

Multimodal Dialog Systems with Dual Knowledge-enhanced Generative Pretrained Language Model

Xiaolin Chen, Xuemeng Song, *Senior Member, IEEE*, Liqiang Jing, Shuo Li, Linmei Hu, and Liqiang Nie, *Senior Member, IEEE*

Abstract—Text response generation for multimodal task-oriented dialog systems, which aims to generate the proper text response given the multimodal context, is an essential yet challenging task. Although existing efforts have achieved compelling success, they still suffer from two pivotal limitations: 1) *overlook the benefit of generative pre-training*, and 2) *ignore the textual context related knowledge*. To address these limitations, we propose a novel dual knowledge-enhanced generative pretrained language model for multimodal task-oriented dialog systems (DKMD), consisting of three key components: *dual knowledge selection*, *dual knowledge-enhanced context learning*, and *knowledge-enhanced response generation*. To be specific, the dual knowledge selection component aims to select the related knowledge according to both textual and visual modalities of the given context. Thereafter, the dual knowledge-enhanced context learning component targets seamlessly integrating the selected knowledge into the multimodal context learning from both global and local perspectives, where the cross-modal semantic relation is also explored. Moreover, the knowledge-enhanced response generation component comprises a revised BART decoder, where an additional dot-product knowledge-decoder attention sub-layer is introduced for explicitly utilizing the knowledge to advance the text response generation. Extensive experiments on a public dataset verify the superiority of the proposed DKMD over state-of-the-art competitors.

Index Terms—Multimodal Task-oriented Dialog Systems; Text Response Generation; Generative Pretrained Language Model; Dual Knowledge Selection



1 INTRODUCTION

ACCORDING to the report of Salesforce¹, roughly 68% of customers prefer dialog agents rather than waiting for human services because dialog agents can provide quick answers. Due to its substantial economic value, task-oriented dialog systems, which aim to conduct specific tasks in certain vertical domains, such as ticket booking and restaurant table reserving, have attracted increasing research attention. Although existing research efforts have attained impressive results, most of them work purely on the single-modality (*i.e.*, textual modality) dialog system, neglecting that both the user and the agent may need to employ certain visual clues (*i.e.*, images) to deliver their needs or services. As depicted in Figure 1, the agent shows special dishes for the user via images in the utterance u_4 , while the user describes his/her desired shopping mall with the image in the utterance u_7 . Therefore, multimodal task-oriented dialog systems merit our specific attention.

In general, multimodal task-oriented dialog systems mainly involve two tasks [1]: the text response generation and the image response selection. As compared with the image response selection task, the text response generation task is more challenging, whose performance is far from satisfactory. Existing multimodal task-oriented dialog systems mainly adopt the encoder-decoder framework for text response generation. In particular, recent studies have recognized the pivotal role of the knowledge base for multimodal dialog systems, and designed various schemes for incorporating knowledge to enhance the user’s intention understanding [2], [3], [4], [5], [6], [7], [8], [9]. Although they have achieved significant progress, these research efforts suffer from two key limitations. 1) **Overlook the benefit of generative pre-training**. Previous studies follow the conventional train-from-scratch paradigm and fail to leverage the generative pre-training technique, ignoring the powerful text generation ability of generative pretrained language models (GPLMs) [10], [11], [12]. 2) **Ignore the textual context related knowledge**. Previous studies only refer to the knowledge base according to the images provided by the user (*e.g.*, the picture associated with the utterance u_7 in Figure 1). Namely, they only involve the visual context related knowledge to enhance the user intention modeling. Nevertheless, they overlook that the textual context plays the dominant role in the dialog, and could also be used for fetching related knowledge from the knowledge base to enhance the text response generation.

To address these limitations, in this work, we target at comprehensively utilizing the multimodal context in

X. Chen is with the School of Software, Shandong University, Jinan 250101, China (e-mail: cxlicd@gmail.com).

X. Song, L. Jing and S. Li are with the School of Computer Science and Technology, Shandong University, Qingdao 266237, China (e-mail: sxmstc@gmail.com, jingliqiang6@gmail.com, zile020401@gmail.com).

L. Hu is with the Beijing Key Lab of Intelligent Telecommunications Software and Multimedia, School of Computer Science, Beijing University of Posts and Telecommunications, Beijing 100876, China (e-mail: hulinmei@bupt.edu.cn).

L. Nie is with the School of Computer Science and Technology, Harbin Institute of Technology (Shenzhen campus), Shenzhen 518055, China (e-mail: nieliqiang@gmail.com).

1. <https://startupbonsai.com/chatbot-statistics>.

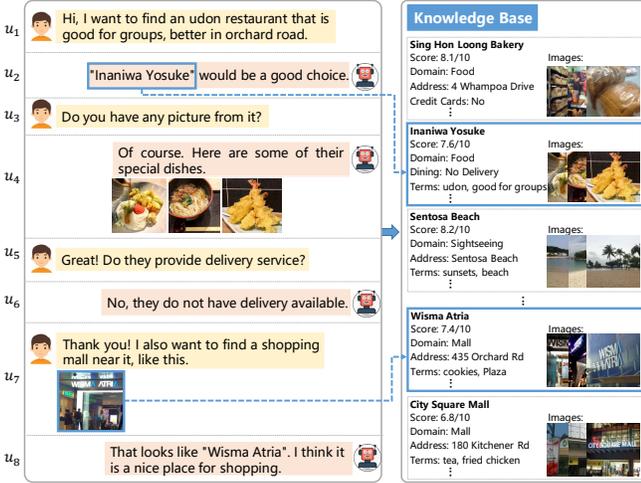


Fig. 1. Illustration of a multimodal dialog system between a user and an agent. “u”: utterance.

knowledge selection with the backbone of GPLMs to improve the performance of text response generation for multimodal task-oriented dialog systems. This is, however, non-trivial due to the following three challenges. 1) The multimodal dialog context cannot fit well with the GPLMs that are pretrained with only the textual corpus, and thus directly feeding the multimodal context into GPLMs deteriorates the text generation capability of GPLMs. Therefore, how to subtly adapt GPLMs to cope with the multimodal dialog context constitutes the main challenge. 2) As aforementioned, the context related knowledge is of crucial importance to the text response generation. For example, as shown in Figure 1, the agent can generate the proper response (e.g., u_6) only conditioned on the attribute knowledge (i.e., “No Delivery”) of “Inaniwa Yosuke”. Hence, how to accurately select the knowledge concerning the given multimodal context and properly inject knowledge to enhance the user intention modeling and text response generation with GPLMs is another crucial challenge. 3) Both textual context and visual context serve to demonstrate the user’s intention, where they are closely related and mutually reinforce each other. As shown in Figure 1, the user demonstrates his/her intention of finding a restaurant and a shopping mall with not only the textual description (e.g., ‘an udon restaurant’, ‘good for groups’, ‘better in orchard road’ and ‘shopping mall near it’), but also images of his/her desired shopping mall. Therefore, how to mine the context cross-modality semantic relation based on GPLMs and thus accurately capture the user’s intention is a tough challenge.

To address the challenges mentioned above, we propose a novel dual knowledge-enhanced generative pretrained language model for multimodal task-oriented dialog systems, DKMD for short, where BART [12] is adopted as the backbone. As illustrated in Figure 2, DKMD contains three vital components: *dual knowledge selection*, *dual knowledge-enhanced context learning*, and *knowledge-enhanced response generation*. To be specific, the dual knowledge selection component is devised to select the context related knowledge from the whole knowledge base according to both the textual and visual

modality of the given context. Thereafter, the dual knowledge-enhanced context learning component aims to properly incorporate dual knowledge (i.e., both textual and visual context related knowledge) to the multimodal context modeling and hence accurately captures the user’s intention. In particular, considering different roles of multimodal context in conveying the user’s intention, we design the knowledge-enhanced context representation module with the global knowledge-enhanced textual representation learning and local knowledge-enhanced visual representation learning. Moreover, we introduce the dual cross-modal representation refinement module, comprising vision-oriented representation refinement and text-oriented representation refinement, to capture the semantic relation hidden in the multimodal context and facilitate the user intention modeling. Ultimately, the knowledge-enhanced response generation component targets at explicitly using the knowledge to advance the text response generation, where a revised BART decoder with an additional dot-product knowledge-decoder attention (DKDA) sub-layer is introduced. Extensive experiments on one public dataset have fully validated the effectiveness of our proposed DKMD.

Our main contributions can be summarized as follows:

- To the best of our knowledge, we are among the first to incorporate the GPLMs into multimodal task-oriented dialog systems. In particular, we propose a novel dual knowledge-enhanced generative pretrained language model for the text response generation task.
- We propose the dual knowledge-enhanced context learning component, which seamlessly integrates the selected dual knowledge into the multimodal context learning from global and local perspectives and also explores the context cross-modality semantic relation to facilitate the user intention modeling.
- We devise the knowledge-enhanced decoder that can utilize knowledge to stimulate the precise text response generation explicitly. As a byproduct, we have released codes and involved parameters to facilitate the research community².

2 RELATED WORK

In this section, we briefly introduce the studies of task-oriented dialog systems and pretrained language models, respectively.

2.1 Task-oriented Dialog Systems

Traditional task-oriented dialog systems mainly adopt a pipeline structure and usually contain the following functional modules: natural language understanding, dialogue state tracking, policy learning, and natural language generation. Specifically, the natural language understanding module aims to classify the user’s intentions, and then the dialogue state tracking module can track the current state and fill in the predefined slots. Thereafter, the policy learning module predicts the following action on the basis

2. <https://multimodaldialog.wixsite.com/website>.

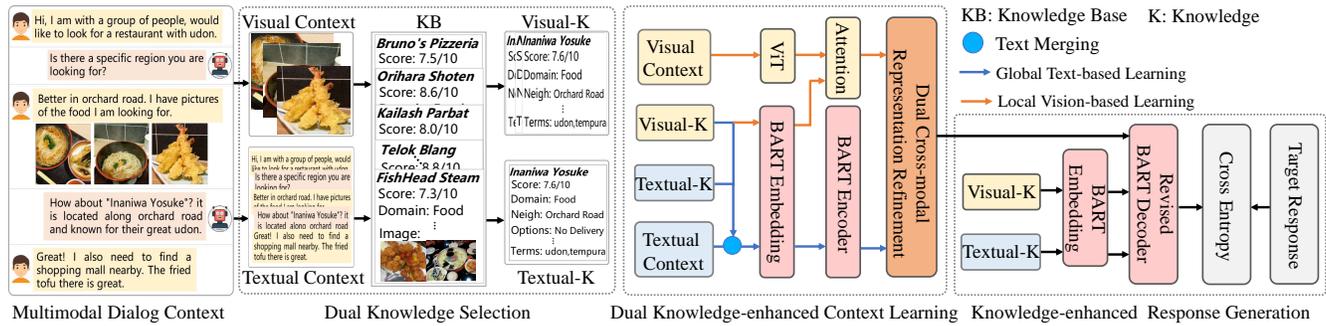


Fig. 2. Illustration of the proposed model. DKMD consists of three vital components: *Dual Knowledge Selection*, *Dual Knowledge-enhanced Context Learning*, and *Knowledge-enhanced Response Generation*.

of the current state representation, and the natural language generation module returns the response through generation methods [13], [14], [15] or predefined templates. Despite the remarkable success of the pipeline-based methods, they are prone to suffer from error propagation [16] and heavy dependence on the sequential modules [17].

With the evolution of deep neural networks, several efforts have been made toward end-to-end task-oriented dialog systems [18], [19], [20]. Although these efforts have achieved compelling success, they focus on the pure textual modality, *i.e.*, the single-modality task-oriented dialog system. In reality, both the user and the agent may need to refer to certain images to deliver their needs or services. Therefore, Saha et al. [1] investigated the multimodal dialog system, and proposed a multimodal hierarchical encoder-decoder model (MHRED) for addressing the two primary tasks of the multimodal dialog system: text response generation and image response selection. Moreover, they released a large-scale multimodal dialog dataset in the context of online fashion shopping, named MMD, which significantly promotes the research progress on multimodal dialog systems. In particular, several efforts further explore the semantic relation in the multimodal dialog context and incorporate knowledge based on the framework of MHRED [2], [3], [4], [5], [6], [7]. For example, Liao et al. [5] developed a taxonomy-based visual semantic learning module to capture the fine-grained semantics (*e.g.*, the category and attributes of a product) in product images, and introduced a memory network to integrate the knowledge of fashion style tips. In addition, Nie et al. [7] devised a multimodal dialog system with multiple decoders, which can generate diverse responses according to the user's intention and adaptively integrate the related knowledge. Recently, some studies have resorted to Transformer [21] to investigate the multimodal dialog systems [8], [9] due to its impressive results in natural language processing (NLP) tasks [10], [11], [12], [22], [23]. For example, He et al. [8] introduced a Transformer-based element-level encoder, which can capture the semantic dependencies of multimodal elements (*i.e.*, words and images) via the attention mechanism.

As compared with the image response selection task, the text response generation task is more challenging, whose performance is far from satisfactory. Therefore, in this work, we particularly study the task of text response generation in the context of multimodal task-oriented

dialog systems. Notably, although the pioneer studies have achieved tremendous strides on this task, they overlook the benefit of pre-training and only utilize the attribute knowledge concerning the visual context of the dialog. Beyond that, in this work, we aim to generate a precise response by utilizing pretrained techniques and capturing related knowledge from both the textual context and visual context perspectives.

2.2 Pretrained Language Models

As an emerging technique, pretrained language models have been arresting much research attention [10], [11], [12], [24], [25] and achieve remarkable success in plenty of NLP tasks. Initially, Word2vec [26] and GloVe [27] are proposed to obtain pretrained word embeddings based on shallow architectures. Thereafter, with the flourish of Transformer, considerable studies make efforts to devise Transformer-based pretrained models [10], [11], [12]. For example, Devlin et al. [10] proposed the bidirectional encoder representation from transformer (BERT) to capture the accurate textual representation via two pre-training tasks: masked language model and next sentence prediction. In addition, Lewis et al. [12] presented a Transformer-based denoising autoencoder (BART) for the language generation task, with the bidirectional encoder and the autoregressive decoder. With the remarkable progress of generative pretrained language models, a surge of follow-up works solve diverse tasks by adapting publicly available pretrained language models. For example, Yu et al. [24] designed a vision-guided generative pretrained language model based on BART and text-to-text transfer transformer [11] for the multimodal abstractive summarization task.

Although generative pretrained language models have shown compelling success in many tasks, limited efforts have been devoted to conducting the text response generation in multimodal task-oriented dialog systems. To fill the research gap, we adapt the publicly available pretrained BART to integrate the multimodal context and corresponding knowledge to enhance the response generation capability of our model.

3 PRELIMINARY

We choose BART as our backbone for the text response generation due to its superior performance in many text

generation tasks, such as multimodal abstractive summarization [24] and community question answering [28]. In particular, BART is a Transformer-based denoising autoencoder, consisting of a position-wise embedding layer, a bidirectional encoder, and an autoregressive decoder.

Position-wise Embedding Layer. Suppose we have a text $t = [x_1, x_2, \dots, x_M]$, where x_q represents the q -th token and M is the total number of tokens in the text. Each token x_q is assigned with an initial embedding \mathbf{e}_q by a linear transformation as follows,

$$\mathbf{e}_q = \mathbf{W}_e^\top \mathbf{g}_q, q = 1, 2, \dots, M, \quad (1)$$

where $\mathbf{W}_e \in \mathbb{R}^{|\mathcal{U}| \times D}$ is the token embedding matrix to be fine-tuned, $|\mathcal{U}|$ is the number of tokens in the token vocabulary \mathcal{U} , and D is the dimension of the token embeddings. $\mathbf{g}_q \in \mathbb{R}^{|\mathcal{U}|}$ is the one-hot vector, indicating the index of x_q in the token vocabulary.

To encode the order information among input tokens, position encodings [29] are further inserted as follows,

$$\mathbf{Z}_0^{enc} = [\mathbf{e}_1; \mathbf{e}_2; \dots; \mathbf{e}_M]^\top + \mathbf{E}_{pos}, \quad (2)$$

where $\mathbf{E}_{pos} \in \mathbb{R}^{M \times D}$ is the positional embedding matrix, each row of which corresponds to a token in the given text. $\mathbf{Z}_0^{enc} \in \mathbb{R}^{M \times D}$ is the matrix containing all the final embeddings of tokens in the input text. $[\cdot]$ refers to the concatenation operation.

Bidirectional Encoder. The bidirectional encoder of BART, denoted as \mathcal{B}_e , is composed of L encoder layers, and used to encode the input text. To be specific, each layer has two sub-layers: 1) multi-head self-attention mechanism (MSA), which aims to model the semantic dependencies among tokens in the input text; and 2) feed-forward network (FFN), used for the nonlinear transformation. Notably, each sub-layer is followed by a residual connection and layer normalization (LN) operations to enhance the model generalization as follows,

$$\begin{cases} \mathbf{Z}_l^S = LN(MSA(\mathbf{Z}_{l-1}^{enc}) + \mathbf{Z}_{l-1}^{enc}), \\ \mathbf{Z}_l^{enc} = LN(FFN(\mathbf{Z}_l^S) + \mathbf{Z}_l^S), \end{cases} l = 1, 2, \dots, L, \quad (3)$$

where $\mathbf{Z}_l^{enc} \in \mathbb{R}^{M \times D}$ refers to the output of l -th encoder layer, and \mathbf{Z}_0^{enc} is obtained by the aforementioned position-wise embedding layer in Eqn. (2). $\mathbf{Z}_l^S \in \mathbb{R}^{M \times D}$ is the intermediate output of MSA in the l -th encoder layer. Ultimately, the output of the L -th layer is treated as the final encoded context representation, namely $\mathbf{Z}_L^{enc} \in \mathbb{R}^{M \times D}$.

Autoregressive Decoder. The decoder \mathcal{B}_d of BART also contains L decoder layers, which can generate the response based on the encoded representation. To be specific, each layer consists of three sub-layers: 1) masked multi-head self-attention mechanism (MMSA), combined the mask mechanism and the operation making the output embeddings offset by one position, which ensures that the current output only depends on the known outputs; 2) multi-head encoder-decoder attention mechanism (MEDA), which can distinguish the informative output of the encoder and adaptively assign weights to different previous outputs; and 3) FFN. Similar to the

encoder, each sub-layer is followed by a residual connection and layer normalization operations as follows,

$$\begin{cases} \mathbf{q}_l^S = LN(MMSA(\mathbf{q}_{l-1}^{dec}) + \mathbf{q}_{l-1}^{dec}), \\ \mathbf{q}_l^E = LN(MEDA(\mathbf{q}_l^S, \mathbf{Z}_L^{enc}) + \mathbf{q}_l^S), \\ \mathbf{q}_l^{dec} = LN(FFN(\mathbf{q}_l^E) + \mathbf{q}_l^E) \end{cases} l = 1, 2, \dots, L, \quad (4)$$

where $\mathbf{q}_l^S \in \mathbb{R}^D$ and $\mathbf{q}_l^E \in \mathbb{R}^D$ refer to the intermediate output of MMSA and MEDA in the l -th decoder layer, respectively. $\mathbf{q}_l^{dec} \in \mathbb{R}^D$ denotes the final output of l -th decoder layer. Thereafter, the decoder \mathcal{B}_d of BART employs the linear transformation and softmax function to project the decoder output into the probability space as follows,

$$\tilde{\mathbf{y}} = softmax(\mathbf{q}_L^{dec} \mathbf{W}_y + \mathbf{b}_y), \quad (5)$$

where \mathbf{W}_y and \mathbf{b}_y represent the weight matrix and bias vector, respectively. $\tilde{\mathbf{y}} \in \mathbb{R}^{|\mathcal{U}|}$ denotes the predicted token distribution. The predicted token of the current time step can be obtained according to the largest element of $\tilde{\mathbf{y}}$.

4 MODEL

In this section, we first formulate the research task of text response generation in multimodal task-oriented dialog systems, and then detail the proposed model illustrated in Figure 2, which comprises three vital components: *dual knowledge selection*, *dual knowledge-enhanced context learning*, and *knowledge-enhanced response generation*.

4.1 Problem Formulation

In this work, we aim to investigate the task of text response generation conditioned on multimodal task-oriented dialog systems. Suppose we have a set of N training dialog pairs $\mathcal{D} = \{(\mathcal{C}_1, \mathcal{R}_1), (\mathcal{C}_2, \mathcal{R}_2), \dots, (\mathcal{C}_N, \mathcal{R}_N)\}$, where each pair comprises a multimodal dialog context \mathcal{C}_i (*i.e.*, sequence of historical dialog utterances between the user and the agent) and a target text response \mathcal{R}_i . Notably, apart from the common textual modality, each utterance in \mathcal{C}_i can also involve certain related images, as the user/agent may sometimes use images to facilitate the request/response expression. Accordingly, each multimodal dialog context \mathcal{C}_i can be decomposed into a sequence of historical textual utterances $\mathcal{T}_i = [t_g^i]_{g=1}^{N_T^i}$ (*i.e.*, a sequence of tokens) and a set of images $\mathcal{V}_i = \{v_j^i\}_{j=1}^{N_V^i}$, where t_g^i is the g -th token and v_j^i is the j -th image of \mathcal{C}_i . N_T^i and N_V^i refer to the total number of tokens and images, respectively. Notably, N_V^i may be zero, *i.e.*, there is no image in the context \mathcal{C}_i . The target text response of \mathcal{C}_i can be denoted as $\mathcal{R}_i = [r_n^i]_{n=1}^{N_R^i}$, where r_n^i denotes the n -th token and N_R^i is the total number of tokens in the response.

Besides, the multimodal dialog system is equipped with a knowledge base, containing rich knowledge of N_K entities $\mathcal{K} = \{e_p\}_{p=1}^{N_K}$. Specifically, for each entity e_p , the knowledge base provides a set of attributes \mathcal{A}_p and images \mathcal{I}_p characterizing it. The attributes (*e.g.*, score, domain, and location) reveal the semantic information of the entity, while the images intuitively describe the entity, like the photos showing the appearance or food of a restaurant entity.

In a sense, we aim to devise a novel model \mathcal{F} which can accurately generate the appropriate text response given the multimodal context and the knowledge base as follows,

$$\mathcal{F}(C_i, \mathcal{K} | \Theta_F) \rightarrow \mathcal{R}_i, \quad (6)$$

where Θ_F represents the model parameters.

4.2 Dual Knowledge Selection

To effectively leverage the entity knowledge, the premise is to correctly select the related knowledge entities from the whole knowledge base for the given multimodal context. Considering the multimodal nature of the given context, we devise the dual knowledge selection with the *text-based knowledge selection* and *vision-based knowledge selection*. Specially, the text-based and vision-based knowledge selections aim to retrieve the related knowledge entities according to the textual and visual modality of the given context, respectively.

Text-based Knowledge Selection. To capture the related knowledge entities concerning the textual context, we directly judge which knowledge entity in the knowledge base is mentioned in the given textual context. Namely, for each entity e_p in the knowledge base, we check whether it appears in the given textual context. If it appears, we select its attributes \mathcal{A}_p as the related knowledge. Notably, we here only consider attributes rather than images due to the fact that the attribute knowledge is essential to understanding user’s intentions and generating the text response [2]. In this vein, we can obtain the overall knowledge set involved with the textual context, denoted as $\mathcal{K}_t^A = \mathcal{A}_1^t \cup \mathcal{A}_2^t \cup \dots \cup \mathcal{A}_{N_k^t}^t$, where \mathcal{A}_m^t is the attribute set of the m -th related knowledge entity and N_k^t is the number of knowledge entities appearing in the textual context.

Vision-based Knowledge Selection. As aforementioned, the goal of the vision-based knowledge selection is to find the related knowledge entities for the given dialog context with its visual modality. As for the same entity, there can be various images characterizing it, and thus we employ the visual feature similarity to select the related knowledge for the visual context. To be specific, we first extract the visual features of entities in the knowledge base \mathcal{K} and images in \mathcal{V} of the given context with ViT-B/32 [30] pretrained by CLIP [31], due to its superior performance in various computer vision tasks [32], [33]. Thereafter, for each image v_j in \mathcal{V} , we measure its similarity to each image of entities in \mathcal{K} based on the cosine similarity between their corresponding visual features, and select the top k most similar knowledge entities. Similar to the text-based knowledge selection, we also only consider the semantic knowledge (*i.e.*, attributes) of them. In this way, we can acquire the related knowledge set conditioned on the visual context as $\mathcal{K}_v^A = \mathcal{A}_1^v \cup \mathcal{A}_2^v \cup \dots \cup \mathcal{A}_{N_v}^v$, where \mathcal{A}_j^v refers to the related semantic knowledge of the image v_j (*i.e.*, attributes of the related knowledge entities of image v_j).

4.3 Dual Knowledge-enhanced Context Learning

To accurately capture the user’s intention hidden in the multimodal context, we design the dual knowledge-enhanced context learning scheme with

two modules: *knowledge-enhanced context representation* and *dual cross-modal representation refinement*, where the semantic relation between the textual context and the visual context is mined in the latter module. For simplicity, we temporarily omit the subscript i that indexes the training samples.

4.3.1 Knowledge-enhanced Context Representation

In the multimodal dialog, the textual context tends to convey the user’s intention from a global perspective, while the visual context would exert roles from the local perspective by reinforcing certain local intention via intuitive images. As shown in Figure 1, the textual context generally indicates the user’s intention of finding a restaurant and a shopping mall with detailed requirements (*e.g.*, domain and delivery), while the visual context (*i.e.*, the image in u_7) only exhibits the desired shopping mall. Therefore, we conduct the *global knowledge-enhanced textual representation learning* and *local knowledge-enhanced visual representation learning*.

Global Knowledge-enhanced Textual Representation Learning. Considering the global role of the textual context, we jointly utilize the related knowledge of both textual and visual context to promote the textual context learning. In particular, we merge the textual context \mathcal{T} and the related knowledge of both textual and visual context (*i.e.*, \mathcal{K}_t^A and \mathcal{K}_v^A) as a whole $\mathcal{X}_t = [\mathcal{T}, \mathcal{K}_t^A, \mathcal{K}_v^A] = [x_t^1, x_t^2, \dots, x_t^{N_t}]$. Here, x_t^q denotes the q -th token and N_t refers to the total number of tokens. In particular, we first obtain the initial embedding of \mathcal{X}_t , denoted as $\mathbf{E}_t \in \mathbb{R}^{N_t \times D}$, by the position-wise embedding layer of BART in Eqns. (1) and (2). Thereafter, to capture the semantic representation, we feed the initial embedding \mathbf{E}_t into the bidirectional encoder \mathcal{B}_e of BART defined in Eqn. (3) as follows,

$$\mathbf{T}_t = \mathcal{B}_e(\mathbf{E}_t), \quad (7)$$

where $\mathbf{T}_t \in \mathbb{R}^{N_t \times D}$ is the knowledge-enhanced representation of the textual context.

Local Knowledge-enhanced Visual Representation Learning. As aforementioned, each image of the multimodal context can be associated with certain knowledge entities. In light of this, we aim to utilize the corresponding knowledge (*i.e.*, attributes of the related knowledge entity) to enhance the visual context representation.

In particular, given the set of images $\mathcal{V} = \{v_1, v_2, \dots, v_{N_v}\}$, we first employ ViT-B/32 pretrained by CLIP to encode each image v_j and obtain the visual representation as follows,

$$\begin{cases} \mathbf{v}_j = \mathcal{B}_v(v_j), j = 1, 2, \dots, N_v, \\ \mathbf{E}_v = [\mathbf{v}_1; \mathbf{v}_2; \dots; \mathbf{v}_{N_v}]^T, \end{cases} \quad (8)$$

where $\mathbf{E}_v \in \mathbb{R}^{N_v \times D}$ refers to the initial representation of the visual context.

Considering heterogeneity between images and their related semantic knowledge, instead of the direct merging operation used in the textual context learning, we resort to the dot-product attention mechanism [29], which has been proven to be effective in many multimodal tasks, such as multimodal abstractive summarization [24], task-oriented language grounding [34], and video editing [35]. To be specific, given the related knowledge \mathcal{A}_j^v of the image v_j , we first acquire the knowledge embeddings $\mathbf{K}_v^j \in \mathbb{R}^{N_j^v \times D}$ by

the position-wise embedding layer of BART in Eqns. (1) and (2). N_v^j is the total number of tokens in \mathcal{A}_j^v . Thereafter, we adopt the dot-product attention mechanism to distinguish informative knowledge tokens towards the representation of v_j . Formally, we can obtain the knowledge-enhanced visual representation $\tilde{\mathbf{v}}_j$ of the image v_j as follows,

$$\begin{cases} \tilde{\mathbf{v}}_j = \mathbf{v}_j^\top \mathbf{W}_v^k, \\ \tilde{\mathbf{K}}_v^j = \mathbf{K}_v^j \mathbf{W}_k^k, \\ \mathbf{a}_j = \text{softmax}(\tilde{\mathbf{v}}_j (\tilde{\mathbf{K}}_v^j)^\top), \\ \tilde{\mathbf{v}}_j = \text{LN}(\mathbf{v}_j + (\mathbf{a}_j \mathbf{K}_v^j)^\top), \end{cases} \quad (9)$$

where \mathbf{W}_v^k and \mathbf{W}_k^k are the to-be-learned transformation matrices, which aim to project the visual representation (*i.e.*, \mathbf{v}_j) and knowledge embeddings (*i.e.*, \mathbf{K}_v^j) into the same space, and obtain the corresponding latent representations (*i.e.*, $\tilde{\mathbf{v}}_j$ and $\tilde{\mathbf{K}}_v^j$). $\mathbf{a}_j \in \mathbb{R}^{N_v^j}$ is the confidence vector, which denotes different confidence levels of tokens in the knowledge \mathcal{A}_j^v towards the image representation v_j . $\text{softmax}(\cdot)$ denotes the softmax activation function. $\text{LN}(\cdot)$ represents the layer normalization operation, which contributes to enhancing the model generalization ability. Ultimately, we use $\tilde{\mathbf{E}}_v = [\tilde{\mathbf{v}}_1; \tilde{\mathbf{v}}_2; \dots; \tilde{\mathbf{v}}_{N_v}]^\top \in \mathbb{R}^{N_v \times D}$ to denote the knowledge-enhanced representation of all images in the dialog context.

4.3.2 Dual Cross-modal Representation Refinement

In multimodal dialog systems, as both modalities serve to express the same user’s intention, it is promising to learn the context of one modality by referring to the context of the other modality. For example, as depicted in Figure 1, the user exhibits his/her intention of finding a restaurant and a shopping mall with multimodal input, including the textual description (*e.g.*, ‘an udon restaurant’, ‘in orchard road’, and ‘a shopping mall near it’), and the image for intuitively showing his/her desired shopping mall. To fully leverage the semantic relation between the textual context and the visual context to enhance the user intention understanding, we devise the dual cross-modal representation refinement component, with both the *vision-oriented representation refinement* and the *text-oriented representation refinement* modules.

Vision-oriented Representation Refinement. In this module, we aim to enhance the visual context representation by referring to the textual modality. Towards this end, we utilize the dot-product attention mechanism to highlight the informative tokens in the textual context to refine the visual representation. Specifically, we first obtain the embedding matrix of the textual context $\mathbf{E}_c = [\mathbf{t}_1; \mathbf{t}_2; \dots; \mathbf{t}_{N_T}]^\top \in \mathbb{R}^{N_T \times D}$ by the position-wise embedding layer of BART in Eqns. (1) and (2). \mathbf{t}_g is the embedding of t_g , and N_T is the number of tokens in the context. Then, the vision-oriented representation refinement can be denoted as follows,

$$\begin{cases} \mathbf{o}_j = \text{softmax}(\tilde{\mathbf{v}}_j^\top \mathbf{W}_v^v (\mathbf{E}_c \mathbf{W}_c)^\top), \\ \mathbf{P}_j = [[\mathbf{t}_1; \mathbf{v}_j]; [\mathbf{t}_2; \mathbf{v}_j]; \dots; [\mathbf{t}_{N_T}; \mathbf{v}_j]], \quad j = 1, 2, \dots, N_v, \\ \tilde{\mathbf{v}}_j = \mathbf{o}_j \mathbf{P}_j^\top, \end{cases} \quad (10)$$

where \mathbf{W}_v^v and \mathbf{W}_c are to-be-learned weight matrices to project different modalities representations into the same

semantic space. $\mathbf{o}_j \in \mathbb{R}^{N_T}$ is the confidence vector to indicate the confidence of tokens in the textual context towards the image v_j . Inspired by [36], instead of simply using the textual embedding vector (*i.e.*, \mathbf{t}_g) to derive the refined visual representation of v_j , we integrate the original visual representation \mathbf{v}_j to get the enhanced representation of v_j , *i.e.*, $\tilde{\mathbf{v}}_j$. Let $\hat{\mathbf{E}}_v = [\hat{\mathbf{v}}_1; \hat{\mathbf{v}}_2; \dots; \hat{\mathbf{v}}_{N_v}] \in \mathbb{R}^{N_v \times 2D}$ denote the refined visual context representation matrix.

Text-oriented Representation Refinement. Analogously, to refine the text context representation learning by referring to the visual modality, we conduct the text-oriented representation refinement. To be specific, we also employ the dot-product attention mechanism to distinguish informative images in the visual context to refine the textual representation. Thereafter, we concatenate the textual context representation \mathbf{T}_t and corresponding distinguished vision representation, and then utilize a fully connected layer to get the final context representation $\mathbf{T}_E \in \mathbb{R}^{N_t \times D}$ as follows,

$$\begin{cases} \bar{\mathbf{T}}_t = \mathbf{T}_t \mathbf{W}_t, \\ \bar{\mathbf{E}}_v = \hat{\mathbf{E}}_v \mathbf{W}_v^k, \\ \mathbf{S}_E = \text{softmax}(\bar{\mathbf{T}}_t \bar{\mathbf{E}}_v^\top), \\ \mathbf{T}_E = [\mathbf{T}_t; \mathbf{S}_E \bar{\mathbf{E}}_v] \mathbf{W}_f, \end{cases} \quad (11)$$

where \mathbf{W}_t , \mathbf{W}_v^k , and \mathbf{W}_f are the to-be-learned matrices, and $\bar{\mathbf{T}}_t \in \mathbb{R}^{N_t \times D}$ and $\bar{\mathbf{E}}_v \in \mathbb{R}^{N_v \times D}$ are the transferred representation of the textual context and visual context, respectively. $\mathbf{S}_E \in \mathbb{R}^{N_t \times N_v}$ is the confidence matrix, whose (q, j) -th entry denotes the confidence of the j -th image v_j towards reflecting the q -th token x_t^q .

4.4 Knowledge-enhanced Response Generation

By now, we have obtained the knowledge-enhanced context representation and can move forward to generate the target response. In particular, we employ the decoder of BART (*i.e.*, \mathcal{B}_d) in Eqn. (4) as our decoder. Although the origin decoder is feasible, the obvious drawback lies in that it neglects to explicitly exploit the knowledge for the response generation in the multimodal task-oriented dialog systems.

In light of this, we revise the original decoder \mathcal{B}_d by introducing a dot-product knowledge-decoder attention (DKDA) sub-layer, which can distinguish the informative tokens of the context related knowledge and adaptively utilizes the knowledge to facilitate the text response generation. To be specific, we insert the DKDA sub-layer between the MMSA and the MEDA as follows,

$$\begin{cases} \mathbf{q}_l^S = \text{LN}(\text{MMSA}(\mathbf{q}_{l-1}^{dec} + \mathbf{q}_{l-1}^{dec}), \\ \mathbf{q}_l^K = \text{LN}(\text{DKDA}(\mathbf{q}_l^S, \mathbf{E}_k) + \mathbf{q}_l^S), \\ \mathbf{q}_l^E = \text{LN}(\text{MEDA}(\mathbf{q}_l^K, \mathbf{Z}_L^{enc}) + \mathbf{q}_l^K), \\ \mathbf{q}_l^{dec} = \text{LN}(\text{FFN}(\mathbf{q}_l^E) + \mathbf{q}_l^E) \end{cases}, \quad l = 1, 2, \dots, L, \quad (12)$$

where $\mathbf{E}_k \in \mathbb{R}^{N_k^d \times D}$ denotes the embedding of related knowledge $\mathcal{K}_d = \mathcal{K}_t^A \cup \mathcal{K}_v^A$, which can be derived by the position-wise embedding layer of BART in Eqns. (1) and (2). N_k^d is the total number of tokens in \mathcal{K}_d . Notably, $\mathbf{Z}_L^{enc} = \mathbf{T}_E$, which can be obtained by Eqn. (11).

TABLE 1
Detailed statistics of the MMConv dataset.

Entry	Number
#dialogues	5, 106
#turns	39, 759
#single-modality dialogues	751
#multi-modality dialogues	4, 355
#single-domain dialogues	808
#multi-domain dialogues	4, 298
#entities in the knowledge base	1, 771

Thereinto, considering different knowledge tokens may contribute differently in promoting the target response generation, we define the DKDA sub-layer as follows,

$$\begin{cases} \bar{\mathbf{q}} = \mathbf{q}^\top \mathbf{W}_d, \\ \bar{\mathbf{E}}_k = \mathbf{E}_k \mathbf{W}_d^k, \\ \mathbf{a}_d = \text{softmax}(\bar{\mathbf{q}}(\bar{\mathbf{E}}_k)^\top), \\ \mathbf{q}_k = LN(\mathbf{q} + (\mathbf{a}_d \mathbf{E}_k)^\top), \end{cases} \quad (13)$$

where for simplicity, we temporarily omit the subscripts l and n that index the decoder layer and the decoding step, respectively. $\mathbf{q} \in \mathbb{R}^D$ refers to the output of the masked multi-head self-attention layer (i.e., \mathbf{q}_i^S in Eqn.(12)). \mathbf{W}_d and \mathbf{W}_d^k are to-be-learned matrices projecting \mathbf{q} and \mathbf{E}_k into the same space. $\bar{\mathbf{q}} \in \mathbb{R}^{1 \times D}$ and $\bar{\mathbf{E}}_k \in \mathbb{R}^{N_k^d \times D}$ are the projected latent representation of \mathbf{q} and \mathbf{E}_k , respectively. $\mathbf{a}_d \in \mathbb{R}^{N_k^d}$ is the confidence vector to indicate different levels confidences of tokens in \mathcal{K}_d . \mathbf{q}_k (i.e., the output of the DKDA sub-layer) denotes the knowledge-enhanced decoder representation.

In a nutshell, for each time step, we can obtain its corresponding predicted distribution $\tilde{\mathbf{y}}$ according to Eqn.(5), and thus capture the predicted token of the current time step based on the largest element of $\tilde{\mathbf{y}}$. Ultimately, we adopt the cross entropy loss [37] to supervise the response generation as follows,

$$\mathcal{L}_{CE} = -\frac{1}{N_R} \sum_{n=1}^{N_R} \log(\tilde{\mathbf{y}}_n[t*]), \quad (14)$$

where $\tilde{\mathbf{y}}_n[t*]$ refers to the element of $\tilde{\mathbf{y}}_n$ that corresponds to the n -th token of the ground truth response \mathcal{R} , and N_R is the total number of tokens in \mathcal{R} . Notably, the loss is defined for a single sample.

5 EXPERIMENT

In this section, we first introduce the dataset as well as the experiment setting, and then detail the experiments by answering the following research questions:

- **RQ1:** Does DKMD surpass state-of-the-art methods?
- **RQ2:** How does the knowledge affect the DKMD?
- **RQ3:** How does the dual cross-modal representation refinement influence the DKMD?
- **RQ4:** Is DKMD sensitive to the location of the encoder layer incorporating the dual knowledge-enhanced context learning?

5.1 Dataset

As a matter of fact, existing efforts evaluate their models on the publicly available dataset MMD built by Saha et al. [1] from the fashion domain. However, in this work,

TABLE 2
Performance comparison among different methods in terms of BLEU- N (%) and Nist.

Models	BLEU-1	BLEU-2	BLEU-3	BLEU-4	Nist
MHRED	15.02	6.66	4.24	2.94	0.9529
KHRED	18.29	8.28	4.98	3.36	1.1189
LARCH	20.86	11.33	7.58	5.58	1.3400
MATE	30.45	22.06	17.05	13.41	2.3426
UMD	31.14	21.87	17.12	13.82	2.5290
TREASURE	34.75	24.82	18.67	14.53	2.4398
DKMD	39.59	31.95	27.26	23.72	4.0004

we did not choose this dataset for evaluation due to the fact that MMD only allows the knowledge referring by the visual dialog context. Instead, towards comprehensive knowledge referring, we employed the more recently released public dataset MMConv [38], which is constructed from the general domain and supports knowledge selection from both modalities. The MMConv dataset contains 5,106 conversations between users and agents spanning five domains: *Food*, *Hotel*, *Nightlife*, *Shopping mall* and *Sightseeing*. Thereinto, the number of single-modality and multi-modality dialogues in the MMConv dataset are 751 and 4,355, respectively, where the corresponding average number of turns are 7.1 and 7.9. In addition, the knowledge base of MMConv involves 1,771 knowledge entities, each of which involves a set of attributes and a few images. More detailed information about the MMConv dataset is summarized in Table 1.

5.2 Experiment Setting

We followed the original setting in MMConv [38], which divides dialogues into three chunks: 3,500 for training, 606 for validation, and 1,000 for testing. Following the former studies [5], [7], we treated every utterance of agents in the conversations as a target response and utilized its former two-turn utterances as the given context. We employed the pretrained BART-large³ model with 12 layers for encoder and decoder, respectively. For optimization, we utilized the adaptive moment estimation (Adam) optimizer and set the learning rate as $1e-5$. Moreover, we fine-tuned the proposed DKMD on the basis of the training and validation dataset with 100 epochs, and reported the performance on the testing dataset. In addition, we implemented our DKMD by Pytorch [39] and conducted all experiments on a server equipped with 8 NVIDIA A100 GPUs. Following existing methods [2], [7], [8], we adopted BLEU- N [40] where N varies from 1 to 4, and Nist [41] as evaluation metrics. In particular, both BLEU- N and Nist can measure the similarity between the generated and target responses, and higher BLEU- N and Nist scores denote the more n -gram overlap between the generated and target responses [8].

5.3 Model Comparison (RQ1)

To verify the effectiveness of our proposed DKMD, we chose the following state-of-the-art methods on multimodal dialog systems as baselines.

- **MHRED** [1]. This is the first work on the multimodal task-oriented dialog systems, which consists of a

3. <https://huggingface.co/facebook/bart-large>.

hierarchical encoder and a GRU-based decoder. The hierarchical encoder contains two levels of the gated recurrent units (GRU) [42], corresponding to encode the utterance and the context, respectively. Notably, it neglects the knowledge base and the semantic relation in the multimodal context.

- **KHRED**. Considering the vital role of knowledge in multimodal task-oriented dialog systems, we designed this baseline by incorporating the knowledge into MHRED [1]. To be specific, following the knowledge integration way [3], we utilized the memory network to encode the attribute knowledge, and fed the knowledge representation into the GRU-based decoder to generate text responses.
- **LARCH** [3] designs a hierarchical graph-based neural network to model the semantic relation in the given multimodal context, where each word, image, sentence, utterance, dialog pair, and the entire session are treated as nodes. Similar to KHRED, it also integrates the attribute knowledge via the memory network and adopts the GRU-based decoder for response generation.
- **MATE** [5] introduces the Transformer network [29] to capture the context semantic relation (*i.e.*, semantic dependencies) between the textual context and the visual context, and devises the Transformer-based decoder to generate the text response. Notably, it neglects the attribute knowledge.
- **UMD** [6] adopts the common multimodal hierarchical encoder-decoder model. In particular, it designs a hierarchy-aware tree encoder to learn the attribute-level visual representation, and devises the multimodal factorized bilinear pooling layer to model the semantic relation hidden in the multimodal context. Notably, the attribute knowledge stimulates the visual representation learning.
- **TREASURE** [2] introduces an attribute-enhanced textual encoder, which can adaptively focus on attribute-related keywords and obtain the utterance representation. Besides, this baseline designs a sparse graph attention network to learn the semantic relation and adaptively aggregate the context information. Notably, this method utilizes attribute-related keywords to enhance the textual representation learning.

Table 2 illustrates the performance comparison among different models with regard to different evaluation metrics. From this table, we make the following observations: 1) DKMD consistently outperforms all the baselines across different evaluation metrics, indicating the effectiveness of our proposed DKMD. This suggests that it is reasonable to incorporate the generative pretrained language model as well as the multimodal context related knowledge in multimodal dialog systems. 2) DKMD surpasses all the baselines that also consider knowledge (*i.e.*, TREASURE, UMD, LARCH, and KHRED), which confirms the advantage of our knowledge incorporation manner, *i.e.*, the local-wise and global-wise knowledge enhancement in the multimodal context encoding part, as well as the explicit enhancement in the decoding part. 3) MHRED gets the

TABLE 3
Ablation study results on knowledge in terms of BLEU- N (%) and Nist.

Models	BLEU-1	BLEU-2	BLEU-3	BLEU-4	Nist
w/o-GlobalK-All	34.03	25.34	20.61	17.25	3.1755
w/o-GlobalK-OnlyV	38.77	30.97	26.25	22.73	3.9112
w/o-LocalK	38.68	31.13	26.55	23.11	3.8015
w-LocalK-AddT	30.32	22.64	18.51	15.67	3.7723
w/o-DKDA	38.17	30.56	26.03	22.64	3.8678
w/o-K-All	29.99	21.67	17.30	14.23	2.7053
w/o-TextualK-All	34.05	26.21	21.78	18.54	3.2022
w/o-VisualK-All	38.46	30.98	26.46	23.07	3.8442
DKMD	39.59	31.95	27.26	23.72	4.0004

worst performance compared to other methods. This may be attributed to that MHRED not only neglects context related knowledge but also overlooks the semantic relation in the multimodal context, which limits its capability of accurately capturing the user’s intention. 4) DKMD, TREASURE, UMD, MATE, and LARCH exceed all the baselines neglecting the context semantic relation (*i.e.*, KHRED and MHRED), which reflects the necessity of exploring the semantic relation hidden in the multimodal dialog context.

5.4 Effect of Knowledge (RQ2)

To thoroughly verify the effect of knowledge in multimodal task-oriented dialog systems, we devised eight derivations as follows.

1) **w/o-GlobalK-All**. To demonstrate the effect of the knowledge in the global text-based learning, we removed all the related knowledge and only used the textual context as the input of Eqn.(7).

2) **w/o-GlobalK-OnlyV**. To illustrate the necessity of visual context related knowledge in the global text-based learning, we disabled the visual context related knowledge input of Eqn.(7).

3) **w/o-LocalK**. To show the benefit of the knowledge in the local vision-based learning, we removed the knowledge refined visual representation obtained by Eqn.(9) and only used the original visual representation extracted by CLIP (*i.e.*, Eqn.(8)).

4) **w-LocalK-AddT**. To verify that the textual context related knowledge is redundant in the local vision-based learning, we modified the input of Eqn.(9) (*i.e.*, \mathbf{K}_v^j) by concatenating the textual context related knowledge with the visual context related knowledge.

5) **w/o-DKDA**. To indicate the necessity of injecting the explicit knowledge into the decoder, we discarded the related knowledge (*i.e.*, the DKDA sub-layer) in the knowledge-enhanced response generation. Namely, we utilized the original decoder of BART (*i.e.*, \mathcal{B}_d).

6) **w/o-K-All**. To verify the importance of the knowledge, we disabled all the knowledge in our proposed DKMD.

7) **w/o-TextualK-All**. To illustrate the importance of textual context related knowledge in the whole network, we disabled the text-based knowledge and only kept the vision-based knowledge in the global text-based learning (*i.e.*, removing \mathcal{K}_t^A in \mathcal{X}_t for Eqn.(7)) and knowledge-enhanced response generation (*i.e.*, disabling \mathcal{K}_t^A in \mathcal{K}_d for Eqn.(13)).

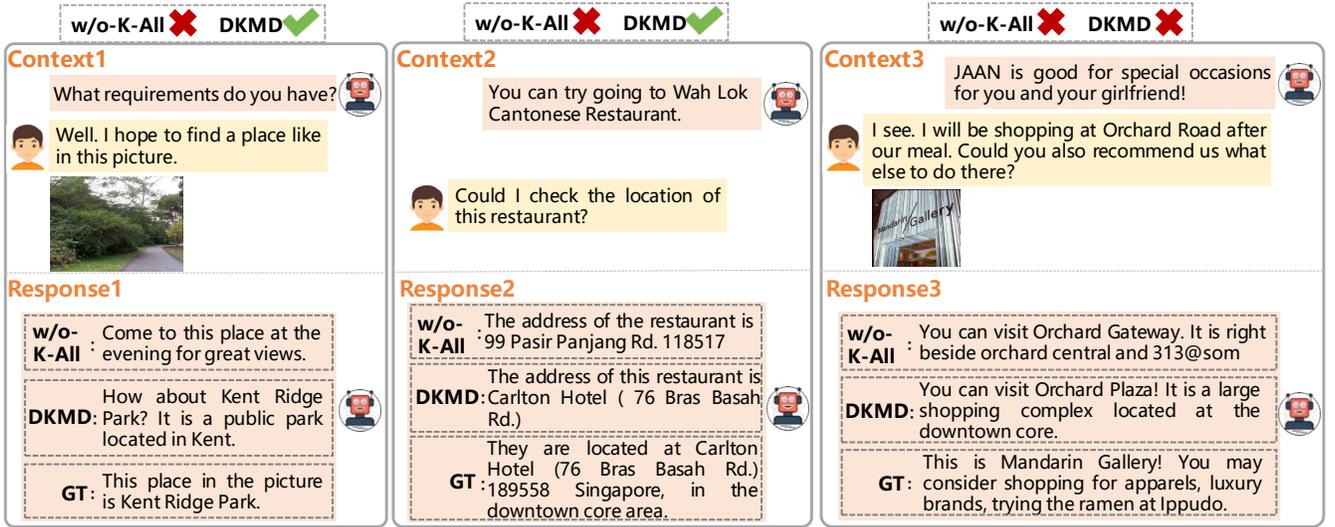


Fig. 3. Comparison between DKMD and w/o-K-All on several testing dialog pairs. “GT” refers to the groundtruth. We represent the correct response of the model with the green tick and the wrong one with the red cross.

8) **w/o-VisualK-All**. To show the roles of visual context related knowledge, we removed the vision-based knowledge in the global text-based learning (*i.e.*, removing \mathcal{K}_v^A in \mathcal{X}_t for Eqn.(7)), local vision-based learning (*i.e.*, Eqn.(9)), and knowledge-enhanced response generation (*i.e.*, disabling \mathcal{K}_v^A in \mathcal{K}_d for Eqn.(13)).

Table 3 shows the performance of DKMD and its above derivations. From this table, we have the following observations. 1) DKMD outperforms all w/o-GlobalK-All, w/o-LocalK, w/o-DKDA, and w/o-K-All. In addition, removing all the knowledge (*i.e.*, w/o-K-All) leads to the worst performance. It indicates that disabling the knowledge anywhere (*i.e.*, global knowledge-enhanced textual representation learning, local knowledge-enhanced visual representation learning, or knowledge-enhanced response generation) hurts the performance of DKMD. This may be attributed to that the knowledge anywhere complement each other and all contribute to the text response generation in the context of multimodal task-oriented dialog systems. 2) w/o-GlocalK-OnlyV performs worse than DKMD. It demonstrates the effectiveness of the visual context related knowledge in the global text-based learning. This may be due to that the visual context related knowledge can enhance the global textual context learning and supplement the user intention modeling from the visual perspective, thus boosting the performance of text response generation. 3) w-LocalK-AddT underperforms DKMD, which suggests the redundancy of considering the textual context related knowledge in the local vision-based learning. One plausible explanation is that merging the

textual context related knowledge may bring the noise (*e.g.*, the visual context irrelevant knowledge entities appeared in the textual context) to the visual context learning and thus hurt the performance. 4) Both w/o-TextualK-All and w/o-VisualK-All perform worse than DKMD, confirming the superiority of considering the knowledge from both the text and vision perspectives. 5) w/o-TextualK-All underperforms w/o-VisualK-All, which implies that the text-based knowledge contributes more than vision-based knowledge. One possible explanation is that the text-based knowledge derived from the global textual context may be more comprehensive and exert a larger role in facilitating the user intention modeling than the vision-based knowledge.

To gain more deep insights into the influence of knowledge, we showed the comparison between DKMD and w/o-K-All on three testing dialog pairs in Figure 3. As we can see, DKMD performs better than w/o-K-All in *context1* and *context2* when the knowledge is indispensable to understanding the user’s intention. For example, we found that the *context1* and *context2* involve vision-based knowledge and text-based knowledge, respectively, and DKMD can generate proper responses and provide accurate

TABLE 4
Ablation study results on dual cross-modal representation refinement in terms of BLEU-N (%) and Nist.

Models	BLEU-1	BLEU-2	BLEU-3	BLEU-4	Nist
DKMD-w/o-V	38.25	30.79	26.34	23.00	3.8046
DKMD-w/o-Dual	38.24	30.17	25.14	21.88	3.8662
DKMD-w/o-VR	38.50	30.74	26.09	22.61	3.9252
DKMD-w/o-TR	38.52	31.05	26.56	23.19	3.8961
DKMD	39.59	31.95	27.26	23.72	4.0004

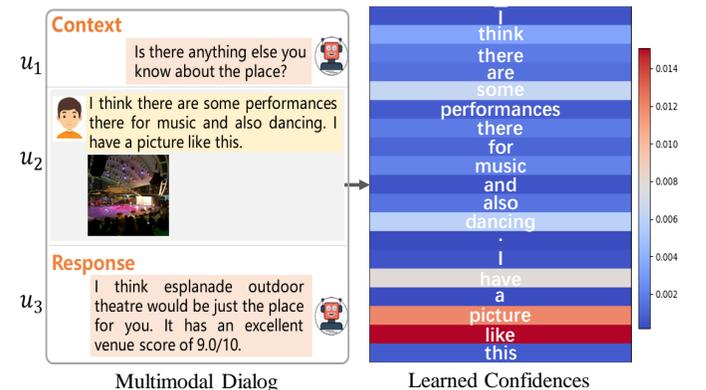


Fig. 4. Visualization of the learned confidences in the vision-oriented representation refinement concerning the given multimodal dialog.

information while w/o-K-All fails that. This phenomenon validates the advantage of incorporating the knowledge into the text response generation in the context of multimodal task-oriented dialog systems. Nevertheless, DKMD can also yield the failed cases, such as the *context3* in Figure 3. To be specific, we found that DKMD recommends the wrong shopping mall “Orchard Plaza” in the *context3*, which shares similar properties (e.g., shopping mall, at Orchard Road) with the accurate one “Mandarin Gallery”.

5.5 Effect of Dual Cross-modal Representation Refinement (RQ3)

To explore roles of the dual cross-modal representation refinement, we conduct the comparative experiment with the following derivatives: **DKMD-w/o-Dual**, **DKMD-w/o-VR** and **DKMD-w/o-TR**, where dual cross-modal (i.e., both vision and text), vision-oriented and text-oriented representation refinement are removed, respectively. In addition, we also introduce the derivative **DKMD-w/o-V** by removing all the visual information (i.e., visual context and knowledge) to verify the necessity of integrating visual information.

Table 4 shows the performance comparison between DKMD and its derivatives. As can be seen, DKMD-w/o-Dual deviates more than DKMD, as compared with DKMD-w/o-VR, and DKMD-w/o-TR. Based on the phenomenon, we had the following two observations. 1) The dual cross-modal representation refinement does benefit the user intention modeling and promotes the text response generation in the context of multimodal task-oriented dialog systems. 2) The vision-oriented representation refinement and text-oriented representation refinement complement each other and both contribute to modeling the user’s intention. In addition, we found that DKMD-w/o-VR underperforms DKMD-w/o-TR, denoting that the vision-oriented representation refinement contributes more than the text-oriented representation refinement. This may be due to that the textual context may convey massive information and play a vital role in delivering the user’s intention. Therefore, introducing the vision-oriented refinement by referring to the textual modality can significantly improve the visual context understanding. In contrast, the visual context usually expresses some specific user’s intention, like the desired food, thus contributing little to refining the textual context representation and promoting the user’s intention learning. Last but not least, DKMD-w/o-V performs worse than DKMD, denoting that the visual information plays a pivotal role in the textual response generation of multimodal task-oriented dialog systems. The underlying philosophy is that the visual information can convey essential clues and thus facilitates the user intention understanding.

To gain a better understanding of the dual cross-modal representation refinement, as shown in Figure 4, we randomly selected a testing multimodal dialog and illustrated the learned confidences in the vision-oriented representation refinement (i.e., \mathbf{o}_j in Eqn.(10)) of the utterance u_2 with a thermodynamic diagram. The color of the bar denotes the confidence of tokens towards the image, where the color approaching the orange refers to the larger weight. As can be seen from Figure 4,

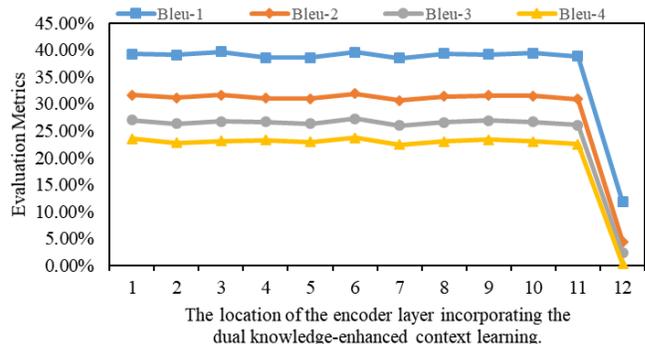


Fig. 5. Sensitivity analysis on the location of the encoder layer incorporating the dual knowledge-enhanced context learning.

the vision-oriented representation refinement does assign different levels of confidences to different tokens for the given image. As we can see, our model does identify the informative tokens, such as “dancing”, “have”, “picture” and “like”, and assigns smaller weights to tokens possessing smaller semantic relation with the image (e.g., “I”, “and”, “.”, and “a”). This suggests that the semantic relation does exist in the multimodal context and the dual cross-modal representation refinement can well capture it.

5.6 Sensitivity Analysis

As there is a stack of layers in the BART encoder, in this part, we performed the sensitivity analysis on the location of the encoder layer incorporating the dual knowledge-enhanced context learning.

As shown in Figure 5, we enumerated the performance of each layer (i.e., from 1-st to 12-th) in the encoder to perform the dual knowledge-enhanced context learning. As can be seen, our proposed DKMD performs relatively stably when it integrates the dual knowledge-enhanced context learning at an arbitrary layer between the 1-st encoder layer to 11-th encoder layer. Thereinto, our proposed DKMD achieves the optimal performance when it integrates the dual knowledge-enhanced context learning in the 6-th encoder layer. Interestingly, we observed that the proposed DKMD performs the worst when it integrates the dual knowledge-enhanced context learning at the 12-th layer. This suggests that incorporating the dual knowledge-enhanced context learning in the last BART encoder layer will hurt the original data distribution and harm the capacity of the BART encoder.

6 CONCLUSION

In this paper, we tackle the textual response generation task in multimodal task-oriented dialog systems based on GPLMs. In particular, we propose a novel dual knowledge-enhanced generative pretrained language model for multimodal dialog systems, named DKMD, which consists of three pivotal components: *dual knowledge selection*, *dual knowledge-enhanced context learning*, and *knowledge-enhanced response generation*. Extensive experiments on a public dataset well validate our proposed DKMD and demonstrate the necessity of incorporating the GPLMs and the multimodal context related knowledge in multimodal task-oriented dialog systems. In addition,

we observe that the textual context related knowledge and the visual context related knowledge complement each other and both contribute to the textual response generation. Besides, the semantic relation in the multimodal context does exist and should be taken into account. Currently, we mainly investigate the potential of GPLMs in the textual response generation task of multimodal task-oriented dialog systems, but ignore the cross-domain relation among different domains. In the future, we plan to turn to cross-domain multimodal dialog systems to explore the semantic relation among domains.

REFERENCES

- [1] A. Saha, M. M. Khapra, and K. Sankaranarayanan, "Towards building large scale multimodal domain-aware conversation systems," in *Proceedings of the AAAI Conference on Artificial Intelligence*. AAAI Press, 2018, pp. 696–704.
- [2] H. Zhang, M. Liu, Z. Gao, X. Lei, Y. Wang, and L. Nie, "Multimodal dialog system: Relational graph-based context-aware question understanding," in *Proceedings of the ACM International Conference on Multimedia*. ACM, 2021, pp. 695–703.
- [3] L. Nie, F. Jiao, W. Wang, Y. Wang, and Q. Tian, "Conversational image search," *IEEE Transactions on Image Processing*, vol. 30, pp. 7732–7743, 2021.
- [4] H. Chauhan, M. Firdaus, A. Ekbal, and P. Bhattacharyya, "Ordinal and attribute aware response generation in a multimodal dialogue system," in *Proceedings of the Conference of the Association for Computational Linguistics*. Association for Computational Linguistics, 2019, pp. 5437–5447.
- [5] L. Liao, Y. Ma, X. He, R. Hong, and T. Chua, "Knowledge-aware multimodal dialogue systems," in *Proceedings of the ACM International Conference on Multimedia*. ACM, 2018, pp. 801–809.
- [6] C. Cui, W. Wang, X. Song, M. Huang, X. Xu, and L. Nie, "User attention-guided multimodal dialog systems," in *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 2019, pp. 445–454.
- [7] L. Nie, W. Wang, R. Hong, M. Wang, and Q. Tian, "Multimodal dialog system: Generating responses via adaptive decoders," in *Proceedings of the ACM International Conference on Multimedia*. ACM, 2019, pp. 1098–1106.
- [8] W. He, Z. Li, D. Lu, E. Chen, T. Xu, B. Huai, and J. Yuan, "Multimodal dialogue systems via capturing context-aware dependencies of semantic elements," in *Proceedings of the ACM International Conference on Multimedia*. ACM, 2020, pp. 2755–2764.
- [9] Z. Ma, J. Li, G. Li, and Y. Cheng, "UniTranSeR: A unified transformer semantic representation framework for multimodal task-oriented dialog system," in *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2022, pp. 103–114.
- [10] J. Devlin, M. Chang, K. Lee, and K. Toutanova, "BERT: pre-training of deep bidirectional transformers for language understanding," in *Proceedings of the Conference of the North American Chapter of the Association for Computational Linguistics*. Association for Computational Linguistics, 2019, pp. 4171–4186.
- [11] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, and P. J. Liu, "Exploring the limits of transfer learning with a unified text-to-text transformer," *Journal of Machine Learning Research*, vol. 21, pp. 140:1–140:67, 2020.
- [12] M. Lewis, Y. Liu, N. Goyal, M. Ghazvininejad, A. Mohamed, O. Levy, V. Stoyanov, and L. Zettlemoyer, "BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension," in *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2020, pp. 7871–7880.
- [13] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," in *Proceedings of the International Conference on Neural Information Processing Systems*. MIT Press, 2014, p. 3104–3112.
- [14] L. Liao, R. Takano, Y. Ma, X. Yang, M. Huang, and T. Chua, "Topic-guided conversational recommender in multiple domains," *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 5, pp. 2485–2496, 2022.
- [15] Y. Li, R. Zhang, W. Li, and Z. Cao, "Hierarchical prediction and adversarial learning for conditional response generation," *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 1, pp. 314–327, 2022.
- [16] W. Lei, X. Jin, M. Kan, Z. Ren, X. He, and D. Yin, "Sequicity: Simplifying task-oriented dialogue systems with single sequence-to-sequence architectures," in *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. Association for Computational Linguistics, 2018, pp. 1437–1447.
- [17] Y. Zhang, Z. Ou, and Z. Yu, "Task-oriented dialog systems that consider multiple appropriate responses under the same context," in *Proceedings of the AAAI Conference on Artificial Intelligence*. AAAI Press, 2020, pp. 9604–9611.
- [18] X. Li, Y. Chen, L. Li, J. Gao, and A. Celikyilmaz, "End-to-end task-completion neural dialogue systems," in *Proceedings of the Eighth International Joint Conference on Natural Language Processing*. Asian Federation of Natural Language Processing, 2017, pp. 733–743.
- [19] J. Wang, J. Liu, W. Bi, X. Liu, K. He, R. Xu, and M. Yang, "Dual dynamic memory network for end-to-end multi-turn task-oriented dialog systems," in *Proceedings of the International Conference on Computational Linguistics*. International Committee on Computational Linguistics, 2020, pp. 4100–4110.
- [20] Z. Liu, D. Zhou, H. Liu, H. Wang, Z.-Y. Niu, H. Wu, W. Che, T. Liu, and H. Xiong, "Graph-grounded goal planning for conversational recommendation," *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–15, 2022.
- [21] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, 2017, pp. 5998–6008.
- [22] H. Chen, C. Zhang, J. Li, P. S. Yu, and N. Jing, "Kggen: A generative approach for incipient knowledge graph population," *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 5, pp. 2254–2267, 2022.
- [23] M. Zhang, S. Wu, M. Gao, X. Jiang, K. Xu, and L. Wang, "Personalized graph neural networks with attention mechanism for session-aware recommendation," *IEEE Transactions on Knowledge and Data Engineering*, vol. 34, no. 8, pp. 3946–3957, 2022.
- [24] T. Yu, W. Dai, Z. Liu, and P. Fung, "Vision guided generative pre-trained language models for multimodal abstractive summarization," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2021, pp. 3995–4007.
- [25] F. Zhang, J. Tang, X. Liu, Z. Hou, Y. Dong, J. Zhang, X. Liu, R. Xie, K. Zhuang, X. Zhang, L. Lin, and P. Yu, "Understanding wechat user preferences and "wow" diffusion," *IEEE Transactions on Knowledge and Data Engineering*, pp. 1–14, 2021.
- [26] T. Mikolov, I. Sutskever, K. Chen, G. Corrado, and J. Dean, "Distributed representations of words and phrases and their compositionality," in *Proceedings of the International Conference on Neural Information Processing Systems*. Curran Associates Inc., 2013, p. 3111–3119.
- [27] J. Pennington, R. Socher, and C. Manning, "GloVe: Global vectors for word representation," in *Proceedings of the Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 2014, pp. 1532–1543.
- [28] S. Gao, Y. Zhang, Y. Wang, Y. Dong, X. Chen, D. Zhao, and R. Yan, "Heteroqa: Learning towards question-and-answering through multiple information sources via heterogeneous graph modeling," in *Proceedings of the ACM International Conference on Web Search and Data Mining*. ACM, 2022, pp. 307–315.
- [29] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," *CoRR*, 2017.
- [30] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, "An image is worth 16x16 words: Transformers for image recognition at scale," in *Proceedings of the International Conference on Learning Representations*. OpenReview.net, 2021.
- [31] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark, G. Krueger, and I. Sutskever, "Learning transferable visual models from natural language supervision," in *Proceedings of the International Conference on Machine Learning*. PMLR, 2021, pp. 8748–8763.
- [32] Y. Shi, X. Yang, H. Xu, C. Yuan, B. Li, W. Hu, and Z. Zha, "Emscore:

Evaluating video captioning via coarse-grained and fine-grained embedding matching," *CoRR*, 2021.

- [33] J. Sun, Q. Deng, Q. Li, M. Sun, M. Ren, and Z. Sun, "Anyface: Free-style text-to-face synthesis and manipulation," *CoRR*, 2022.
- [34] D. S. Chaplot, K. M. Sathyendra, R. K. Pasumarthi, D. Rajagopal, and R. Salakhutdinov, "Gated-attention architectures for task-oriented language grounding," in *Proceedings of the AAAI Conference on Artificial Intelligence*. AAAI Press, 2018, pp. 2819–2826.
- [35] O. Fried, A. Tewari, M. Zollhöfer, A. Finkelstein, E. Shechtman, D. B. Goldman, K. Genova, Z. Jin, C. Theobalt, and M. Agrawala, "Text-based editing of talking-head video," *ACM Transactions on Graphics*, vol. 38, no. 4, pp. 68:1–68:14, 2019.
- [36] T. Fu, X. E. Wang, S. T. Grafton, M. P. Eckstein, and W. Y. Wang, "Language-based video editing via multi-modal multi-level transformer," *CoRR*, 2021.
- [37] C. H. Li and C. K. Lee, "Minimum cross entropy thresholding," *Pattern Recognition*, vol. 26, no. 4, pp. 617–625, 1993.
- [38] L. Liao, L. H. Long, Z. Zhang, M. Huang, and T. Chua, "Mmconv: An environment for multimodal conversational search across multiple domains," in *Proceedings of the International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM, 2021, pp. 675–684.
- [39] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimeshein, L. Antiga, A. Desmaison, A. Köpf, E. Z. Yang, Z. DeVito, M. Raison, A. Tejani, S. Chilamkurthy, B. Steiner, L. Fang, J. Bai, and S. Chintala, "Pytorch: An imperative style, high-performance deep learning library," in *Proceedings of the Annual Conference on Neural Information Processing Systems*, 2019, pp. 8024–8035.
- [40] K. Papineni, S. Roukos, T. Ward, and W. Zhu, "Bleu: a method for automatic evaluation of machine translation," in *Proceedings of the Annual Meeting of the Association for Computational Linguistics*. ACL, 2002, pp. 311–318.
- [41] G. Doddington, "Automatic evaluation of machine translation quality using n-gram co-occurrence statistics," in *Proceedings of the Second International Conference on Human Language Technology Research*. Morgan Kaufmann Publishers Inc., 2002, p. 138–145.
- [42] J. Chung, Ç. Gülçehre, K. Cho, and Y. Bengio, "Empirical evaluation of gated recurrent neural networks on sequence modeling," *CoRR*, 2014.



Xiaolin Chen received the M.S. degree from Shandong University in 2021. She is currently pursuing the Ph.D. degree with the School of Software, Shandong University. Her research primarily focuses on information retrieval and multimedia computing. She has published papers in the top venues, including ACM SIGIR, TOIS and IEEE TMM. Moreover, she has served as reviewers for conferences and journals, such as ACM MM and Neurocomputing.



Xuemeng Song received the B.E. degree from University of Science and Technology of China in 2012, and the Ph.D. degree from the School of Computing, National University of Singapore in 2016. She is currently an associate professor of Shandong University, Jinan, China. Her research interests include the information retrieval and social network analysis. She has published several papers in the top venues, such as ACM SIGIR, MM and TOIS. In addition, she has served as reviewers for many top conferences

and journals.



Liqiang Jing is a graduate student with the School of Computer Science and Technology, Shandong University. He received the B.E. degree in School of Computer Science and Technology from Hefei University of Technology, Anhui, in 2020. His research interests include multimodal learning and natural language processing.



Shuo Li is currently an undergraduate student with the School of Computer Science and Technology, Shandong University. His research interests include natural language processing, social media computing, and software engineering.



Linmei Hu received the PhD degree from Tsinghua University, in 2018. She is an associate professor with the School of Computer Sciences, Beijing University of Posts and Telecommunications. Her research interests include natural language processing and data mining. She was awarded Beijing Excellent PhD Student in 2018.



Liqiang Nie is currently the dean with the School of Computer Science and Technology, Harbin Institute of Technology (Shenzhen campus). He received his B.Eng. and Ph.D. degree from Xi'an Jiaotong University and National University of Singapore (NUS), respectively. His research interests lie primarily in multimedia content analysis and information retrieval. Dr. Nie has co-authored more than 100 CCF-A papers and 5 books, with 15k plus Google Scholar citations. He is IAPR Fellow and an AE of IEEE TKDE,

IEEE TMM, IEEE TCSVT, ACM ToMM, and Information Science. Meanwhile, he is the regular area chair or SPC of ACM MM, NeurIPS, IJCAI and AAAI. He is a member of ICME steering committee. He has received many awards over the past three years, like ACM MM and SIGIR best paper honorable mention in 2019, the AI 2000 most influential scholars 2020, SIGMM rising star in 2020, MIT TR35 China 2020, DAMO Academy Young Fellow in 2020, SIGIR best student paper in 2021, first prize of the provincial science and technology progress award in 2021 (rank 1), and provincial youth science and technology award in 2022. Some of his research outputs have been integrated into the products of Alibaba, Kwai, and other listed companies.