung, "Retrieval-free knowledge-grounded dialogue response generation with adapters," *arXiv preprint arXiv:2105.06232*, 2021.
- [35] Tiancheng Zhao, Ran Zhao, and Maxine Eskenazi, "Learning discourse-level diversity for neural dialog models using conditional variational autoencoders," *arXiv preprint arXiv:1703.10960*, 2017.
- [36] Soujanya Poria, Navonil Majumder, Rada Mihalcea, and Eduard Hovy, "Emotion recognition in conversation: Research challenges, datasets, and recent advances," *IEEE Access*, vol. 7, pp. 100943–100953, 2019.
- [37] Mengshi Duan, Qing Li, and Le Xiao, "Topic-extended emotional conversation generation model based on joint decoding," *IEEE Access*, vol. 9, pp. 89934–89940, 2021.
- [38] Mengjuan Liu, Xiaoming Bao, Jiang Liu, Pei Zhao, and Yuchen Shen, "Generating emotional response by conditional variational auto-encoder in open-domain dialogue system," *Neurocomputing*, vol. 460, pp. 106–116, 2021.
- [39] Yiping Song, Rui Yan, Xiang Li, Dongyan Zhao, and Ming Zhang, "Two are better than one: An ensemble of retrieval-and generation-based dialog systems," *arXiv preprint arXiv:1610.07149*, 2016.
- [40] Gaurav Pandey, Danish Contractor, Vineet Kumar, and Sachindra Joshi, "Exemplar encoder-decoder for neural conversation generation," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, 2018, pp. 1329–1338.
- [41] Ning Miao, Hao Zhou, Lili Mou, Rui Yan, and Lei Li, "Cgmh: Constrained sentence generation by metropolis-hastings sampling," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2019, vol. 33, pp. 6834–6842.
- [42] Ziqiang Cao, Wenjie Li, Sujian Li, and Furu Wei, "Retrieve, rerank and rewrite: Soft template based neural summarization," in *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics*, 2018, pp. 152–161.
- [43] Amirhossein Kazemnejad, Mohammadreza Salehi, and Mahdieh Soleymani Baghshah, "Paraphrase generation by learning how to edit from samples," in *proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, 2020, pp. 6010–6021.
- [44] Kulothunkan Palasundram, Nurfadhilina Mohd Sharef, Khairul Azhar Kasmiran, and Azreen Azman, "Seq2seq++: A multitasking-based seq2seq model to generate meaningful and relevant answers," *IEEE Access*, vol. 9, pp. 164949–164975, 2021.
- [45] Sixing Wu, Ying Li, Dawei Zhang, and Zhonghai Wu, "Improving knowledge-aware dialogue response generation by using human-written prototype dialogues," in *Findings of the Association for Computational Linguistics: EMNLP 2020*, 2020, pp. 1402–1411.
- [46] Jing Xu, Arthur Szlam, and Jason Weston, "Beyond goldfish memory: Long-term open-domain conversation," *arXiv preprint arXiv:2107.07567*, 2021.
- [47] Yingce Xia, Fei Tian, Lijun Wu, Jianxin Lin, Tao Qin, Nenghai Yu, and Tie-Yan Liu, "Deliberation networks: Sequence generation beyond one-pass decoding," *Advances in neural information processing systems*, vol. 30, 2017.
- [48] Juncen Li, Robin Jia, He He, and Percy Liang, "Delete, retrieve, generate: a simple approach to sentiment and style transfer," *arXiv preprint arXiv:1804.06437*, 2018.
- [49] Liu Yang, Junjie Hu, Minghui Qiu, Chen Qu, Jianfeng Gao, W Bruce Croft, Xiaodong Liu, Yelong Shen, and Jingjing Liu, "A hybrid retrieval-generation neural conversation model," in *Proceedings of the 28th ACM international conference on information and knowledge management*, 2019, pp. 1341–1350.
- [50] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova, "Bert: Pre-training of deep bidirectional transformers for language understanding," *arXiv preprint arXiv:1810.04805*, 2018.
- [51] Yu Wu, Furu Wei, Shaohan Huang, Zhoujun Li, and Ming Zhou, "Response generation by context-aware prototype editing," *arXiv preprint arXiv:1806.07042*, 2018.
- [52] Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov, "Multimodal transformer for unaligned multimodal language sequences," in *Proceedings of the conference. Association for Computational Linguistics. Meeting. NIH Public Access*, 2019, vol. 2019, p. 6558.
- [53] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin, "Attention is all you need," *Advances in neural information processing systems*, vol. 30, 2017.
- [54] Cristian Danescu-Niculescu-Mizil and Lillian Lee, "Chameleons in imagined conversations: A new approach to understanding coordination of linguistic style in dialogs," *arXiv preprint arXiv:1106.3077*, 2011.
- [55] Chia-Wei Liu, Ryan Lowe, Iulian V Serban, Michael Noseworthy, Laurent Charlin, and Joelle Pineau, "How not to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation," *arXiv preprint arXiv:1603.08023*, 2016.
- [56] Ryan Lowe, Nissam Pow, Iulian Serban, and Joelle Pineau, "The ubuntu dialogue corpus: A large dataset for research in unstructured multi-turn dialogue systems," *arXiv preprint arXiv:1506.08909*, 2015.
- [57] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al., "Retrieval-augmented generation for knowledge-intensive nlp tasks," *Advances in Neural Information Processing Systems*, vol. 33, pp. 9459–9474, 2020.
- [58] Shaden Smith, Mostofa Patwary, Brandon Norick, Patrick LeGresley, Samyam Rajbhandari, Jared Casper, Zhun Liu, Shrimai Prabhumoye, George Zerveas, Vijay Korthikanti, et al., "Using deepspeed and megatron to train megatron-turing nlg 530b, a large-scale generative language model," *arXiv preprint arXiv:2201.11990*, 2022.
- [59] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al., "Palm: Scaling language modeling with pathways," *arXiv preprint arXiv:2204.02311*, 2022.



CHANG SHU received his M.Sc. degree in Machine Learning and Data Mining at University of Bristol in 2011. He is currently a PhD student at the University of Nottingham Ningbo China, and working in Ping An Technology (Shenzhen) Co., Ltd. He has published some research papers in conferences and journals including NAACL, TACL, COLING, IJCNN, etc.

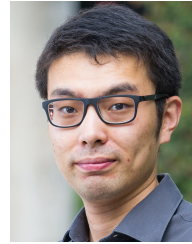


ZIJIAN ZHANG received the B.E. degree from Ocean University of China in 2019 and M.E. degree from Tongji University in 2022. He is now an Engineer with Meituan-Tianping, Shanghai, China. His current research interests include the areas of natural language processing and multimodal learning.



understanding, speech recognition, etc.

YOUXIN CHEN received his M.Sc. degree in Signal and Information Processing in 2003 at Tsinghua University, Beijing, China. From 2012 to 2016, he was the Principal Engineer, R&D Director in China Communications Research Institute of Samsung Electronics. He joined the department of Bigdata & AI, Ping An Technology (Shenzhen) Co., Ltd., as a Senior Director & Expert of AI in 2019. His interests include the area of pattern



machine learning, computer vision as well as natural language processing.

ZHENG LU received his PhD from the National University of Singapore in 2011. He is currently an Assistant Professor at the University of Nottingham Ningbo China. Prior to his current position, he was a Postdoctoral Research Fellow at the University of Texas at Austin, and an Assistant Professor at the City University of Hong Kong. He has published over 30 research papers in conferences and journals including CVPR, TPAMI, IJCV, etc. His current research interests include the area of

...



since 1995, covering a broad range of application areas such as healthcare, autonomous driving, 3D printing and display, biometrics, web search, and finance. He is now leading research and development in AI-related technologies and their applications on finance, healthcare, and smart-city in Ping An.

JING XIAO is the Group Chief Scientist of Ping An Insurance (Group) Company of China, LTD. He received his PhD degree from School of Computer Science, Carnegie Mellon University, and has published over 180 academic papers and 102 granted US patents. Before joining Ping An, he worked as Principal Applied Scientist Lead in Microsoft Corp. and Manager of Algorithm Group in Epson Research and Development, Inc. He started R&D in artificial intelligence and related fields



text generation and misinformation detection.

JEY HAN LAU received his PhD from the University of Melbourne in 2013. He is currently a Lecturer at the University of Melbourne. Prior to joining the University of Melbourne, he spent over 3.5 years as an industry scientist at IBM Research. Over his career he has published over 60 peer-reviewed papers with a H-index of 24. Jey Han's research interest is in natural language processing, and his research guided by a diverse flavour of applications, e.g. topic models, lexical semantics,



QIAN ZHANG received her B.Sc. from the Hong Kong Baptist University in 2009. She obtained an M.Sc. in Computer Science at University of Nottingham in 2010, where she obtained her Ph.D. in 2015. She joined University of Nottingham Ningbo China in 2016 as an assistant professor.