

DAM: Deliberation, Abandon and Memory Networks for Generating Detailed and Non-repetitive Responses in Visual Dialogue

Xiaoze Jiang^{1,2*}, Jing Yu^{1,3*†}, Yajing Sun^{1,3},
Zengchang Qin^{2,4}, Zihao Zhu^{1,3}, Yue Hu^{1,3} and Qi Wu⁵

¹Institute of Information Engineering, Chinese Academy of Sciences, Beijing, China

²Intelligent Computing and Machine Learning Lab, School of ASEE, Beihang University, Beijing, China

³School of Cyber Security, University of Chinese Academy of Sciences, Beijing, China

⁴AI Research, Codemao Inc.

⁵University of Adelaide, Australia

{yujing02, sunyajing, zhuzihao, huyue}@iie.ac.cn, {xzjiang, zcquin}@buaa.edu.cn,
qi.wu01@adelaide.edu.au

Abstract

Visual Dialogue task requires an agent to be engaged in a conversation with human about an image. The ability of generating detailed and non-repetitive responses is crucial for the agent to achieve human-like conversation. In this paper, we propose a novel generative decoding architecture to generate high-quality responses, which moves away from decoding the whole encoded semantics towards the design that advocates both transparency and flexibility. In this architecture, word generation is decomposed into a series of attention-based information selection steps, performed by the novel recurrent Deliberation, Abandon and Memory (DAM) module. Each DAM module performs an adaptive combination of the response-level semantics captured from the encoder and the word-level semantics specifically selected for generating each word. Therefore, the responses contain more detailed and non-repetitive descriptions while maintaining the semantic accuracy. Furthermore, DAM is flexible to cooperate with existing visual dialogue encoders and adaptive to the encoder structures by constraining the information selection mode in DAM. We apply DAM to three typical encoders and verify the performance on the VisDial v1.0 dataset. Experimental results show that the proposed models achieve new state-of-the-art performance with high-quality responses. The code is available at <https://github.com/JXZe/DAM>.

1 Introduction

Visual Dialogue [Das *et al.*, 2017] is a task that requires an agent to answer a series of questions grounded in an image, demanding the agent to reason about both visual content and

*Equal contribution. This work is done when Xiaoze Jiang is an intern in IIE, CAS.

†Corresponding author.

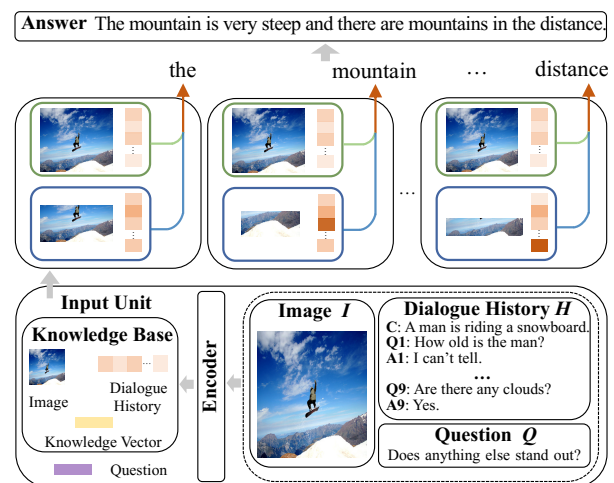


Figure 1: An illustration of DAM. The encoder encodes visual dialogue information into Knowledge Base (KB). DAM adaptively composites the information from response-level (green block) and word-level (blue block) to generate word at each decoding step.

dialogue history. There are two kinds of typical approaches to this task [Das *et al.*, 2017]: *discriminative* and *generative*. Discriminative approach learns to select the best response in a candidate list, while generative approach may generate new responses that are not provided in the pre-constructed repository. The discriminative approach is relatively easier since the grammaticality and accuracy are guaranteed in the human-written responses. However, the retrieved responses are limited by the capacity of the pre-constructed repository. Even the best matched response may not be exactly appropriate since most cases are not tailored for the on-going questions [Qi *et al.*, 2020]. Therefore, the generative ability is crucial to achieve human-like conversation by synthesizing more factual and flexible responses accordingly. The typical solution for the generative visual dialogue system is based on the encoder-decoder framework [Yang *et al.*, 2019]. The encoder aims to capture the semantics of the image, question and dialogue history by embeddings, while the decoder decodes these embeddings to a response by recurrent neural net-

works (RNN) [Hopfield, 1982]. Due to the difficulty of generation, the majority of previous works [Niu *et al.*, 2019] have focused on designing more comprehensive encoder structures to make use of different aspects of information from the input. Though these methods achieve promising improvement, they still have obvious limitations, such as generating inaccurate details and repetitive words or phrases.

To tackle the above problems, we propose to adaptively incorporate more detailed information from the encoder for generating each word in the decoding process. Specifically, we propose a recurrent **Deliberation, Abandon and Memory (DAM)** module, a novel architecture of generative decoder to address the above two issues. As shown in Figure 1, on the one hand, DAM incorporates the global information in the response-level to keep semantic coherence. On the other hand, DAM pays attention to capture the related and unique details in the word-level by designing Deliberation Unit guided by the current generated word. To further reduce repetition, we devise Abandon Unit to select the unique information for the current word. In the end, Memory Unit integrates the derived word-level and response-level semantics into the memory state for word generation, which contributes to the unification of semantic coherence and the richness of details. With recurrent connections between the DAM cells inspired by LSTM [Hochreiter and Schmidhuber, 1997], the network is capable of generating visual-grounded details in a progressive manner and remarkably eliminates repetition by coverage control. Note that DAM is a universal architecture that can be combined with existing visual dialogue models by adapting the Deliberation Unit to the corresponding encoder. To show the effectiveness of DAM, we propose three models by combining DAM with three typical visual dialogue encoders, including Late Fusion encoder [Das *et al.*, 2017] for general feature fusion, Memory Network encoder [Das *et al.*, 2017] for dialogue history reasoning, and DualVD encoder [Jiang *et al.*, 2020] for visual-semantic image understanding. We show that the performance of baseline models is consistently improved by combining with DAM.

The main contributions are summarized as follows: (1) We propose a novel generative decoder DAM to generate more detailed and less repetitive responses. DAM contains a compositional structure that leverages the complementary information from both response-level and word-level, which guarantees the accuracy and richness of the responses. (2) DAM is universal to cooperate with existing visual dialogue encoders by constraining the information selection mode to adapt to different encoder structures. (3) We demonstrate the module’s capability, generality and interpretability on the VisDial v1.0 dataset. DAM consistently improves the performance of existing models and achieves a new state-of-the-art 60.93% on NDCG for the generative task.

2 Related Work

Visual Dialogue. Most previous works focused on discriminative approaches [Zheng *et al.*, 2019; Schwartz *et al.*, 2019; Kottur *et al.*, 2018; Kong and Wu, 2018] and achieved great progress. However, generative approaches, which are more practical in realistic applications, typically perform inferior

to the discriminative approaches. [Wu *et al.*, 2018] combined reinforcement learning with generative adversarial networks [Goodfellow *et al.*, 2014] to generate human-like answers. [Zhang *et al.*, 2019] introduced negative responses to generative model to reduce safe responses. [Chen *et al.*, 2020] proposed a multi-hop reasoning model to generate more accurate responses. However, how to generate less repetitive and more detailed responses has been less studied. Our work devotes to reducing the repetition and improving the richness in responses via designing a universal generative decoder by selecting more related information for generating the current word from response-level and word-level semantics.

Generation-based Dialogue Systems. The typical solution adopts the sequence-to-sequence (seq2seq) framework [Madotto *et al.*, 2018; Xu *et al.*, 2019] and uses RNN to generate responses. Existing works studied diverse aspects of generation, including expressing specific emotions [Song *et al.*, 2019; Rashkin *et al.*, 2019], introducing new topics [Xu *et al.*, 2019], generating robust task-oriented responses [Peng *et al.*, 2018; Lei *et al.*, 2018], improving the richness [Tian *et al.*, 2019] and reducing repetition [Shao *et al.*, 2017], *etc.* [See *et al.*, 2017] assigned *pointing* to copy words from the source text to improve the richness of sentences and used *coverage mechanism* to reduce repetition. The problem of reducing repetition of response has been less studied in visual dialogue. What’s more, the methods in visual dialogue cannot adopt *pointing* to copy words directly, since the *pointing* clues come from image and dialogue history in visual dialogue. One limitation of coverage mechanism is that it reduces repetition by rigid constraints of the loss function, which may result in the missing of essential words. Intuitively, understanding the input information comprehensively and capturing word-specific semantics can also reduce repetition. Inspired by this intuition, we propose a novel visual dialogue decoder to generate less repetitive and more detailed responses by considering the encoder structure and adaptively selecting and decoding information from the encoder.

3 Methodology

The visual dialogue task can be described as follows: given an image I and its caption C , a dialogue history till round $t-1$, $H_t = \{C, (Q_1, A_1), \dots, (Q_{t-1}, A_{t-1})\}$, and the current question Q_t , the task aims to generate an accurate response A_t . Our work mainly focuses on the design of a novel generative decoder architecture DAM. To prove the effectiveness of DAM, we combine it with three typical encoders: Late Fusion (LF), Memory Network (MN) and the state-of-the-art Dual-coding Visual Dialogue (DualVD). In this section, we will introduce (1) the typical encoder-decoder generative model in visual dialogue, (2) the structure of our proposed generative decoder, and (3) the combination strategies of our decoder with the three typical encoders.

3.1 Encoder-Decoder Generative Model

Our proposed DAM network is an advancement of the typical generative decoder with deliberation and control abilities. In this section, we first introduce the typical generative visual dialogue encoder-decoder model. Encoder encodes the

image I , dialogue history H_t and current question Q_t by a hidden state called knowledge vector K_t (for conciseness, t is omitted below). On each decoding step τ , the decoder, typically using a single-layer unidirectional LSTM, receives the word embedding of previous generated word $x_{\tau-1}$ and previous hidden state $s_{\tau-1}$ (the output knowledge vector K from encoder serves as the initial hidden state) and outputs a decoded vector a_τ . Then the probability distribution P_τ over the vocabulary can be computed by:

$$P_\tau = \text{softmax}(w_p^T a_\tau + b_p) \quad (1)$$

The word with the highest probability is selected as the predicted word and the model is trained by log-likelihood loss.

3.2 The DAM Decoder

DAM is a novel composite decoder that can be incorporated with standard sequence-to-sequence generation framework. It helps to improve the richness of semantic details as well as discouraging repetition in the responses. As shown in Figure 2, DAM consists of response-level semantic decode layer (RSL), word-level detail decode layer (WDL) and information fusion module (Memory Unit). RSL is responsible for capturing the global information to guarantee the response’s fluency and correctness. However, the global information lacks the detailed semantics, for the current word and the rigid-decoding mode in LSTM tends to generate repeated words. WDL incorporates the essential and unique visual dialogue contents (i.e. question, dialogue history and image) into the generation of current word to enrich the word-level details. The structure of WDL consists of an LSTM, Deliberation Unit and Abandon Unit. Finally, Memory Unit is responsible for adaptively fusing both the response-level and word-level information.

Response-Level Semantic Decode Layer (RSL)

When answering a question about an image, human needs to capture the global semantic information to decide the main ideas and content for the responses. In our model, we regard the embedded information from the encoder as global semantic information, and denote it as knowledge vector K . K is used for providing the response-level semantics in the generation process. The response-level information r_τ for generating the current word is computed as:

$$r_\tau = LSTM_r(x_{\tau-1}, s_{\tau-1}^r) \quad (2)$$

where $x_{\tau-1}$ is the previous generated word and $s_{\tau-1}^r$ is the memory state of $LSTM_r$.

Word-Level Detail Decode Layer (WDL)

On the one hand, the response-level information lacks the details of the image and dialogue history, providing rigid clues for generating different words. On the other hand, response-level information changes slightly with the recurrent word generation process and results in repetitive words or phrases. To solve these problems, it’s critical to enrich the decoding vector with more detailed question-relevant information that is unique for current generated word.

For generating the τ^{th} word, we first adaptively capture word-relevant information from the encoded knowledge information along with previous generated word and previous

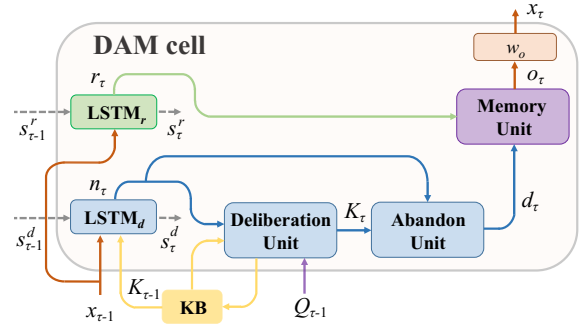


Figure 2: Overview structure of DAM. It consists of RSL (green part), WDL (blue part) and information fusion module (purple part). For the τ generation step, the inputs of DAM contain the question embedding $Q_{\tau-1}$, the knowledge embedding $K_{\tau-1}$, the previous generated word embedding $x_{\tau-1}$ and the previous hidden states $s_{\tau-1}^r$ and $s_{\tau-1}^d$.

hidden state via $LSTM_d$:

$$n_\tau = LSTM_d([x_{\tau-1}, K_{\tau-1}], s_{\tau-1}^d) \quad (3)$$

where “[,]” denotes concatenation, $K_{\tau-1}$ is the updated knowledge vector in the $\tau-1$ step and $s_{\tau-1}^d$ is the memory state of $LSTM_d$. Since n_τ only capture the global semantics from the encoder, we further incorporate the structure-adaptive local semantics from the encoder via the Deliberation Unit. Finally, we propose the Abandon Unit to filter out the redundant information while enhancing the word-specific information from both global and local clues. The Deliberation Unit and the Abandon Unit are detailed below.

Deliberation Unit. It aims to adaptively leverage the encoder structure to extract the most related and detailed information for current word generation. Specifically, we first capture the significant information in the question under the guidance of the global semantic vector n_τ . Guided by the upgraded question representation, we adopt structure-adaptive strategies to different encoder structures to select image and dialogue history information. In the end, we get the detailed question-related information by fusing the information of question, dialogue history and image. Compared with most existing decoders that merely use the encoded embedding without considering the diverse encoder structures, our proposed Deliberation Unit provides a flexible strategy to derive more detailed information by taking the advantages of the elaborate encoders. To prove the effectiveness of DAM, we combine it with three typical encoders, including LF encoder for the general feature fusion, MN encoder for dialogue history reasoning and DualVD encoder for visual-semantic image understanding. The details of Deliberation Unit adaptive to these three encoders will be introduced in Section 3.3.

Abandon Unit. It further filters out the redundant information while enhancing the word-specific information from both the global and local encoded clues. Specifically, Abandon Unit updates current generated decoding vector n_τ by combining detailed knowledge information K_τ with n_τ via a gate operation and achieves the final word-level embedding d_τ :

$$gate_\tau^a = \sigma(\mathbf{W}_a[n_\tau, K_\tau] + b_a) \quad (4)$$

$$d_\tau = gate_\tau^a \circ [n_\tau, K_\tau] \quad (5)$$

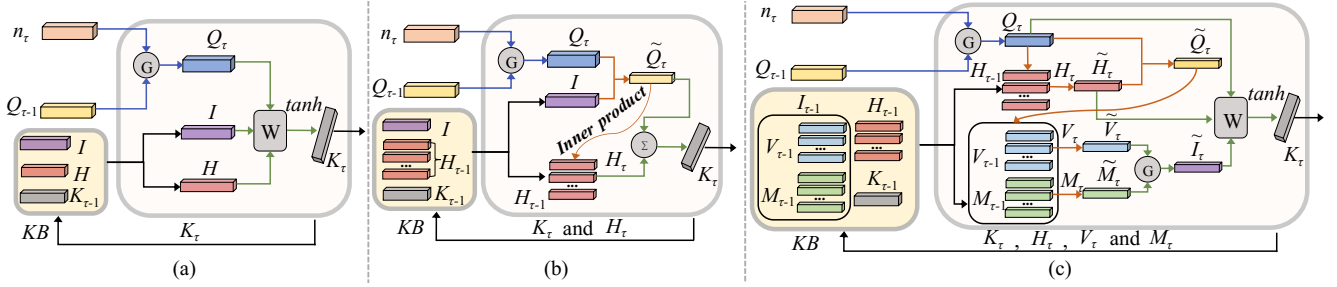


Figure 3: The illustration of Deliberation Unit adaptive to LF Encoder (a), MN Encoder (b) and DualVD Encoder (c), where ‘‘G’’: gate operation, V : visual representation of image, M semantic representation of image. Different colored lines represent different update steps: *blue lines*: word-guided question information update, *orange lines*: question-guided information update, *green lines*: general feature fusion.

where ‘‘ \circ ’’ denotes the element-wise product.

Two Level Information Fusion

The information from RSL and WDL is complementary to each other. We design *Memory Unit* to combine the two kinds of information for word prediction. Memory Unit selects response-level information to control the global semantics in response and tracks the word-level information for generating more detailed and less repeated response via a gate operation:

$$gate_{\tau}^m = \sigma(\mathbf{W}_m[r_{\tau}, d_{\tau}] + b_m) \quad (6)$$

$$o_{\tau} = gate_{\tau}^m \circ [r_{\tau}, d_{\tau}] \quad (7)$$

The generated word x_{τ} with the maximum value in the probability distribution P_{τ}^o is selected as the predicted word. P_{τ}^o is computed as:

$$P_{\tau}^o = softmax(\mathbf{w}_o^T o_{\tau} + b_o) \quad (8)$$

3.3 Variants of Deliberation Unit

Guided by the question and current generated word state n_{τ} , Deliberation Unit captures more detailed information from encoder-specific structures. The Deliberation Unit mainly contains three steps: (1) word-guided question information update, (2) question-guided information update, and (3) general feature fusion. The last two steps are adaptive to different encoders while the first step keeps unchanged. To select the most related information for current generated word, we first update question information $Q_{\tau-1}$ with n_{τ} :

$$gate_{\tau}^q = \sigma(\mathbf{W}_q[Q_{\tau-1}, n_{\tau}] + b_q) \quad (9)$$

$$Q_{\tau} = \mathbf{W}_1(gate_{\tau}^q \circ [Q_{\tau-1}, n_{\tau}]) \quad (10)$$

We will introduce the next two steps adaptive to LF, MN and DualVD encoders below. It should be noted that the parameters of the Deliberation Unit are independent of its encoder.

Deliberation Unit Adaptive to LF Encoder. LF Encoder focuses on multi-modal information fusion without complex information reasoning. In our decoder, we merely fuse the updated question information Q_{τ} with dialogue history H and image I from the encoder without question-guided information update step as shown in Figure 3(a).

Deliberation Unit Adaptive to MN Encoder. MN Encoder focuses on the dialogue history reasoning. Compared with Deliberation Unit for LF Encoder, we further add question-guided information update step to reason over dialogue history via attention mechanism before general feature fusion as shown in Figure 3(b).

Model	MRR	R@1	R@5	R@10	Mean	NDCG
HCIAE-G [Lu <i>et al.</i> , 2017]	49.07	39.72	58.23	64.73	18.43	59.70
CoAtt-G [Wu <i>et al.</i> , 2018]	49.64	40.09	59.37	65.92	17.86	59.24
Primary-G [Guo <i>et al.</i> , 2019]	49.01	38.54	59.82	66.94	16.69	-
ReDAN-G [Gan <i>et al.</i> , 2019]	49.60	39.95	59.32	65.97	17.79	59.41
DMRM [Chen <i>et al.</i> , 2020]	50.16	40.15	60.02	67.21	15.19	-
LF-G [Das <i>et al.</i> , 2017]	44.67	34.84	53.64	59.69	21.11	52.23
MN-G [Das <i>et al.</i> , 2017]	45.51	35.40	54.91	61.20	20.24	51.86
DualVD-G [Jiang <i>et al.</i> , 2020]	49.78	39.96	59.96	66.62	17.49	60.08
LF-DAM (ours)	45.08	35.01	54.48	60.57	20.83	52.68
MN-DAM (ours)	46.16	35.87	55.99	62.45	19.57	52.82
DualVD-DAM (ours)	50.51	40.53	60.84	67.94	16.65	60.93

Table 1: Result comparison on validation set of VisDial v1.0.

Deliberation Unit Adaptive to DualVD Encoder. DualVD Encoder focuses on the visual-semantic image understanding. As shown in Figure 3(c), for the question-guided information update step, we first concatenate updated question and dialogue history to form the query vector \tilde{Q}_{τ} , and assign \tilde{Q}_{τ} to guide the update of image from visual and semantic aspects respectively. For the feature fusion step, we utilize the gate operation between visual and semantic representation (\tilde{V}_{τ} and \tilde{T}_{τ}) to obtain the updated image representation.

4 Experiments

Dataset. We conduct extensive experiments on VisDial v1.0 dataset [Das *et al.*, 2017] constructed based on MSCOCO images and captions. VisDial v1.0 is split into training, validation and test sets. The training set consists of dialogues on 120k images from COCO-trainval while the validation and test sets are consisting of dialogues on an additional 10k COCO-like images from Flickr.

Evaluation Metrics. Following [Das *et al.*, 2017], we rank the 100 candidate answers based on their posterior probabilities and evaluate the performance by retrieval metrics: mean reciprocal rank (MRR), recall@ k ($k=1, 5, 10$), mean rank of human response (Mean) and normalized discounted cumulative gain (NDCG). Lower value for Mean and higher value for other metrics are desired.

Implementation Details. To build the vocabulary, we retain words in the dataset with word frequency greater than 5. The vocabulary contains 10366 words. The hidden states and cell states of $LSTM_d$ are randomly initialized while $LSTM_r$ is using the output knowledge vector K from encoder as the initial hidden state and randomly initializing cell state. The maximum sentence length of the responses is set to 20. The

Base Model	Model	MRR	R@1	R@5	R@10	Mean	NDCG
LF-DAM	2LSTM	44.43	34.53	53.55	59.48	21.38	51.99
	2L-M	44.77	34.85	54.06	60.03	21.13	52.04
	2L-DM	45.06	34.90	54.24	60.39	20.87	52.58
	2L-DAM	45.08	35.01	54.48	60.57	20.83	52.68
MN-DAM	2LSTM	45.58	35.27	55.38	61.54	19.96	52.38
	2L-M	45.67	35.29	55.57	61.97	19.91	52.11
	2L-DM	45.77	35.53	55.40	62.05	19.95	52.51
	2L-DAM	46.16	35.87	55.99	62.45	19.57	52.82
DualVD-DAM	2LSTM	49.72	40.04	59.52	66.41	17.62	59.79
	2L-M	50.09	40.38	59.94	66.77	17.31	59.85
	2L-DM	50.20	40.33	60.22	67.48	17.15	59.72
	2L-DAM	50.51	40.53	60.84	67.94	16.65	60.93

Table 2: Ablation study of each unit on VisDial v1.0 validation set.

hidden state size of all the LSTM blocks is set to 512 and the dimension of each gate is set to 1024. The Adam optimizer [Kingma and Ba, 2015] is used with the initial learning rate of 1e-3 and final learning rate of 3.4e-4 via cosine annealing strategy with 16 epochs. The batch size is set to 15.

4.1 State-of-the-Art Comparison

As shown in Table 1, we compare our models (third block) with SOTA generative models (first block) and baseline models (second block, re-trained by us). ReDAN-G and DMRM adopted complex multi-step reasoning, while HCIAE-G, CoAtt-G and Primary-G are attention-based models. For fairness, we only compare the original generative ability without re-ranking. We just replace the decoders in baseline models by our proposed DAM. Compared with the baseline models, our models outperform them on all the metrics, which indicates the complementary advantages between DAM and existing encoders in visual dialogue. Though DualVD-G performs lower than DMRM on *Mean*, DualVD-DAM outperforms DMRM on all the other metrics without multi-step reasoning, which is the advantages in DMRM over our models.

4.2 Ablation Study

The Effectiveness of Each Unit

We consider the following ablation models to illustrate the effectiveness of each unit of our model: 1) **2L-DAM**: this is our full model that adaptively selects related information for decoding. 2) **2L-DM**: full model w/o Abandon Unit. 3) **2L-M**: 2L-DM w/o Deliberation Unit. 4) **2-LSTM**: 2L-M w/o Memory Unit. As shown in Table 2, taking DualVD-DAM for example, the MRR values increase by 0.37%, 0.11% and 0.31% respectively when introducing the Memory Unit (2L-M), Deliberation Unit (2L-DM) and Abandon Unit (2L-DAM) to the baseline model (2-LSTM) progressively. Similar trend exists in LF-DAM and MN-DAM, which indicates the effectiveness of each unit in DAM. Since the space limitation and similar observations, we show the ablation studies on DualVD-DAM in the following experiments.

The Effectiveness of Two-Level Decode Structure

To prove the complementary advantages of the response-level semantic decode layer (RSL) and the word-level detail decode layer (WDL), and to figure out the information selection mode, we first conduct **Human Study** on the effectiveness of RSL, WDL and the full model DualVD-DAM, and then visualize the gate values of Memory Unit to reveal the information selection mode.

Model	M1↑	M2↑	Repetition↓	Richness↑
RSL(DualVD-G): RSL only	0.60	0.47	0.20	0.03
WDL: WDL only	0.69	0.54	0.07	0.15
DualVD-DAM	0.75	0.61	0.01	0.13

Table 3: Human evaluation of 100 sample responses on VisDial v1.0 validation set. M1: percentage of responses that pass the Turing Test. M2: percentage of responses that are evaluated better or equal to human responses. Repetition: percentage of responses that have meaningless repeated words. Richness: percentage of responses that contain detailed content to answer the question.

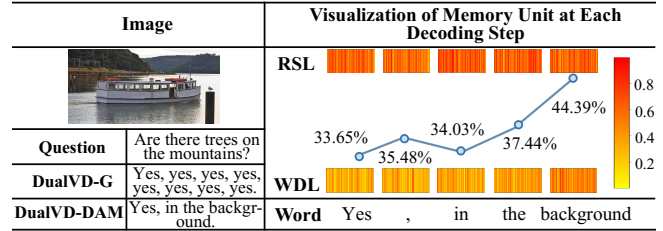


Figure 4: Visualization of gate values in the Memory Unit. Yellow thermogram: gate values of RSL and WDL computed in Eq. (6). Blue polyline: ratio of total gate values of WDL.

Complementary Advantages. In the human study, we follow [Wu *et al.*, 2018] to sample 100 results from VisDial v1.0 validation set and ask 3 persons to evaluate the quality of the last response in the dialogue. Distinct from previous works, we add *Repetition* and *Richness* metrics, and for all metrics, we record the score when at least 2 persons agree. As shown in Table 3, WDL performs best on *Richness* and reduces the *Repetition* by 0.13 compared to RSL, which indicates that WDL contributes to the increase of detailed information and the decrease of repetition in the response. After incorporating RSL and Memory Unit with WDL, the *repetition* further reduces by 0.06 while M1 and M2 outperform by 0.06 and 0.07 respectively, which proves the complementary advantages between these two level information. We also notice that *Richness* decreases slightly. This is mainly because the information from RSL concentrates more attention on the global information, rather than detailed information.

Information Selection Mode. We visualize the gate values in the Memory Unit for DualVD-DAM to demonstrate the information selection mode of the two level information. As shown in Figure 4, we can observe that the ratio of gate values for RSL is always higher than that for WDL. It indicates that the response-level information in RSL plays the predominant role in guiding the response generation. Another obvious phenomenon is that the ratio of gate values for WDL increases rapidly when generating the last word, which can be viewed as a signal to stop the response generation in time when the response already covers the complete semantics. It may due to the fact that WDL captures word-level information and is sensitive to the word repetition, which is beneficial to avoid repetitive word generation.

The Effectiveness of Each Operation in Deliberation Unit

We conduct experiments on DualVD-DAM to reveal the influence of essential operations in the Deliberation Unit: 1) **I-S** only uses semantic-level image information for information selection. 2) **I-V** only utilizes visual-level image information for information selection. 3) **I-SV** jointly exploits semantic



Image + Caption	Question	Human Responses	DualVD-G	DualVD-DAM
 C: A cheerful older gentleman with his hands on his hips is smiling at the camera.	Q1: How old is the gentleman?	A1: About 60.	A1: Maybe in his face is in his face, ..., is in his face.	A1: He 's in his 60s.
	Q2: Is he happy or sad?	A2: Looks very happy.	A2: He looks like he looks like he looks, ..., like he looks.	A2: He looks happy.
	Q3: Is he giving attitude or having fun?	A3: Having fun.	A3: He looks like he looks like he looks like, ..., he looks.	A3: Having fun and he looks like he is posing for the picture.
 C: A woman playing tennis with the ball in mid-air in front of her and the racket raised with her arm across her body in front of a sponsored grandstand.	Q1: Do you see a lot of people?	A1: Just 3.	A1: I see 3.	A1: I see 3 people.
	Q2: What is the tennis player wearing?	A2: White tennis dress.	A2: White shirt and a white shirt and a white shirt, ..., and a.	A2: A white tennis dress and tennis shoes.
	Q3: What color is her tennis racket?	A3: Black.	A3: Red and red and red and red and red and red, ..., and red and.	A3: Black and white.

Figure 5: Qualitative results of DualVD-DAM comparing to Human Responses and DualVD-G, where “...” are omitted repeated words.





Image	Dialogue History	Visualization		
	C: A tarmac with a lot of large blue and white planes parked. Q1: Are there people? A1: I see 2 people. Q2: Is it sunny? A2: It looks like a clear day, yes. Q3: Are there clouds? A3: A couple of clouds, yes. Q4: Are the planes big? A4: They look like large passenger planes. Q5: Are there people boarding? A5: No. Q6: Are there any bags? A6: No. Q7: Are there signs? A7: No.			
	Question Is there a building?			
	DualVD-G No.			
	DualVD-DAM There are some buildings in the background.			

Figure 6: Visualization of the evidence when generating the response by DualVD-DAM. The essential visual regions and dialogue history for answering the question are highlighted in the last three columns. The attention weights of visual regions and dialogue history are visualized, where clearer region and darker orange color indicates higher attention weight.

Model	MRR	R@1	R@5	R@10	Mean	NDCG
I-S	50.01	40.25	59.78	66.76	17.67	59.09
I-V	50.03	40.30	59.34	66.90	17.34	58.93
I-SV	50.13	40.34	60.09	67.06	17.34	59.51
H	50.19	40.36	60.09	66.96	17.27	59.92
DualVD-DAM	50.51	40.53	60.84	67.94	16.65	60.93

Table 4: Ablation study of Deliberation Unit on VisDial v1.0.

and visual information for information selection. 4) **H** only leverages dialogue history for information selection.

As shown in Table 4, I-S and I-V update information from visual and semantic aspect respectively, while I-SV updates information from both two aspects which achieves the best performance compared to the above two models. The relatively higher results of H model indicate that the history information plays a more important role in the decoder. By jointly incorporating all the structure-aware information from the encoder, DualVD-DAM achieves the best performance on all the metrics. It proves the advantages of DAM via fully utilizing the information from the elaborate encoder, which is beneficial for enhancing the existing generation models by incorporating their encoders with DAM adaptively.

4.3 Qualitative Analysis

Response generation quality. Figure 5 shows two examples of three-round dialogues and the corresponding responses generated by DualVD-G and DualVD-DAM. When answering Q3 in the first example, DualVD-DAM generates accurate and non-repetitive response “*having fun*” compared with DualVD-G. Comparing to the human response, DualVD-DAM further provides detailed description “*he is posing for the picture*” so as to increase the richness of the response. Similar observation exists in the second example.

Information selection quality. We further visualize the evidence captured from the image and dialogue history for generating the essential words, i.e. *there*, *buildings* and *background*. As shown in Figure 6, it is difficult to answer the question of “*Is there a building?*” accurately, since the *buildings* are distant and small. DualVD-DAM accurately focuses on the visual and dialogue clues. Taking the word *background* for example, our model focuses on the background in the image and highlights the *clear day* in the dialogue history. It proves that DAM can adaptively focus on the exact visual and textual clues for generating each word, which contributes to the high quality of the responses.

5 Conclusion

In this paper, we propose a novel generative decoder DAM consisting of the *Deliberation* Unit, *Abandon* Unit and *Memory* Unit. The novel decoder adopts a compositive decoding mode in order to model information from both response-level and word-level, so as to discourage repetition in the generated responses. DAM is a universal decoding architecture which can be incorporated with existing visual dialogue encoders to improve their performance. The extensive experiments of combining DAM with LF, MN and DualVD encoders verify that our proposed DAM can effectively improve the generation performance of existing models and achieve new state-of-the-art results on the popular benchmark dataset.

Acknowledgements

This work is supported by the National Key Research and Development Program (Grant No.2017YFB0803301).

References

- [Chen *et al.*, 2020] Feilong Chen, Fandong Meng, Jiaming Xu, Peng Li, Bo Xu, and Jie Zhou. Dmrm: A dual-channel multi-hop reasoning model for visual dialog. In *AAAI*, 2020.
- [Das *et al.*, 2017] Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José M. F. Moura, Devi Parikh, and Dhruv Batra. Visual dialog. In *CVPR*, pages 1080–1089, 2017.
- [Gan *et al.*, 2019] Zhe Gan, Yu Cheng, Ahmed Ei Kholy, Linjie Li, and Jianfeng Gao. Multi-step reasoning via recurrent dual attention for visual dialog. In *ACL*, page 6463–6474, 2019.
- [Goodfellow *et al.*, 2014] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. In *NIPS*, pages 2672–2680, 2014.
- [Guo *et al.*, 2019] Dalu Guo, Chang Xu, and Dacheng Tao. Image-question-answer synergistic network for visual dialog. In *CVPR*, pages 10434–10443, 2019.
- [Hochreiter and Schmidhuber, 1997] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [Hopfield, 1982] John J. Hopfield. Neural networks and physical systems with emergent collective computational abilities. *Proceedings of the National Academy of Sciences of the United States of America*, 79(8):2554–2558, 1982.
- [Jiang *et al.*, 2020] Xiaoze Jiang, Jing Yu, Zengchang Qin, Yingying Zhuang, Xingxing Zhang, Yue Hu, and Qi Wu. Dualvd: An adaptive dual encoding model for deep visual understanding in visual dialogue. In *AAAI*, 2020.
- [Kingma and Ba, 2015] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *ICLR*, 2015.
- [Kong and Wu, 2018] Dejiang Kong and Fei Wu. Visual dialog with multi-turn attentional memory network. In *PCM*, pages 611–621, 2018.
- [Kottur *et al.*, 2018] Satwik Kottur, José MF Moura, Devi Parikh, Dhruv Batra, and Marcus Rohrbach. Visual coreference resolution in visual dialog using neural module networks. In *ECCV*, pages 153–169, 2018.
- [Lei *et al.*, 2018] Wenqiang Lei, Xisen Jin, Min-Yen Kan, Zhaochun Ren, Xiangnan He, and Dawei Yin. Sequicity: Simplifying task-oriented dialogue systems with single sequence-to-sequence architectures. In *ACL*, pages 1437–1447, 2018.
- [Lu *et al.*, 2017] Jiasen Lu, Anitha Kannan, Jianwei Yang, Devi Parikh, and Dhruv Batra. Best of both worlds: Transferring knowledge from discriminative learning to a generative visual dialog model. In *NIPS*, pages 314–324, 2017.
- [Madotto *et al.*, 2018] Andrea Madotto, Chien-Sheng Wu, and Pascale Fung. Mem2Seq: Effectively incorporating knowledge bases into end-to-end task-oriented dialog systems. In *ACL*, pages 1468–1478, 2018.
- [Niu *et al.*, 2019] Yulei Niu, Hanwang Zhang, Manli Zhang, Jianhong Zhang, Zhiwu Lu, and Ji-Rong Wen. Recursive visual attention in visual dialog. In *CVPR*, pages 6679–6688, 2019.
- [Peng *et al.*, 2018] Baolin Peng, Xiujun Li, Jianfeng Gao, Jingjing Liu, and Kam-Fai Wong. Deep Dyna-Q: Integrating planning for task-completion dialogue policy learning. In *ACL*, pages 2182–2192, 2018.
- [Qi *et al.*, 2020] Jiaxin Qi, Yulei Niu, Jianqiang Huang, and Hanwang Zhang. Two causal principles for improving visual dialog. In *CVPR*, 2020.
- [Rashkin *et al.*, 2019] Hannah Rashkin, Eric Michael Smith, Margaret Li, and Y-Lan Boureau. Towards empathetic open-domain conversation models: A new benchmark and dataset. In *ACL*, pages 5370–5381, 2019.
- [Schwartz *et al.*, 2019] Idan Schwartz, Seunghak Yu, Tamir Hazan, and Alexander G Schwing. Factor graph attention. In *CVPR*, pages 2039–2048, 2019.
- [See *et al.*, 2017] Abigail See, Peter J Liu, and Christopher D Manning. Get to the point: Summarization with pointer-generator networks. In *ACL*, page 1073–1083, 2017.
- [Shao *et al.*, 2017] Louis Shao, Stephan Gouws, Denny Britz, Anna Goldie, Brian Strope, and Ray Kurzweil. Generating high-quality and informative conversation responses with sequence-to-sequence models. In *EMNLP*, page 2210–2219, 2017.
- [Song *et al.*, 2019] Zhenqiao Song, Xiaoqing Zheng, Lu Liu, Mu Xu, and Xuanjing Huang. Generating responses with a specific emotion in dialog. In *ACL*, pages 3685–3695, 2019.
- [Tian *et al.*, 2019] Zhiliang Tian, Wei Bi, Xiaopeng Li, and Nevin L. Zhang. Learning to abstract for memory-augmented conversational response generation. In *ACL*, pages 3816–3825, 2019.
- [Wu *et al.*, 2018] Qi Wu, Peng Wang, Chunhua Shen, Ian Reid, and Anton van den Hengel. Are you talking to me? reasoned visual dialog generation through adversarial learning. In *CVPR*, pages 6106–6115, 2018.
- [Xu *et al.*, 2019] Can Xu, Wei Wu, Chongyang Tao, Huang Hu, Matt Schuerman, and Ying Wang. Neural response generation with meta-words. In *ACL*, pages 5416–5426, 2019.
- [Yang *et al.*, 2019] Tianhao Yang, Zheng-Jun Zha, and Hanwang Zhang. Making history matter: Gold-critic sequence training for visual dialog. In *ICCV*, pages 2561–2569, 2019.
- [Zhang *et al.*, 2019] Heming Zhang, Shalini Ghosh, Larry Heck, Stephen Walsh, Junting Zhang, Jie Zhang, and C-C Jay Kuo. Generative visual dialogue system via adaptive reasoning and weighted likelihood estimation. In *IJCAI*, pages 1025–1031, 2019.
- [Zheng *et al.*, 2019] Zilong Zheng, Wenguan Wang, Siyuan Qi, and Song-Chun Zhu. Reasoning visual dialogs with structural and partial observations. In *CVPR*, pages 6669–6678, 2019.