

Multimodal Dialogue Response Generation

Qingfeng Sun¹ Yujing Wang² Can Xu¹ Kai Zheng¹ Yaming Yang²
Huang Hu¹ Fei Xu¹ Jessica Zhang¹ Xiubo Geng¹ Daxin Jiang^{1*}

¹Microsoft STC Asia ²Microsoft Research Asia

{qins, yujwang, caxu, zhengkai, yayaming, huahu, fexu,
jessicaz, xigeng, djiang}@microsoft.com

Abstract

Responding with image has been recognized as an important capability for an intelligent conversational agent. Yet existing works only focus on exploring the multimodal dialogue models which depend on retrieval-based methods, but neglecting generation methods. To fill in the gaps, we first present a new task: multimodal dialogue response generation (MDRG) - given the dialogue context, one model needs to generate a text or an image as response. Learning such a MDRG model often requires multimodal dialogues containing both texts and images which are difficult to obtain. Motivated by the challenge in practice, we consider MDRG under a natural assumption that only limited training examples are available. Under such a low-resource setting, we devise a novel conversational agent, Divter, in order to isolate parameters that depend on multimodal dialogues from the entire generation model. By this means, the major part of the model can be learned from a large number of text-only dialogues and text-image pairs respectively, then the whole parameters can be well fitted using just a few training examples. Extensive experiments demonstrate our method achieves state-of-the-art results in both automatic and human evaluation, and can generate informative text and high-resolution image responses.

1 Introduction

With the development of instant messaging technology in the recent decades, the intermediary of online conversation has also changed from pure text to a variety of visual modalities (e.g., image, gif animation, short video). Similar to communicating by the messenger tools (e.g., Facebook, WhatsApp, WeChat) in reality, an excellent intelligent conversational agent should not only be able to converse freely with plain text, but also have the ability to perceive and share the real visual physical world.

* Corresponding author.

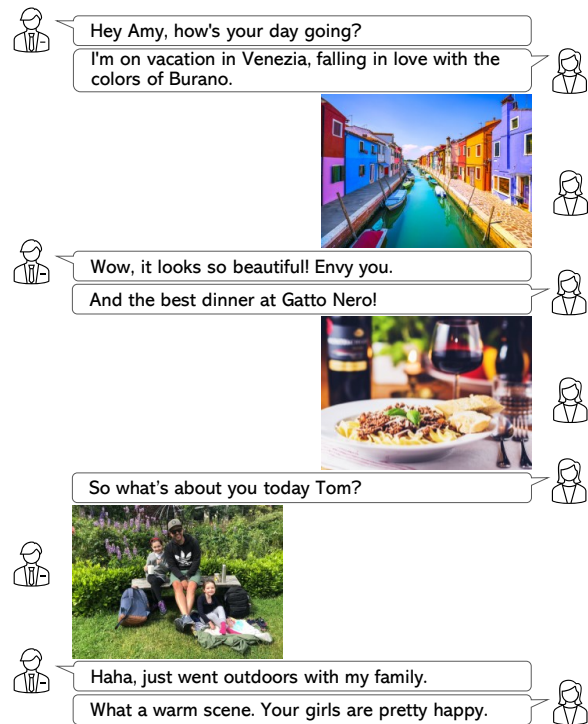


Figure 1: An example of human conversations. They are talking about vacation and outdoors with both text and various images.

Although recently some large-scale pre-trained text-only dialogue generation models, such as DialoGPT (Zhang et al., 2020), Blender (Roller et al., 2021), Meena (Adiwardana et al., 2020), have shown excellent performance, they still cannot rely exclusively on plain text to completely simulate the rich experience of visual perception. Recently, various vision-language tasks have been introduced and attracted widespread attention, such as visual question answering (Ren et al., 2015; Lu et al., 2016; Anderson et al., 2018; Li et al., 2019a; Huang et al., 2020), image captioning (Xu et al., 2015; Anderson et al., 2016; Ghanimifard and Dobnik, 2019; Cornia et al., 2020), image-grounded dialogue (Das et al., 2017; Yang et al., 2021; Agarwal et al., 2020; Qi et al., 2020; Chen et al., 2021; Liang et al., 2021).

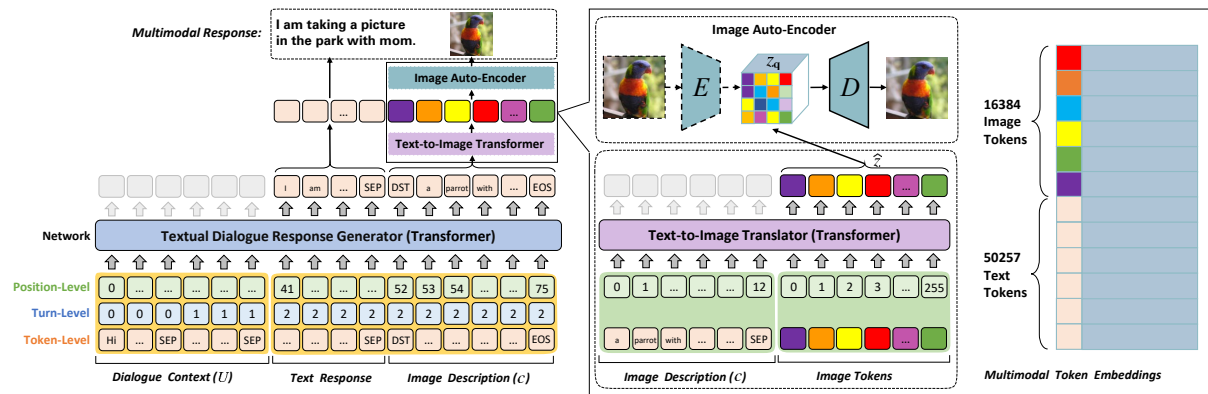


Figure 2: The overview of our multimodal dialogue response generation model. The Textual Dialogue Response Generator takes the text dialogue context U as input and generates a sequence contains text response and a image description (e.g., “a parrot with red belly and green back is standing on the railing.”). With the description as a condition, the Text-to-Image Translator generates image representation \hat{z} . The Image Decoder \mathcal{V}_D reconstructs \hat{z} to a realistic and consistent high resolution image.

Specifically, in human conversations, the images can easily show rich visual perception, which is hard to be expressed by plain text. As the example shown in Figure 1, images are required in at least three circumstances: (i) the other speaker has little knowledge (e.g., colorful Burano, in the 1st image) of the objects only you had seen; (ii) to share more details (e.g., red wine and pasta, in the 2nd image) of the objects even you have common knowledge of them; (iii) to express your emotions (e.g., happy, in the 3rd image) about a specific event. An existing related task is photo sharing (Zang et al., 2021), which aims to select and share the image based on the textual context, is a challenging task that requires models to understand the background story which complemented by human imaginations, rather than to locate related visual objects or explicitly mention main visible content in the image as the previous works do. Zang et al. (2021) propose a retrieval-based method to resolve the above challenge. However, the performance of the retrieval-based method is limited in specific domains by the size of the pre-constructed conversational history repository, especially for long-tail contexts that are not covered in the history, where the set of image responses of a retrieval system is also fixed. On the other hand, a better way is to generate a new one accordingly.

In this paper, we formulate a new problem: **Multimodal Dialogue Response Generation (MDRG)**, that is, given the dialogue context, the model should not only generate a pure text response but also have the capacity to generate a multimodal response (e.g., containing both image and text). We

argue that there are still some hindrances to application, since (1) the sophisticated neural end-to-end architecture will overfit to very few well-annotated training data (e.g., a few existing 10k multimodal dialogues). Evidence is that when discussing the topics outside the training data domain, its performance drops dramatically; and (2) as human effort is expensive, it is not easy to collect enough training data for a new domain. Based on the above facts, we take a step further to extend the assumption of MDRG to a low-resource setting where only a few multimodal dialogues are available.

To tackle the above challenges, our key idea is to make parameters that rely on multimodal dialogues small and independent by disentangling textual response generation and image response generation, and thus we can learn the major part of the generation model from text-only dialogues and <image description, image> pairs that are much easier to be obtained. Specifically, we present **Divter**, a novel conversational agent powered by large-scale visual world experiences. As shown in Figure 2, our Divter is made up of two Transformer-based (Vaswani et al., 2017a) components: a multimodal dialogue response generator, and a text-to-image translator. Divter takes the dialogue context as input, then generates a textual sequence which may contains a text response or a textual image description or both of them. The text-to-image translator takes above image description as condition, then generates a realistic and consistent high resolution image. Both components are independent with the opposite knowledge, and thus can be pre-trained using a large number of text-only dialogues and

the <image description, image> pairs respectively. The end-to-end Divter depends on the multimodal dialogues constructed as the tuple: (*dialogue context, text response / <image description, image>*), but the joint learning and estimation of the two components just require a few training examples depending on specific domains.

Contributions of this work are three-fold:

- To the best of our knowledge, it is the first work on the multimodal dialogue response generation. We explore the task under a low-resource setting where only a few multimodal dialogues are assumed available.
- We present Divter, a novel conversational agent which can effectively understand dialogue context and generate informative text and high-resolution image responses.
- Extensive experiments on PhotoChat Corpus (Zang et al., 2021) indicate the effectiveness of Divter, it achieves a significant improvement with pure text dialogue generation model and retrieval-based image sharing method.

2 Related Work

2.1 Textual Dialogue Response Generation

End-to-end response generation for textual open-domain dialogues is inspired by the successful application of neural sequence-to-sequence models on machine translation (Sutskever et al., 2014). On top of the basic architecture (Shang et al., 2015; Vinyals and Le, 2015), the vanilla encoder-decoder method is widely extended to address the critical challenges in open-domain dialogue systems, including improving the diversity of responses (Li et al., 2016a; Zhao et al., 2017; Tao et al., 2018), modeling conversation contexts (Serban et al., 2016; Xing et al., 2017; Zhang et al., 2019; Zhao et al., 2020), controlling attributes of responses (See et al., 2019; Zhou et al., 2018; Xu et al., 2019), biasing responses to some specific personas (Li et al., 2016b; Zhang et al., 2018), incorporating extra knowledge into generation (Dinan et al., 2019; Ghazvininejad et al., 2018; Kim et al., 2020; Li et al., 2020), and building general pre-trained agents (Adiwardana et al., 2020; Zhang et al., 2020; Roller et al., 2021; Qi et al., 2021). Different from the previous works on open-domain dialogue response generation that converse freely with plain text, our work lies in the research of multimodal response generation.

2.2 Text-to-Image Generation

In the research of text-to-image generation, various works have been extensively studied. Mansimov et al. (2016) shown the Draw generative model (Gregor et al., 2015) could generate images from natural language descriptions. Reed et al. (2016) proposed a generative adversarial network to improve the image fidelity. Then some improvement methods continue to optimize the generation architecture, such as stacked generators (Zhang et al., 2017), attentional network (Xu et al., 2018), and extra knowledge (Li et al., 2019b). Nguyen et al. (2017) provided a unified probabilistic interpretation of related activation maximization methods to produce high-quality images at higher resolutions. Separately, Cho et al. (2020) used uniform masking with a large range of masking ratios and align the suitable pre-training datasets to the proper objectives. More recently, Ramesh et al. (2021) and (Ding et al., 2021) adopt transformer-based methods which autoregressively model the text and image tokens as a single stream of data. For this multimodal response generation scenario, we use the textual image description to bridge above textual dialogue generation and text-to-image generation models, where the image description is the output of the former and input of the latter in a low-resource setting.

3 Problem Formalization

Suppose that we have dataset $\mathcal{D}_S = \{(U_i, R_i)\}_{i=1}^n$, where $\forall i \in \{1, \dots, n\}$, $U_i = \{u_{i,1}, \dots, u_{i,n_i}\}$ is the dialogue context with $u_{i,j}$ the j -th utterance, and R_i is the response regarding to U_i . $u_{i,j}$ and R_i could contain two modalities: text, and image. The goal is to learn a **generation** model $P(R|U; \theta)$ (θ denotes the parameters of the model) with \mathcal{D}_S . Thus, given a new dialogue context U , one can generate a multimodal response R following $P(R|U; \theta)$.

4 Approach

This section first formulates the unified tokenization method for multimodal dialogues. We then introduce the two important components in our proposed multimodal dialogue response generation model (**Divter**) under low-resource scenario, including (i) textual dialogue response generator; (ii) text-to-image translator. Figure 2 shows the overall of our **Divter**.

4.1 Multimodal Tokenization

To learn a multimodal generation model, we should first model the unified representations of both text and image. Inspired by the success of DALLE (Esser et al., 2020) and VQGAN (Ramesh et al., 2021), to utilize the highly expressive transformer architecture for text-to-image generation, we need to express an image in the form of a sequence, similar to what we usually do for pure text tokenization.

4.1.1 Text Tokenization

The tokenization for text is already well-studied, e.g., BPE (Gage, 1994). This work uses 50257 BPE-encoded tokens and distributed embedding of Transformer architecture (Vaswani et al., 2017b) to model the texts in a dialogue.

4.1.2 Image Tokenization

The tokenizer for image is a discrete Auto-Encoder (VQGAN¹) \mathcal{V} as shown in Figure 2. \mathcal{V} uses an encoder \mathcal{V}_E to compress each image r^v of shape $H \times W \times 3$ into \hat{z} of shape $h \times w \times d_z$, then each vector of dimension d_z would be quantized to its closest embedding z_k in a learned, discrete codebook $\mathcal{Z} = \{z_k\}_{k=1}^K \in \mathbb{R}^{d_z}$ under the action of element-wise quantization $\mathbf{q}(\cdot)$

$$z_{\mathbf{q}} = \mathbf{q}(\hat{z}) := \left(\arg \min_{z_k \in \mathcal{Z}} \|\hat{z}_{ij} - z_k\| \right) \in \mathbb{R}^{h \times w \times d_z} \quad (1)$$

Thus r^v can be represented by a spatial collection of codebook entries $z_{\mathbf{q}} \in \mathbb{R}^{h \times w \times d_z}$. The decoder \mathcal{V}_D maps the $z_{\mathbf{q}}$ back to a image \hat{r}^v to reconstruct the input. In this work, $H = W = 256$, $h = w = 16$, $K = 16384$, $d_z = 256$. The learning details of \mathcal{V} and \mathcal{Z} could be found in Ramesh et al. (2021).

4.2 Low-resource Learning Model

Learning an effective multimodal generation model with a single sequence-to-sequence model often requires a large number of training instances. However, only very few multimodal dialogues are available due to the privacy restrictions on social media and the expensive human effort. On the other hand, as shown in Figure 3, there existed a large number of open source text-only dialogues (e.g. Reddit comments², formulated as $\mathcal{D}_C = \{(U_i, r_i^e)\}_{i=1}^N$ with (U_i, r_i^e) a <text dialogue context, text response> pair), and a large number of <image description, image> pairs (e.g. YFCC100M (Thomee

¹<https://github.com/CompVis/taming-transformers>

²<https://files.pushshift.io/reddit/>

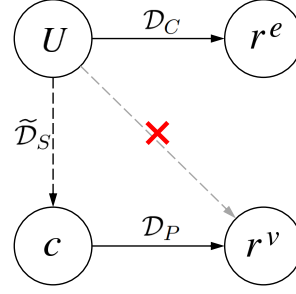


Figure 3: Abstract Logic of the proposed approach. Solid lines mean that there exists large-scale training set to pre-train the generation model, while dotted lines mean that only very few training instances are available, “x” means bad generation quality.

et al., 2016), formulated as $\mathcal{D}_P = \{(c_j, r_j^v)\}_{j=1}^M$ with (c_j, r_j^v) a <textual image-description, image> pair). Based on the above facts and the low-resource challenges on MDRG task, we adapt to incorporate generative text-to-image translation into text-only open domain dialogue response generation. More specifically: (i) if the multimodal dialogue context contains an image, we replace the image with its description to form a text-only context, and take this context as the input of the text-only dialogue generation model \mathcal{G} (pre-trained with \mathcal{D}_C); (ii) if we need to generate an image as a part of response, we could first generation a textual description with \mathcal{G} , then adopt a text-to-image translator module \mathcal{F} (pre-trained with \mathcal{D}_P) to translate the description to a synonymous image. To bridge \mathcal{G} and \mathcal{F} , we further extend the formalization of \mathcal{D}_S to a new $\tilde{\mathcal{D}}_S$ in which each image r^v is paired with its textual description c . Both the (i) and (ii) actions can be independently learned, which becomes the key to aiding the small $\tilde{\mathcal{D}}_S$ with the large \mathcal{D}_C and \mathcal{D}_P .

By this means, the current goal is to learn a generation model $P(R|U; \theta)$ with $\mathcal{D} = \{\tilde{\mathcal{D}}_S, \mathcal{D}_C, \mathcal{D}_P\}$. With the pre-trained \mathcal{G} and \mathcal{F} available, we finally use $\tilde{\mathcal{D}}_S$ to jointly finetune \mathcal{G} and \mathcal{F} to obtain the capacity of generating multimodal responses.

Figure 2 illustrates the architecture of our model. The model is made up of two components: a textual dialogue response generator \mathcal{G} and a text-to-image translator \mathcal{F} . In the rest of this section, we will elaborate these two modules in detail.

4.2.1 Textual Dialogue Response Generator

The textual dialogue response generator \mathcal{G} is a sequence-to-sequence model based on the Trans-

former architecture (Vaswani et al., 2017b), it consists of a 24-layers Transformer with a hidden size of 1024 and 16 heads. Specifically, given a text dialogue context $U = \{u_1, \dots, u_l\}$ from \tilde{D}_S as source, and the target is a text $\tilde{R} = \{w_1, \dots, [\text{SEP}], [\text{DST}], \dots, [\text{SEP}], \dots, w_T\}$ with w_t the t -th word, the [DST] token means the following subsequence is a textual image description c . The generation loss is defined by

$$\mathcal{L}_G = \mathbb{E}_{(U, \tilde{R}) \sim \tilde{D}_S} [-\log p(\tilde{R})] \quad (2)$$

$$p(\tilde{R}) = \prod_t p(w_t | U, w_{1:t-1}) \quad (3)$$

Inference Given a new text dialogue context U , when a generated image description c occurs, it will be fed into the following text-to-image translator, then constructed to the codebook embeddings of its synonymous image.

4.2.2 Text-to-Image Translator

The text-to-image translator \mathcal{F} is also a sequence-to-sequence generation model based on the Transformer architecture, it consists of 24-layers Transformer with a hidden size of 1024 and 16 attention heads. Given an image $r^v \in \mathbb{R}^{H \times W \times 3}$ and its textual description $c = \{w_1, \dots, w_T\}$ from \tilde{D}_S , with the \mathcal{V}_E and \mathcal{Z} available, we can represent r^v in terms of the codebook indices of its encodings. More precisely, the quantized encoding of image r^v is given by $z_q = \mathbf{q}(\mathcal{V}_E(r^v)) \in \mathbb{R}^{h \times w \times d_z}$, and could be transferred to a sequence $s \in \{0, \dots, |\mathcal{Z}| - 1\}^{h \times w}$ of indices from the codebook \mathcal{Z} , which is obtained by replacing each code with its index in the codebook \mathcal{Z}

$$s_{i,j} = k \quad \text{such that} \quad (z_q)_{i,j} = z_k \quad (4)$$

Then we concatenate tokenized c and s to a single stream of tokens

$$x = \{w_1, \dots, w_T, [\text{SEP}], s_1, \dots, s_{h \times w}\} \quad (5)$$

and train an autoregressive transformer to model the joint distribution over the text and image tokens, the generation loss is defined by

$$\mathcal{L}_F = \mathbb{E}_{(c, r^v) \sim \tilde{D}_S} [-\log p(x)] \quad (6)$$

$$p(x) = \prod_t p(w_t | w_{1:t-1}) \prod_i p(s_i | c, s_{1:i-1}) \quad (7)$$

Inference Given a description c , we leverage the text-to-image translator to generate the representations $\hat{z} = \mathcal{F}(c) \in \mathbb{R}^{h \times w \times d_z}$ of its synonymous image.

4.2.3 Learning Details

Let us denote $\{\theta_g, \theta_\pi, \theta_\phi\}$ as the parameters of textual dialogue response generator \mathcal{G} , image tokenizer \mathcal{V} and text-to-image translator \mathcal{F} . In the pre-training stage, we use textual dialogues \mathcal{D}_C to estimate θ_g , use the ImageNet (Deng et al., 2009) to estimate θ_π , use <image description, image> pairs \mathcal{D}_P to estimate θ_ϕ . Then we fix θ_π , and jointly fine-tune θ_g and θ_ϕ with \tilde{D}_S , thus the final objective is to minimize the integrated loss

$$\mathcal{L} = \mathcal{L}_G + \lambda \mathcal{L}_F \quad (8)$$

where λ is a hyper parameter.

Remarks. In this work, we mainly focus on integrating text and image responses generation, but our proposed approach actually provides a recipe for a general solution to low-resource MDRG in which the target modality could be gifs, videos, or speech sounds, etc. To do that, one only needs to modify the text-to-image translator to make it compatible with the specific modality type, then pre-train a new text-to-<target modality> translator.

5 Experiments

5.1 Dataset

To evaluate the performance of Divter, we conduct comprehensive experiments on the PhotoChat dataset released by Zang et al. (2021), which is a multimodal conversational dataset consisting of 10917 images and 12286 dialogues, each of which is paired with a user image that is shared during the conversation, and each image is paired with its text description. The dataset has been split into 10286 train, 1000 dev, and 1000 test instances. More details are described in Appendix A.1.

5.2 Evaluation Metrics

We conduct evaluation with both automatic metrics and human judgements. For automatic evaluation, we focus on four aspects: (1) Image Intent Prediction, the goal of this task is to predict whether a image should be produced in the next turn for given context; (2) Text Description Generation; (3) Image Generation Quality ; (4) Text Response Generation. For (1), we follow Zang et al. (2021), which formulates the problem as a binary classification task, and use **F1** as metric; for (2) and (4), we use **PPL**, **BLEU** (Papineni et al., 2002), **Rouge** (Lin, 2004) and **F1**; for (3) we follow Ramesh et al. (2021) and use Frechet Inception Distance (**FID**) and Inception Score (**IS**).

Models	Intent	Image Description Generation				Image Generation		Text Response Generation				
		F1	PPL	B-1	B-2	Rouge	FID ↓	IS ↑	PPL	B-1	B-2	Rouge
BERT-base	53.2*	–	–	–	–	–	–	–	–	–	–	–
T5-3B	58.9*	–	–	–	–	–	–	–	–	–	–	–
S2S-TF	47.6	213.81	1.65	0.17	1.84	278.63	4.4 ± 0.8	329.43	3.61	0.40	3.05	
Divter	56.2	5.12	15.08	11.42	15.81	29.16	15.8 ± 0.6	59.63	6.52	1.66	5.69	
Divter (w/o \mathcal{G} pre-train)	47.3	122.56	1.99	0.23	2.60	29.78	15.5 ± 0.5	153.62	4.82	0.53	3.83	
Divter (w/o \mathcal{F}_ϕ pre-train)	55.9	5.23	15.01	11.20	15.63	262.09	4.9 ± 0.7	63.76	6.28	1.51	5.40	
Divter (w/o $\mathcal{G}, \mathcal{F}_\phi$ pre-train)	47.1	128.87	1.75	0.21	2.38	254.31	5.2 ± 0.6	163.85	4.53	0.48	3.55	
Divter (w/o joint learning)	55.6	5.20	15.00	11.36	15.73	29.04	15.4 ± 0.6	59.21	6.47	1.58	5.63	

Table 1: Automatic evaluation results of Divter and baselines on the test set. (w/o joint learning) means fine-tuning \mathcal{G} and \mathcal{F}_ϕ respectively rather than using Eq. 8. Numbers in bold mean that the improvement to the best baseline is statistically significant (t-test with p -value < 0.01). * reported by Zang et al. (2021).

Models	Context Coherence	Text Fluency	Image Quality	Background Consistency	Kappa
SCAN	–	–	1.95	0.96	0.65
S2S-TF	0.42	0.58	0.25	0.20	0.67
Divter	1.59	1.95	1.83	1.61	0.63

Table 2: Human evaluation results.

Models	Overall Improvement			Kappa
	W(%)	L(%)	T(%)	
Divter (pure text) vs. DialoGPT	34.4	35.7	29.9	0.64
Divter vs. DialoGPT	53.5	27.4	19.1	0.68

Table 3: Human evaluation results. (W, L, T) means (Win, Lose, Tie).

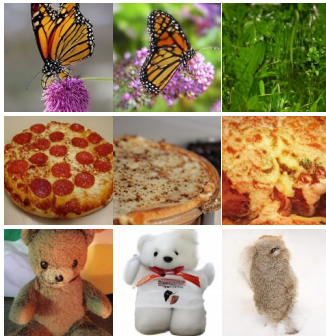


Figure 4: Qualitative assessment of various variants for image generation with same context as input in PhotoChat test set. 1st column: Divter. 2nd column: Divter w/o \mathcal{G} pre-train. 3rd column: Divter w/o \mathcal{F} pre-train.

For human evaluation, we randomly sample 200 dialogue contexts and generate responses from PhotoChat for Divter and baselines. Three human annotators are asked to score the response quality on a scale of {0, 1, 2} from four aspects: (1) **Context Coherence**: Whether the text response is coherent with the context; (2) **Text Fluency**: Whether the text response is natural and fluent; (3) **Image Quality**: The quality (including definition and integrity) of the image response; (4) **Background Consistency of Image**: For each dialogue, We select the top-8 generated/retrieved images group and ask the annotators to decide whether the group is consistent with the dialogue background, a qualitative

assessment is also shown in Figure 5. We report the average scores over three annotators, and the higher score means the better.

We also compare both pure text Divter and multimodal Divter with DialoGPT, respectively. The “pure text Divter” means we block the [DST] token in the vocabulary in the decoding stage, so that the responses would only contain texts. We also randomly sample 200 dialogues. To each annotator, two responses from different models are presented, which are randomly shuffled to hide their sources. The annotators then judge which response is more effective in improving the dialogue experience and attractiveness. The agreement among the annotators is measured by Fleiss’ Kappa (Fleiss, 1971).

5.3 Implementation Details

For the textual dialogue response generator \mathcal{G} , we use DialoGPT (Zhang et al., 2020) as pre-trained model initialization, trained on 147M conversation-like exchanges extracted from Reddit comment chains over a period spanning from 2005 through 2017. In the fine-tuning stage, we concatenate the context turns with the token [SEP] as a single sequence, we adopt Adam optimizer as an initial learning rate of 1e-5, and the batch size is 256, the training of PhotoChat is conducted on 16 Nvidia Tesla V100 32G GPU cards. We use beam search(size=5) to decode the text sequence.

For the image tokenizer \mathcal{V} , we inherit the model released by Ramesh et al. (2021).

For the text-to-image translator \mathcal{F} , we randomly select 5M <category image description, image> pairs from ImageNet, and <image description, image> pairs from YFCC100M (Thomee et al., 2016) as training data. We set the maximum image description length as 32, then pre-train \mathcal{F} for 3.5 million steps with a batch size of 256 accumulated on 16 Nvidia Tesla V100 32G GPUs. In the fine-tuning stage, we train PhotoChat for 50000 steps.

		Generated / Retrieved Images															
1	Generation (Generated Description: objects in the photo: animal, dog, carnivore, grassland)																
	Retrieval (Should contain “dog”.)																
2	Generation (Generated Description: objects in the photo: coffee cup, drink, bottle, mug, tea.)																
	Retrieval (Should contain “coffee cup”.)																
3	Generation (Generated Description: objects in the photo: curtain.)																
	Retrieval (Should contain “curtain”.)																

Figure 5: Examples of the images generated by Divter and the images retrieved by SCAN. The dialogue contexts are presented in Appendix A.2.

In the inference stage, we use CLIP (Radford et al., 2021) to rerank the generated 256 samples.

In the joint learning, we first train \mathcal{F} for 48000 steps, then jointly train \mathcal{G} and \mathcal{F} for 2000 steps. The λ in Eq.8 is 0.2. Early stopping on validation is adopted as a regularization strategy. All the hyper parameters are determined by grid search. More details are described in Appendix A.3.

We implement the image Auto-Encoder using the code <https://github.com/CompVis/taming-transformers>, implement the Textual Dialogue Response Generator using the code <https://github.com/microsoft/DialoGPT>, and implement the Text-to-Image Translator using the code <https://github.com/lucidrains/DALLE-pytorch>.

5.4 Baselines

Two pre-trained models **BERT-base** (Devlin et al., 2019) and **T5-3B** (Raffel et al., 2020) are selected as baselines to measure the “Image Intent Prediction” task in Section 5.2. They takes the text dialogue context as input, and predict “whether a image will be shared in the next turn”.

SCAN is proposed by Lee et al. (2018), the model captures interplay between image regions and text tokens to infer image-text similarity, SCAN achieves state-of-the-art performance of the “Image Retrieval” task on PhotoChat.

S2S-TF is a single sequence-to-sequence model with 24-layers Transformer, we only use PhotoChat

to train this multimodal generation model.

5.5 Evaluation Results

As shown in Table 1, our Divter achieves not only comparable performance with the state-of-the-art retrieval-based image response intent prediction model but also achieves remarkable performance in all the generation metrics. This indicates that Divter can accurately judge the timing of generating image response with the given dialogue context, and produce text responses that are coherent to the context, and generate high-quality image responses. The significant performance gap between Divter and the baseline models (e.g. S2S-TF, Divter variants) without pre-training indicates the superiority of our proposed learning strategy. Table 2 reports the results of human evaluation, our Divter also significantly outperforms the baselines on most of the aspects. The comparison results shown in Table 3 indicates (1): our Divter can achieve comparable performance on pure text response generation with DialoGPT; (2): the multimodal responses generated by Divter achieve a significant improvement on the dialogue experience and attractiveness in contrast to pure text dialogue model (DialoGPT).

5.6 Ablation Study

We conduct extensive ablation experiments over different variants to better understand their relative importance to the MDRG task. As shown in Table 1, all the variants lead to worse performance





Example 1	Example 2
<p>A: OMG...the new ice cream shop is amazing.</p> <p>A: I had the twist chocolate and vanilla but it was so fresh tasting. like you just made it. like you just made it.</p> <p>B: I call it the malado gilato.</p> <p>A: Sam wouldn't let me have another lick bc he thought I'd eat it all.</p> <p>D: That sounds interesting.</p> <p>D: Yes, could you please share it with me?</p> <p>D: Objects in the photo: Chocolate Ice cream, Dairy, Drink.</p> <p>D:</p> <div style="display: flex; justify-content: space-around;">   </div> <p>D: Wow! The ice cream looks so delicious.</p> <p>D: Sure, it tastes pretty good.</p>	<p>A: Have you been out in nature lately?</p> <p>B: Yes.</p> <p>A: I'm sitting at home now looking through some old photographs.</p> <p>B: I see. than.</p> <p>A: Would you like to see one of my favorites It's a cool shot of a honey bee near a beautiful flower.</p> <p>D: Objects in the photo: Honey bee, Insect, Animal, Flower.</p> <p>D:</p> <div style="display: flex; justify-content: space-around;">   </div> <p>D: It is a nice picture. Thank you for sharing.</p> <p>D: Haha, just enjoy the beautiful scenery.</p> <p>D: Yeah, definitely.</p>

Table 4: Examples of PhotoChat test set. In each example, the turns with the prefix of “A”/“B” are the given context; the **blue** text is the text description generated by Divter; the **left** image and the **red** response are generated by Divter, the **right** image is the ground-truth image.

in most of the metrics. For a more intuitive comparison, the qualitative assessment results are also shown in Figure 4. In particular, both quantitative and qualitative results on the ablation study validate that: (1) pre-training is crucial to low-resource multimodal dialogue response generation, since removing any component from pre-training causes performance drop when training data is small; (2) in terms of impact to performance of image generation, $\mathcal{F} > \mathcal{G}$, in terms of impact to performance of text generation, $\mathcal{G} > \mathcal{F}$; (3) The joint learning also has contributions to Divter, indicating that leveraging the integrated learning of textual context and visual image benefits more in contrast to any single one of them.

5.7 Case Study

To further investigate the quality of multimodal responses generated by Divter, we show two examples on the PhotoChat test data in Table 4. The given context of the first one is about “ice cream”, and the second one is about “honey bee”. As we can see, Divter can not only generate a realistic high-resolution image which is coherent to the background, but also generate the informative text responses grounded on the image. Separately, The high-quality generated images are comparable to those real-world ground truths, which demonstrates the practicability of Divter.

5.8 Discussions

Benefits over retrieval-based methods To further investigate and compare the generalization capability between Divter and the retrieval-based method, we also get top-10 generated images from Divter and equivalent retrieved images from SCAN model given the same context. As shown in Figure 5, on the one hand, the diversity and richness of the generated images are desirable, on the other hand, those retrieved results often suffer from wrong consistency with dialogue background. For example in the second case, the dialogue is talking about “coffee”, but the retrieved images contain some uncorrelated objects like “milk”, “cake”, “dog” and “snack”. And in the third example, all the retrieval results are mistaken since there is little “curtain” in the training and retrieval space. This demonstrates the fact that the performance of retrieval-based method is extremely limited in specific domains by the size of the pre-constructed conversational history repository, especially in the low-resource scenario. Furthermore, our proposed generation based method shows better generalization capability to tackle the low-resource challenge.

6 Conclusion

In this paper, we explore multimodal dialogue response generation under a low-resource setting. To overcome the challenges from the new task and insufficient training data, we propose Divter, a neural

conversational agent which incorporates text-to-image generation into text-only dialogue response generation, in which most parameters do not rely on the training data any more and can be estimated from large scale textual open domain dialogues and <image description, image> pairs. Extensive experiments demonstrate Divter achieves state-of-the-art results in automatic and human evaluation. In the future, we will explore more efficient methods to inject more modalities into response generation.

Acknowledgement

We thank anonymous reviewers for their insightful suggestions to improve this paper.

References

- Daniel Adiwardana, Minh-Thang Luong, David R. So, Jamie Hall, Noah Fiedel, Romal Thoppilan, Zi Yang, Apoorv Kulshreshtha, Gaurav Nemade, Yifeng Lu, and Quoc V. Le. 2020. [Towards a human-like open-domain chatbot](#).
- Shubham Agarwal, Trung Bui, Joon-Young Lee, Ioannis Konstas, and Verena Rieser. 2020. [History for visual dialog: Do we really need it?](#) In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8182–8197, Online. Association for Computational Linguistics.
- Peter Anderson, Basura Fernando, Mark Johnson, and Stephen Gould. 2016. [Spice: Semantic propositional image caption evaluation](#). In *ECCV*.
- Peter Anderson, Xiaodong He, Chris Buehler, Damien Teney, Mark Johnson, Stephen Gould, and Lei Zhang. 2018. [Bottom-up and top-down attention for image captioning and visual question answering](#). In *CVPR*.
- Feilong Chen, Fandong Meng, Xiuyi Chen, Peng Li, and Jie Zhou. 2021. [Multimodal incremental transformer with visual grounding for visual dialogue generation](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 436–446, Online. Association for Computational Linguistics.
- Jaemin Cho, Jiasen Lu, Dustin Schwenk, Hannaneh Hajishirzi, and Aniruddha Kembhavi. 2020. [X-lxmert: Paint, caption and answer questions with multi-modal transformers](#). In *EMNLP*.
- Marcella Cornia, Matteo Stefanini, Lorenzo Baraldi, and Rita Cucchiara. 2020. [Meshed-Memory Transformer for Image Captioning](#). In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*.
- Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José M. F. Moura, Devi Parikh, and Dhruv Batra. 2017. [Visual dialog](#). In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, page 326–335.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. [Imagenet: A large-scale hierarchical image database](#). In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Emily Dinan, Stephen Roller, Kurt Shuster, Angela Fan, Michael Auli, and Jason Weston. 2019. [Wizard of wikipedia: Knowledge-powered conversational agents](#). In *International Conference on Learning Representations*.
- Ming Ding, Zhuoyi Yang, Wenyi Hong, Wendi Zheng, Chang Zhou, Da Yin, Junyang Lin, Xu Zou, Zhou Shao, Hongxia Yang, and Jie Tang. 2021. [Cogview: Mastering text-to-image generation via transformers](#).
- Patrick Esser, Robin Rombach, and Björn Ommer. 2020. [Taming transformers for high-resolution image synthesis](#).
- Joseph L. Fleiss. 1971. [Measuring nominal scale agreement among many raters](#). *Psychological Bulletin*, 76:378–382.
- Philip Gage. 1994. [A new algorithm for data compression](#). *The C Users Journal archive*, 12:23–38.
- Mehdi Ghanimifard and Simon Dobnik. 2019. [What goes into a word: generating image descriptions with top-down spatial knowledge](#). In *Proceedings of the 12th International Conference on Natural Language Generation*, pages 540–551, Tokyo, Japan. Association for Computational Linguistics.
- Marjan Ghazvininejad, Chris Brockett, Ming-Wei Chang, Bill Dolan, Jianfeng Gao, Wen-tau Yih, and Michel Galley. 2018. [A knowledge-grounded neural conversation model](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 32.
- Karol Gregor, Ivo Danihelka, Alex Graves, Danilo Rezende, and Daan Wierstra. 2015. [Draw: A recurrent neural network for image generation](#). In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 1462–1471, Lille, France. PMLR.
- Qingbao Huang, Jielong Wei, Yi Cai, Changmeng Zheng, Junying Chen, Ho-fung Leung, and Qing Li. 2020. [Aligned dual channel graph convolutional network for visual question answering](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7166–7176, Online. Association for Computational Linguistics.

- Byeongchang Kim, Jaewoo Ahn, and Gunhee Kim. 2020. [Sequential latent knowledge selection for knowledge-grounded dialogue](#). In *International Conference on Learning Representations*.
- Kuang-Huei Lee, Xi Chen, Gang Hua, Houdong Hu, and Xiaodong He. 2018. Stacked cross attention for image-text matching. *arXiv preprint arXiv:1803.08024*.
- Jiwei Li, Michel Galley, Chris Brockett, Jianfeng Gao, and Bill Dolan. 2016a. [A diversity-promoting objective function for neural conversation models](#). In *Proceedings of the 2016 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 110–119, San Diego, California. Association for Computational Linguistics.
- Jiwei Li, Michel Galley, Chris Brockett, Georgios Spithourakis, Jianfeng Gao, and Bill Dolan. 2016b. [A persona-based neural conversation model](#). In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 994–1003, Berlin, Germany. Association for Computational Linguistics.
- Linjie Li, Zhe Gan, Yu Cheng, and Jingjing Liu. 2019a. Relation-aware graph attention network for visual question answering. *ICCV*.
- Linxiao Li, Can Xu, Wei Wu, Yufan Zhao, Xueliang Zhao, and Chongyang Tao. 2020. Zero-resource knowledge-grounded dialogue generation. *arXiv preprint arXiv:2008.12918*.
- Wenbo Li, Pengchuan Zhang, Lei Zhang, Qiuyuan Huang, Xiaodong He, Siwei Lyu, and Jianfeng Gao. 2019b. Object-driven text-to-image synthesis via adversarial training.
- Zujie Liang, Huang Hu, Can Xu, Chongyang Tao, Xubo Geng, Yining Chen, Fan Liang, and Daxin Jiang. 2021. [Maria: A visual experience powered conversational agent](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 5596–5611, Online. Association for Computational Linguistics.
- Chin-Yew Lin. 2004. [ROUGE: A package for automatic evaluation of summaries](#). In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.
- Jiasen Lu, Jianwei Yang, Dhruv Batra, and Devi Parikh. 2016. Hierarchical question-image co-attention for visual question answering. In *NIPS*, page 289–297.
- Elman Mansimov, Emilio Parisotto, Jimmy Ba, and Ruslan Salakhutdinov. 2016. Generating images from captions with attention. In *ICLR*.
- Anh Nguyen, Jeff Clune, Yoshua Bengio, Alexey Dosovitskiy, and Jason Yosinski. 2017. Plug & play generative networks: Conditional iterative generation of images in latent space. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. IEEE.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. [Bleu: a method for automatic evaluation of machine translation](#). In *Proceedings of the 40th Annual Meeting of the Association for Computational Linguistics*, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Jiaxin Qi, Yulei Niu, Jianqiang Huang, and Hanwang Zhang. 2020. Two causal principles for improving visual dialog. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*.
- Weizhen Qi, Yeyun Gong, Yu Yan, Can Xu, Bolun Yao, Bartuer Zhou, Biao Cheng, Daxin Jiang, Jiusheng Chen, Ruofei Zhang, et al. 2021. Prophetnet-x: Large-scale pre-training models for english, chinese, multi-lingual, dialog, and code generation. *arXiv preprint arXiv:2104.08006*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. 2021. [Learning transferable visual models from natural language supervision](#). In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8748–8763. PMLR.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*, 21(140):1–67.
- Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. 2021. [Zero-shot text-to-image generation](#). In *Proceedings of the 38th International Conference on Machine Learning*, volume 139 of *Proceedings of Machine Learning Research*, pages 8821–8831. PMLR.
- Scott Reed, Zeynep Akata, Xinchen Yan, Lajanugen Logeswaran, Bernt Schiele, and Honglak Lee. 2016. [Generative adversarial text to image synthesis](#). In *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1060–1069, New York, New York, USA. PMLR.
- Mengye Ren, Ryan Kiros, and Richard Zemel. 2015. Exploring models and data for image question answering. In *NIPS*, page 2953–2961.

- Stephen Roller, Emily Dinan, Naman Goyal, Da Ju, Mary Williamson, Yinhan Liu, Jing Xu, Myle Ott, Eric Michael Smith, Y-Lan Boureau, and Jason Weston. 2021. [Recipes for building an open-domain chatbot](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 300–325, Online. Association for Computational Linguistics.
- Abigail See, Stephen Roller, Douwe Kiela, and Jason Weston. 2019. [What makes a good conversation? how controllable attributes affect human judgments](#).
- Iulian Serban, Alessandro Sordoni, Yoshua Bengio, Aaron Courville, and Joelle Pineau. 2016. [Building end-to-end dialogue systems using generative hierarchical neural network models](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 30.
- Lifeng Shang, Zhengdong Lu, and Hang Li. 2015. [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 (Volume 1: Long Papers)*, pages 1577–1586, Beijing, China. Association for Computational Linguistics.
- Ilya Sutskever, Oriol Vinyals, and Quoc V. Le. 2014. Sequence to sequence learning with neural networks. In *NIPS*, page 3104–3112.
- Chongyang Tao, Shen Gao, Mingyue Shang, Wei Wu, Dongyan Zhao, and Rui Yan. 2018. [Get the point of my utterance! learning towards effective responses with multi-head attention mechanism](#). In *IJCAI-18*, pages 4418–4424. International Joint Conferences on Artificial Intelligence Organization.
- Bart Thomee, David A. Shamma, Gerald Friedland, Benjamin Elizalde, Karl Ni, Douglas Poland, Damian Borth, and Li-Jia Li. 2016. [Yfcc100m](#). *Communications of the ACM*, 59(2):64–73.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017a. [Attention is all you need](#). In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017b. [Attention is all you need](#). In *Advances in neural information processing systems*, pages 5998–6008.
- Oriol Vinyals and Quoc Le. 2015. [A neural conversational model](#).
- Chen Xing, Wei Wu, Yu Wu, Ming Zhou, Yalou Huang, and Wei-Ying Ma. 2017. [Hierarchical recurrent attention network for response generation](#).
- Can Xu, Wei Wu, Chongyang Tao, Huang Hu, Matt Schuerman, and Ying Wang. 2019. [Neural response generation with meta-words](#). *arXiv preprint arXiv:1906.06050*.
- Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhutdinov, Richard Zemel, and Yoshua Bengio. 2015. [Show, attend and tell: Neural image caption generation with visual attention](#). In *Proceedings of the 32nd International Conference on Machine Learning*, volume 37 of *Proceedings of Machine Learning Research*, pages 2048–2057, Lille, France. PMLR.
- Tao Xu, Pengchuan Zhang, Qiuyuan Huang, Han Zhang, Zhe Gan, Xiaolei Huang, and Xiaodong He. 2018. [Attngan: Fine-grained text to image generation with attentional generative adversarial networks](#).
- Ze Yang, Wei Wu, Huang Hu, Can Xu, Wei Wang, and Zhoujun Li. 2021. [Open domain dialogue generation with latent images](#). *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(16):14239–14247.
- Xiaoxue Zang, Lijuan Liu, Maria Wang, Yang Song, Hao Zhang, and Jindong Chen. 2021. [PhotoChat: A human-human dialogue dataset with photo sharing behavior for joint image-text modeling](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 6142–6152, Online. Association for Computational Linguistics.
- Hainan Zhang, Yanyan Lan, Liang Pang, Jiafeng Guo, and Xueqi Cheng. 2019. [ReCoSa: Detecting the relevant contexts with self-attention for multi-turn dialogue generation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 3721–3730, Florence, Italy. Association for Computational Linguistics.
- Han Zhang, Tao Xu, Hongsheng Li, Shaoting Zhang, Xiaogang Wang, Xiaolei Huang, and Dimitris Metaxas. 2017. [Stackgan: Text to photo-realistic image synthesis with stacked generative adversarial networks](#). In *ICCV*.
- Saizheng Zhang, Emily Dinan, Jack Urbanek, Arthur Szlam, Douwe Kiela, and Jason Weston. 2018. [Personalizing dialogue agents: I have a dog, do you have pets too?](#) In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2204–2213, Melbourne, Australia. Association for Computational Linguistics.
- Yizhe Zhang, Siqi Sun, Michel Galley, Yen-Chun Chen, Chris Brockett, Xiang Gao, Jianfeng Gao, Jingjing Liu, and Bill Dolan. 2020. [Dialogpt: Large-scale generative pre-training for conversational response generation](#). In *ACL, system demonstration*.
- Tiancheng Zhao, Ran Zhao, and Maxine Eskenazi. 2017. [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 (Volume 1: Long Papers)*, pages 654–664, Vancouver, Canada. Association for Computational Linguistics.

Yufan Zhao, Can Xu, Wei Wu, and Lei Yu. 2020. Learning a simple and effective model for multi-turn response generation with auxiliary tasks. *arXiv preprint arXiv:2004.01972*.

Hao Zhou, Minlie Huang, Tianyang Zhang, Xiaoyan Zhu, and Bing Liu. 2018. Emotional chatting machine: Emotional conversation generation with internal and external memory.

A Appendix

A.1 Dataset

Table 5 reports the statistics of the PhotoChat dataset.

Split	images	dialogues	turns	tokens
Train	8,917	10,286	130,546	827,154
Dev	1,000	1,000	12,701	80,214
Test	1,000	1,000	12,852	80,847
Total	10,917	12,286	156,099	988,215

Table 5: PhotoChat statistics.

A.2 Dialogue Contexts in Figure 5

Table 6 presents the textual dialogue contexts of the examples shown in the Figure 5.

	Textual Dialogue Context
1	A: hows your day going? A: beautiful sky today A: Have you been near a mountain lately B: yes A: beautiful right, just took a hike today with my dog. B: my college placed in mountain area B: super enjoy the lot A: Oh great, do you have an aquarium at your college? B: how is your dog A: He is great. I'll share a pic. B: i want to see your dog
2	B: hi B: hello friend A: hi A: how are you A: i am doing well B: how are yow B: great A: i am having some coffee B: ok A: you should come over for a cup! A: do you like coffee?
3	A: what are you doing? A: great moment B: Just finishing up with some work so I can start fresh tomorrow!! B: Great moment? What's that mean? A: you chat dude B: ??? A: Curtain you have like B: I don't understand. A: why dude?

Table 6: Dialogue contexts of the examples shown in the Figure 5.

A.3 More Implementation Details

The CLIP model assigns a score based on how well the image matches the description, we use CLIP to rerank the generated 256 samples, and select the best image as the final response. To obtain high-quality training set, we discard the instances with the prefix of “The photo has your * #” in descriptions, “*” includes “mom”, “dad”, “daughter”, “sister”, “uncle”, etc. “#” is name of a person. To build the training set for text-to-image translator \mathcal{F} from ImageNet, we combine the text “Objects in the photo:” and textual categorical name of each image to build the <category image description, image> pair. To train the baseline S2S-TF model, we also use the image tokenizer \mathcal{V} to tokenize each image, and combine the image tokens with text tokens to form a single stream as the generation source/target.