

# Towards Building an Open-Domain Dialogue System Incorporated with Internet Memes

Hua Lu\*, Zhen Guo\*, Chanjuan Li†, Yunyi Yang†, Huang He, Siqu Bao

Baidu Inc., China  
{luhua05, guozhenguo}@baidu.com

## Abstract

In recent years, Internet memes have been widely used in online chatting. Compared with text-based communication, conversations become more expressive and attractive when Internet memes are incorporated. This paper presents our solutions for the Meme incorporated Open-domain Dialogue (MOD) Challenge of DSTC10, where three tasks are involved: text response modeling, meme retrieval, and meme emotion classification. Firstly, we leverage a large-scale pre-trained dialogue model for coherent and informative response generation. Secondly, based on interaction-based text-matching, our approach can retrieve appropriate memes with good generalization ability. Thirdly, we propose to model the emotion flow (EF) in conversations and introduce an auxiliary task of emotion description prediction (EDP) to boost the performance of meme emotion classification. Experimental results on the MOD dataset demonstrate that our methods can incorporate Internet memes into dialogue systems effectively.

## Introduction

As Internet memes can make dialogues more vivid and engaging, nowadays, people tend to incorporate memes when chatting online (Kulkarni 2017; Jiang and Vásquez 2020). Despite that Internet memes have become an effective means of expression, they are rarely considered by most open-domain dialogue systems. In DSTC10, the Meme incorporated Open-domain Dialogue (MOD) challenge aims to incorporate Internet memes into open-domain dialogues. It includes the following three tasks: (1) **Text Response Modeling**: given a multi-modal context, the task here is to generate a coherent and informative text response. (2) **Meme Retrieval**: given a multi-modal context and a text response, the task aims to retrieve an appropriate meme. (3) **Meme Emotion Classification**: given a multi-modal context and a text response with a meme, the task here is to predict the emotion type of the Internet meme. Figure 1 shows two examples of conversations in the MOD dataset involving texts, Internet memes, and emotions.

\* First two authors contributed equally to this work.

† Work was done during internship at Baidu.

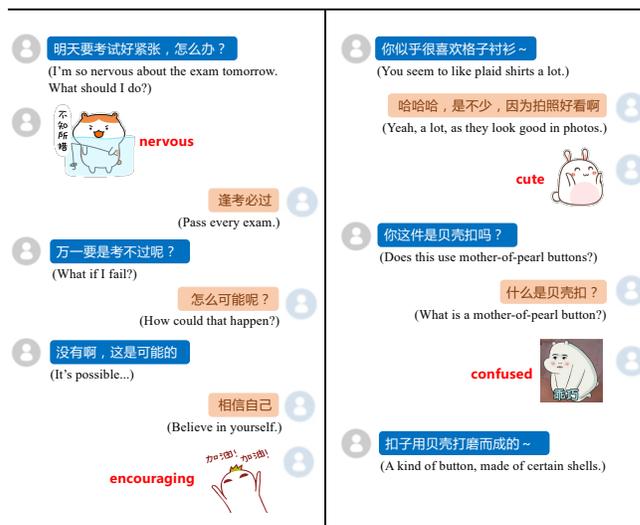


Figure 1: Two examples from the MOD dataset. Corresponding emotion is annotated for each meme in red.

In particular, the test set of MOD is divided into an easy test version and a hard test version. The latter, which contains memes not appearing in the train set, is used to evaluate the generalization ability of the dialogue system. In this work, we introduce the following solutions for the three tasks:

- In Task1, we leverage a powerful pre-trained open-domain dialogue model for coherent and informative text response generation.
- In Task2, we represent memes with textual information consisting of meme titles and OCR texts (extracted from the memes). Based on interaction-based text-matching, our approach can retrieve appropriate memes with good generalization ability.
- In Task3, we propose to model the emotion flow (EF) in conversations and introduce an auxiliary task of emotion description prediction (EDP) to enhance the ability of meme emotion recognition.

Experimental results demonstrate that our methods can effectively incorporate Internet memes into dialogue systems.

Our methods achieve first place in four out of six leaderboards and second place in the others with competitive performance.

## Methodology

Our detailed solutions towards these three tasks will be discussed in the following.

### Text Response Modeling

In open-domain conversation, users are free to talk about any topic, and the system’s replies are expected to meet a high standard in many aspects, including coherence, consistency, informativeness, etc. Incorporated with Internet memes, the dialogue context can be formulated as  $C_t = \{\langle u_1, m_1 \rangle, \langle u_2, m_2 \rangle, \dots, \langle u_t, m_t \rangle\}$ , where  $u_i$  represents the  $i$ -th utterance text and  $m_i$  refers to its associated meme. If no Internet meme is used in the  $i$ -th utterance,  $m_i$  will be denoted as None. The task of text response modeling is to generate the response  $r$  (i.e., next utterance  $u_{t+1}$ ) given the multi-modal dialogue context  $C_t$ .

As suggested in the MOD baseline (Fei et al. 2021), the Internet meme can be represented with visual features extracted by EfficientNet (Tan and Le 2019). While in our preliminary experiments, we observe that incorporating memes, whether as visual features or as textual features, brings little benefit to text response generation. The reasons might be two-fold. Firstly, as memes are usually about emotional expressions with little narrative information, the absence of memes might not remarkably undermine a dialogue system’s text response generation ability. Secondly, given that the MOD dataset is collected by inserting memes into existing conversations, the reliance on these memes might be relatively weak for text response generation. Therefore, we treat this task as a standard text-based response generation problem, with memes in the context set as None.

In this paper, we utilize the pre-trained open-domain dialogue model PLATO-2 (Bao et al. 2020) for text response generation. As illustrated in Figure 2(a), the input to the network is the concatenation of context and response. The input representation is calculated as the sum of the token, segment, and position embeddings. The network employs flexible attention mechanisms, where bi-directional attention is enabled for better contextual understanding and uni-directional attention is utilized for auto-regressive response generation. The training objective of text response generation is to minimize the following negative log-likelihood (NLL) loss:

$$\mathcal{L}_{NLL} = -\log p_{\text{generation}}(u_{t+1}|C_t) \quad (1)$$

### Meme Retrieval

Meme retrieval is a crucial component in meme incorporated dialogue systems. This task is to select an appropriate meme from the Internet meme set, given the multi-modal dialogue context  $C_t$  and text response  $u_{t+1}$ . Formally, we denote Internet meme set as  $M = \{m_1, m_2, \dots, m_n\}$ , where  $m_i$  is the representation of the  $i$ -th meme. In this work, we represent memes with textual information, consisting of meme titles and OCR texts (extracted from the memes). Although there exists plenty of vision or multi-modal pre-training

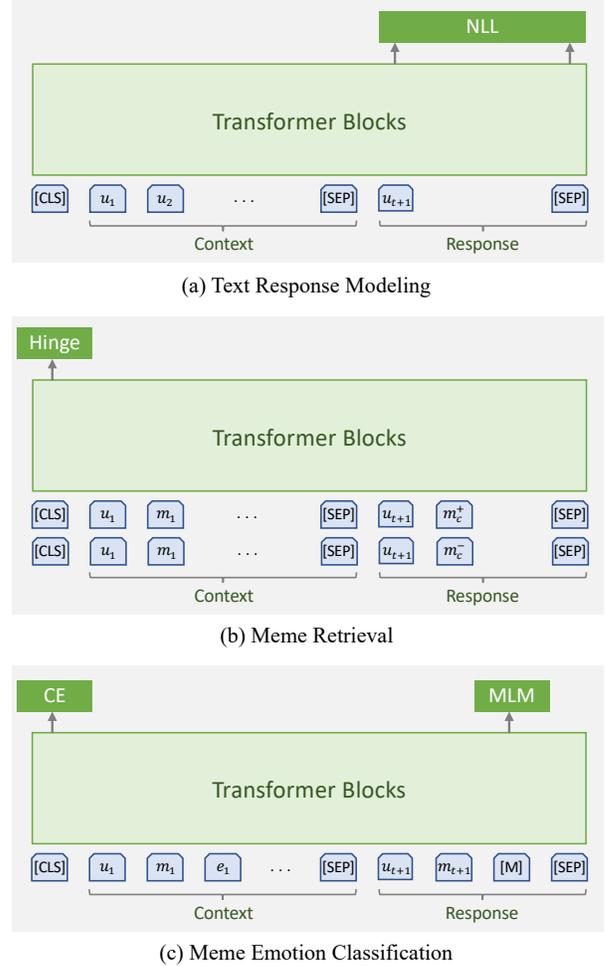


Figure 2: Illustration of network inputs and training objectives for three tasks.

works (Chen et al. 2019; Li et al. 2020b; Gan et al. 2020; Radford et al. 2021), they are less effective at meme feature extraction due to the gap of data distributions between real photos and devised memes. Experimental results also suggest that the meme title and OCR text can sufficiently represent the meaning of Internet memes.

Therefore, we treat the meme retrieval task as a text-matching problem and employ the cross-encoder architecture for relevance estimation. The network overview for meme retrieval is shown in Figure 2(b). The input includes the dialogue context  $C_t$ , the textual response  $u_{t+1}$ , and one candidate meme  $m_c$ . During training, a pair of positive and negative samples are fed into the network. The output of [CLS] token is passed through a fully-connected layer, and a following sigmoid function to obtain the relevance probability  $p_{\text{matching}}(l_{m_c} = 1|C_t, u_{t+1}, m_c)$ , where  $l_{m_c}$  stands for the label to choose meme  $m_c$  or not given the dialogue context and corresponding textual response. The training objective

is to minimize the following margin ranking loss:

$$\mathcal{L}_{\text{Hinge}} = \max(0, p_{\text{matching}}(l_{m_c^-} = 1 | C_t, u_{t+1}, m_c^-) - p_{\text{matching}}(l_{m_c^+} = 1 | C_t, u_{t+1}, m_c^+) + \alpha) \quad (2)$$

where  $\alpha$  is a pre-defined margin parameter,  $m_c^+$  is a positive meme, and  $m_c^-$  is a negative meme. During training, we enable dynamic random negative sampling, which means that as model training progresses, different sets of negative samples are dynamically sampled in each epoch.

During inference, the Internet meme for the given dialogue context and response is selected as:

$$m^* = \operatorname{argmax}_{m_c \in M} p_{\text{matching}}(l_{m_c} = 1 | C_t, u_{t+1}, m_c) \quad (3)$$

### Meme Emotion Classification

Rather than classify the sentiment of an Internet meme, this task aims to predict the meme emotion type situated in the dialogue context. Considering that the emotions of two interlocutors seldom encounter abrupt changes (or to some extent the changes might be traceable), we propose to model the emotion flow (EF) in multi-turn conversations. The textual descriptions of emotions are integrated into the dialogue context. Specifically, the utterance at turn  $i$  is composed of the utterance text, meme text, and textual emotion description, i.e.,  $\langle u_i, m_i, e_i \rangle$ . The meme emotion recognition can be considered as a classical sequence classification task, and the training objective is to minimize the standard cross-entropy loss.

Additionally, we introduce an auxiliary task of emotion description prediction (EDP) to boost meme emotion recognition performance. As shown in Figure 2(c), the auxiliary task is to recover the masked tokens (i.e., textual emotion description in the response) by minimizing the masked language model (MLM) loss (Devlin et al. 2018). In this way, the training objective of meme emotion classification is to minimize the following integrated loss:

$$\mathcal{L} = \mathcal{L}_{CE} + \mathcal{L}_{MLM} \quad (4)$$

where  $\mathcal{L}_{CE}$  is the cross-entropy loss of the classification task and  $\mathcal{L}_{MLM}$  denotes the MLM loss of the emotion description prediction task.

## Experiments

In the DSTC10 MOD challenge, one open-domain dialogue dataset incorporated with Internet memes is constructed. The memes used in the dataset have been annotated with titles. The MOD test set is divided into easy test version and hard test version. The latter one containing unseen memes is used to evaluate the ability of dialogue systems to exploit new Internet memes. Detailed statistics of the dataset are summarized in Table 1.

### Settings

The evaluation of the MOD challenge covers the following three tasks:

|           | # dialogue | # utterance | # Internet meme | # meme emotion |
|-----------|------------|-------------|-----------------|----------------|
| Train     | 41,644     | 558,181     | 274             | 50             |
| Valid     | 1,000      | 13,666      | 274             | 47             |
| Easy Test | 1,000      | 13,999      | 274             | 50             |
| Hard Test | 1,530      | 20,258      | 307             | 50             |

Table 1: Statistics of the MOD dataset.

- **Task1: Text Response Modeling.** Given a dialogue context, the model needs to produce a coherent and informative text response. The automatic evaluation metrics of this task include BLEU-2/4 (Papineni et al. 2002), DIST-1/2 (Li et al. 2015).
- **Task2: Meme Retrieval.** Given a multi-modal dialogue context and a text response, the model needs to retrieve an appropriate Internet meme. The evaluation metrics of this task include Recall\_10@1, Recall\_10@3, Recall\_10@5, and MAP.
- **Task3: Meme Emotion Classification.** The model needs to predict the corresponding emotion type of the used Internet meme. The evaluation metrics of this task include Accuracy@1, Accuracy@3, and Accuracy@5.

**Implementation Details** In the experiments, we utilize dialogue pre-training models of PLATO-2 (Bao et al. 2020) to improve the performance of all three tasks. The models have 32 transformer blocks and 32 attention heads, with up to 1.6 billion parameters. The generation model of PLATO-2 is used for Task1. The evaluation model of PLATO-2 is employed for the fine-tuning of Task2 and Task3.

In Task1, responses are generated using beam search, with a beam size of 5. The maximum sequence length for the context and response is set to 256 and 128, respectively. In Task2, we set the margin parameter  $\alpha$  to 0.2 and the ratio of positive training samples to negative ones to 1:5. During the fine-tuning of Task2 and Task3, we use Adam (Kingma and Ba 2014) optimizer with a learning rate of  $2e-5$  and warmup steps of 4000. All the models are fine-tuned for five epochs with a batch size of 64. The implementation is based on the PaddlePaddle framework, and the experiments are carried out on 4 Nvidia Tesla V100 GPUs (32G RAM).

## Experimental Results

The experimental results on these three tasks are discussed in the following.

**Text Response Modeling** The evaluation results of text response modeling on two test versions are summarized in Table 2, with the best score written in bold. The final ranking for this task is based on human evaluation, where five metrics are considered: grammatical correctness, informativeness, naturalness, relevance to the dialogue history, and overall feeling based on the above four metrics. The score of each metric ranges from 1 to 5. The higher, the better. The final human score is the average of the above five metric scores. We rank second on the easy version and first on the hard version. From Table 2, it can be observed that our automatic evaluation results are relatively poor, especially on

| Rank                | BLEU-2      | BLEU-4      | DIST-1      | DIST-2       | Human Score |
|---------------------|-------------|-------------|-------------|--------------|-------------|
| <i>easy version</i> |             |             |             |              |             |
| 1                   | <b>5.08</b> | <b>4.25</b> | 1.90        | <b>26.7</b>  | <b>3.74</b> |
| 2(ours)             | 3.57        | 1.32        | 1.93        | 21.5         | 3.69        |
| 3                   | 3.78        | 1.89        | <b>2.20</b> | 20.2         | 3.60        |
| <i>hard version</i> |             |             |             |              |             |
| 1(ours)             | 3.65        | 1.30        | 1.17        | 17.68        | <b>3.72</b> |
| 2                   | <b>5.04</b> | <b>3.65</b> | 1.10        | <b>20.00</b> | 3.68        |
| 3                   | 4.03        | 1.65        | <b>1.36</b> | 16.60        | 3.65        |

Table 2: Task1 evaluation results on two test versions, with the best score written in bold.

| Rank                | Recall_10@1 | Recall_10@3 | Recall_10@5 | MAP         |
|---------------------|-------------|-------------|-------------|-------------|
| <i>easy version</i> |             |             |             |             |
| 1(ours)             | <b>56.8</b> | <b>84.7</b> | <b>94.4</b> | <b>72</b>   |
| 2                   | 34.4        | 60.4        | 76.5        | 52.3        |
| 3                   | 34.2        | 59.6        | 76.0        | 52.3        |
| <i>hard version</i> |             |             |             |             |
| 1(ours)             | <b>42.0</b> | <b>69.7</b> | <b>80.9</b> | <b>58.8</b> |
| 2                   | 27.9        | 50.8        | 66.7        | 45.1        |
| 3                   | 27.5        | 51.0        | 67.6        | 50.3        |

Table 3: Task2 evaluation results on two test versions, with the best score written in bold.

the metrics of BLEU-2/4, while the human evaluation results are relatively competitive. This phenomenon further verifies that the correlation between automatic evaluation metrics and human evaluation is weak in open-domain conversations (Liu et al. 2016).

| Rank                | Accuracy@1  | Accuracy@3  | Accuracy@5  |
|---------------------|-------------|-------------|-------------|
| <i>easy version</i> |             |             |             |
| 1(ours)             | <b>62.3</b> | <b>83.4</b> | <b>89.5</b> |
| 2                   | 58.3        | 74.3        | 78.9        |
| 3                   | 57.0        | 72.0        | 77.6        |
| <i>hard version</i> |             |             |             |
| 1                   | <b>29.7</b> | <b>40.6</b> | <b>49.9</b> |
| 2(ours)             | 27.3        | 39.2        | 47.5        |
| 3                   | 27.0        | 37.6        | 47.2        |

Table 4: Task3 evaluation results on two test versions, with the best score written in bold.

**Meme Retrieval** The evaluation results of meme retrieval on two test versions are summarized in Table 3, with the best score written in bold. The final ranking for this task is based on the Recall.10@1 score, which is the fraction of the ground-truth Internet meme ranked first among ten meme

candidates. Our proposed method obtains first place on both versions and outperforms other teams by a large margin. Possible reasons behind such improvements are discussed as follows:

- The modality of dialogue context and memes is unified by representing Internet memes with texts, making it easier for the model to estimate the relevance. Furthermore, with the pre-trained language models, our text-matching strategy has the generalization ability to retrieve the unseen memes that are not available during training.
- Different from the dual-encoder architecture (Mazaré et al. 2018; Zang et al. 2021), which performs self-attention over the input and the candidate separately, we employ the cross-encoder architecture to yield rich interactions between the dialogue context and the meme candidate. With this interaction-based text-matching, our model can retrieve appropriate memes more effectively.

**Meme Emotion Classification** The evaluation results of meme emotion classification on two test versions are summarized in Table 4, with the best score written in bold. The final ranking for this task is based on the Accuracy@1 score, which is the fraction of the ground-truth emotion type that obtains the highest score among all emotion types. We rank first on the easy version and second on the hard version. In particular, our meme emotion classification model, combining the EF and EDP strategies, achieves 4.0% absolute improvement on Accuracy@1 on the easy test set over the second-ranked team (62.3% vs. 58.3%). On the hard test set, the performance of all the models degenerates significantly, revealing that the emotion classification ability for unseen memes is still weak.

## Case Analysis

To further analyze the performance of our proposed methods, several examples are provided in Figure 3. As shown in the left example, our model is able to generate a coherent and informative response. The interlocutor on the left-hand side seems to be a vegetarian, and our model generates a response consistent with the persona. In the middle and right examples, our model is able to retrieve a relevant meme and accurately identify the emotion type contained in the meme. These dialogue examples suggest that our system can generate a natural and informative response incorporated with an appropriate meme.

## Ablation Study

In this section, several ablation studies are carried out on the validation set to better understand the contribution of each component.

**Meme Retrieval** To evaluate the generalization ability of our matching model, we select 20 memes to formulate the unseen validation set and remove corresponding samples from the train set. The experimental results on the unseen validation set by the matching models trained on the original train set and the filtered train set are summarized in Table 5. The results indicate that our model has the generalization ability to exploit unseen memes. However, the gap between

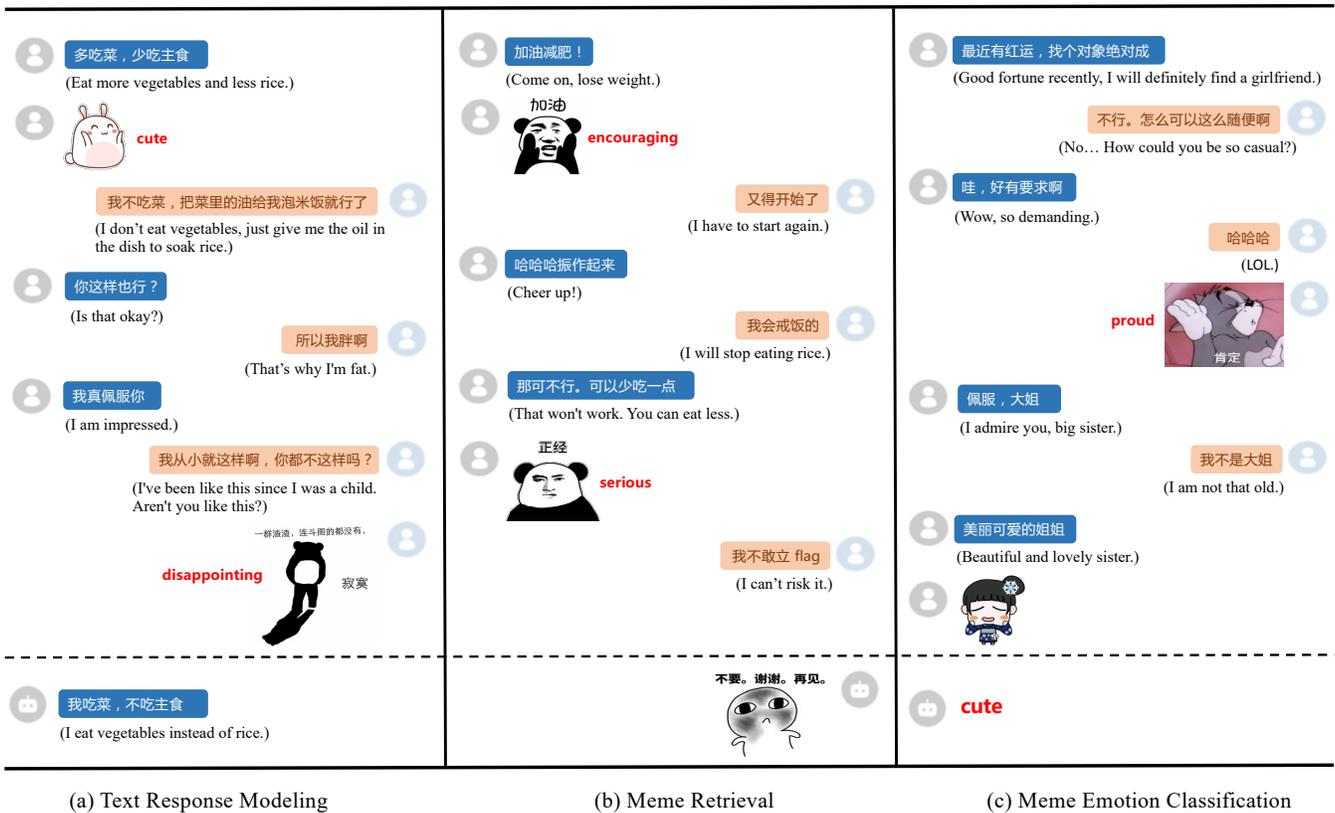


Figure 3: Examples of the input (upper) and our output (bottom) for each task.

|          | Recall_10@1 | Recall_10@3 | Recall_10@5 | MAP         |
|----------|-------------|-------------|-------------|-------------|
| original | <b>50.4</b> | <b>85.0</b> | <b>93.6</b> | <b>68.9</b> |
| filtered | 43.6        | 73.3        | 88.0        | 62.1        |

Table 5: Task2 ablation study on the unseen validation set.

original and filtered (50.4% vs. 43.6% on Recall\_10@1) suggests that there is still some room for improvement on the model’s generalization ability.

|            | Accuracy@1  | Accuracy@3  | Accuracy@5  |
|------------|-------------|-------------|-------------|
| base model | 62.7        | 83.8        | 90.3        |
| + EF       | 64.9        | 84.1        | 90.8        |
| + EF + EDP | <b>65.3</b> | <b>84.4</b> | <b>90.8</b> |

Table 6: Task3 ablation study on the validation set. EF refers to the modeling of emotion flow. EDP refers to the task of emotion description prediction.

**Meme Emotion Classification** To boost the performance of meme emotion classification, we propose to model the emotion flow (EF) in multi-turn conversations and introduce an auxiliary task of the emotion description prediction

(EDP). To analyze the effects of these two components, we carry out the ablation studies, and evaluation results are summarized in Table 6. Compared to the base model, the incorporation of EF modeling gives rise to a significant improvement (+2.2% on Accuracy@1). The combination of EF and EDP obtains +2.6% absolute improvement on Accuracy@1 (65.3% vs. 62.7%), verifying the effectiveness of our proposed strategies.

## Related Work

In this section, we will discuss related works on multi-modal conversation and emotion recognition.

There are several works that attempt to incorporate multi-modal information into conversations. Das et al. introduces the task of VisDial, where the AI agent needs to hold a meaningful conversation with humans and answer questions about the contents of the input image. In addition to the conversational question answering, there are some other tasks where natural and engaging conversations are conducted based on a shared image, such as image-grounded conversations (Mostafazadeh et al. 2017), and image-chat (Shuster et al. 2018). Recently, the PhotoChat dataset (Zang et al. 2021) is presented, which focuses on the photo-sharing behavior in online messaging and aims to improve the photo-sharing experience in conversations. Unlike the above works concentrated on photos, the MOD dataset incorporates Internet memes into open-domain conversations to enhance com-

munication expressiveness.

In open-domain conversation, it is crucial to recognize the emotional state accurately and generate a response appropriately (Rashkin et al. 2018). To boost emotion detection in conversations, HiTrans (Li et al. 2020a) utilizes BERT (Devlin et al. 2018) as the low-level transformer to generate utterance representations and employs another high-level transformer to obtain context representations. In TUCORE-GCN (Lee and Choi 2021), the task of emotion recognition is treated as a dialogue-based relation extraction, where a dialogue graph is constructed and a graph convolution network is employed for relation classification. In addition to the text-based conversations, emotion has been widely analyzed in some areas of computer vision, such as facial expression recognition (Minaee, Minaei, and Abdolrashidi 2021), image emotion classification (Yang, She, and Sun 2017), and so on. In this work, rather than classify the sentiment of an Internet meme, this task aims to predict the meme emotion type situated in the dialogue context.

## Conclusion

In this paper, we introduce our solutions for the DSTC10 MOD challenge. Firstly, we leverage a large-scale pre-trained dialogue model for coherent and informative response generation. Secondly, based on interaction-based text-matching, our approach can retrieve appropriate memes with good generalization ability. Thirdly, we propose to model the emotion flow (EF) in conversations and introduce an auxiliary task of emotion description prediction (EDP) to boost the performance of meme emotion recognition. Comprehensive experiments have been conducted on the MOD dataset. Experimental results demonstrate that our methods can effectively incorporate Internet memes into dialogue systems and accurately recognize the meme emotion. Our methods with competitive performance achieve first place in four out of six leaderboards and second place in the others.

## Acknowledgments

We would like to thank the anonymous reviewers for their constructive suggestions; Xinxian Huang for the helpful discussions.

## References

- Bao, S.; He, H.; Wang, F.; Wu, H.; Wang, H.; Wu, W.; Guo, Z.; Liu, Z.; and Xu, X. 2020. Plato-2: Towards building an open-domain chatbot via curriculum learning. *arXiv preprint arXiv:2006.16779*.
- Chen, Y.-C.; Li, L.; Yu, L.; El Kholy, A.; Ahmed, F.; Gan, Z.; Cheng, Y.; and Liu, J. 2019. Uniter: Learning universal image-text representations.
- Das, A.; Kottur, S.; Gupta, K.; Singh, A.; Yadav, D.; Moura, J. M.; Parikh, D.; and Batra, D. 2017. Visual dialog. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 326–335.
- Devlin, J.; Chang, M.-W.; Lee, K.; and Toutanova, K. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Fei, Z.; Li, Z.; Zhang, J.; Feng, Y.; and Zhou, J. 2021. Towards Expressive Communication with Internet Memes: A New Multimodal Conversation Dataset and Benchmark. *arXiv preprint arXiv:2109.01839*.
- Gan, Z.; Chen, Y.-C.; Li, L.; Zhu, C.; Cheng, Y.; and Liu, J. 2020. Large-scale adversarial training for vision-and-language representation learning. *arXiv preprint arXiv:2006.06195*.
- Jiang, Y.; and Vásquez, C. 2020. Exploring local meaning-making resources: A case study of a popular Chinese internet meme (biaoqingbao). *Internet Pragmatics*, 3(2): 260–282.
- Kingma, D. P.; and Ba, J. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kulkarni, A. 2017. Internet meme and Political Discourse: A study on the impact of internet meme as a tool in communicating political satire. *Journal of Content, Community & Communication Amity School of Communication*, 6.
- Lee, B.; and Choi, Y. S. 2021. Graph Based Network with Contextualized Representations of Turns in Dialogue. *arXiv preprint arXiv:2109.04008*.
- Li, J.; Galley, M.; Brockett, C.; Gao, J.; and Dolan, B. 2015. A diversity-promoting objective function for neural conversation models. *arXiv preprint arXiv:1510.03055*.
- Li, J.; Ji, D.; Li, F.; Zhang, M.; and Liu, Y. 2020a. Hitrans: A transformer-based context-and speaker-sensitive model for emotion detection in conversations. In *Proceedings of the 28th International Conference on Computational Linguistics*, 4190–4200.
- Li, W.; Gao, C.; Niu, G.; Xiao, X.; Liu, H.; Liu, J.; Wu, H.; and Wang, H. 2020b. Unimo: Towards unified-modal understanding and generation via cross-modal contrastive learning. *arXiv preprint arXiv:2012.15409*.
- Liu, C.-W.; Lowe, R.; Serban, I. V.; Noseworthy, M.; Charlin, L.; and Pineau, J. 2016. How not to evaluate your dialogue system: An empirical study of unsupervised evaluation metrics for dialogue response generation. *arXiv preprint arXiv:1603.08023*.
- Mazaré, P.-E.; Humeau, S.; Raison, M.; and Bordes, A. 2018. Training millions of personalized dialogue agents. *arXiv preprint arXiv:1809.01984*.
- Minaee, S.; Minaei, M.; and Abdolrashidi, A. 2021. Deep-emotion: Facial expression recognition using attentional convolutional network. *Sensors*, 21(9): 3046.
- Mostafazadeh, N.; Brockett, C.; Dolan, B.; Galley, M.; Gao, J.; Spithourakis, G. P.; and Vanderwende, L. 2017. Image-grounded conversations: Multimodal context for natural question and response generation. *arXiv preprint arXiv:1701.08251*.
- Papineni, K.; Roukos, S.; Ward, T.; and Zhu, W.-J. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting of the Association for Computational Linguistics*, 311–318.
- Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*.

Rashkin, H.; Smith, E. M.; Li, M.; and Boureau, Y.-L. 2018. Towards empathetic open-domain conversation models: A new benchmark and dataset. *arXiv preprint arXiv:1811.00207*.

Shuster, K.; Humeau, S.; Bordes, A.; and Weston, J. 2018. Image Chat: Engaging Grounded Conversations. *arXiv preprint arXiv:1811.00945*.

Tan, M.; and Le, Q. 2019. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, 6105–6114. PMLR.

Yang, J.; She, D.; and Sun, M. 2017. Joint Image Emotion Classification and Distribution Learning via Deep Convolutional Neural Network. In *IJCAI*, 3266–3272.

Zang, X.; Liu, L.; Wang, M.; Song, Y.; Zhang, H.; and Chen, J. 2021. PhotoChat: A Human-Human Dialogue Dataset with Photo Sharing Behavior for Joint Image-Text Modeling. *arXiv preprint arXiv:2108.01453*.