

# Hybrid Reasoning Network for Video-based Commonsense Captioning

Weijiang Yu<sup>1</sup>, Jian Liang<sup>2</sup>, Lei Ji<sup>3</sup>, Lu Li<sup>4</sup>, Yuejian Fang<sup>2</sup>†, Nong Xiao<sup>1</sup>†, Nan Duan<sup>3</sup>

<sup>1</sup>Sun Yat-sen University, <sup>2</sup>Peking University, <sup>3</sup>Microsoft Research Asia, <sup>4</sup>Zhejiang University  
weijiangyu8@gmail.com, j.liang@stu.pku.edu.cn, leiji@microsoft.com, lu.lee@zju.edu.cn  
fangyj@ss.pku.edu.cn, xiaon6@mail.sysu.edu.cn, nanduan@microsoft.com

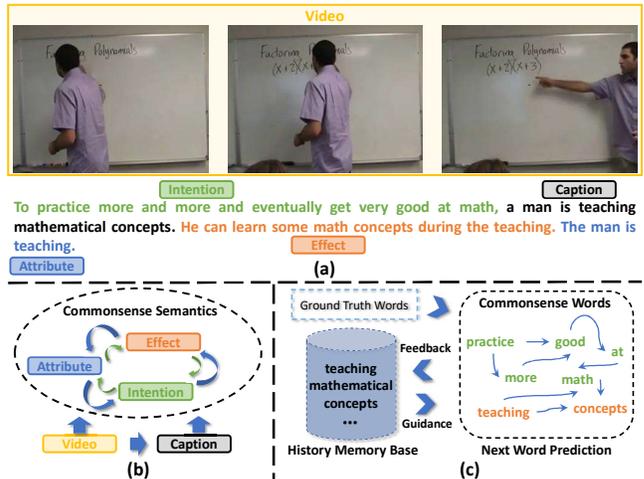
## ABSTRACT

The task of video-based commonsense captioning aims to generate event-wise captions and meanwhile provide multiple commonsense descriptions (e.g., attribute, effect and intention) about the underlying event in the video. Prior works explore the commonsense captions by using separate networks for different commonsense types, which is time-consuming and lacks mining the interaction of different commonsense. In this paper, we propose a Hybrid Reasoning Network (HybridNet) to endow the neural networks with the capability of semantic-level reasoning and word-level reasoning. Firstly, we develop multi-commonsense learning for semantic-level reasoning by jointly training different commonsense types in a unified network, which encourages the interaction between the clues of multiple commonsense descriptions, event-wise captions and videos. Then, there are two steps to achieve the word-level reasoning: (1) a memory module records the history predicted sequence from the previous generation processes; (2) a memory-routed multi-head attention (MMHA) module updates the word-level attention maps by incorporating the history information from the memory module into the transformer decoder for word-level reasoning. Moreover, the multimodal features are used to make full use of diverse knowledge for commonsense reasoning. Experiments and abundant analysis on the large-scale Video-to-Commonsense benchmark show that our HybridNet achieves state-of-the-art performance compared with other methods.

## 1 INTRODUCTION

Recently, research on video-based commonsense captioning [11] has been gaining increasing attention, as it provides a deeper understanding of the video and language and thus facilitates various visual reasoning tasks ranging from fundamental scene understanding [8, 19, 42] to high-level visual-linguistic reasoning tasks [13, 18, 44, 46]. The video-based commonsense captioning task aims to generate captions and three types of commonsense descriptions (intention, effect, attribute) simultaneously given the input video. A show case is presented in Figure 1 (a).

The video-based commonsense captioning is a frontier research topic and the current best performer [11] executes separate networks to learn different types of commonsense separately. It is time-consuming and counterintuitive. When humans infer a specific event, each commonsense is not identified individually. We often consider a global perspective to reason the commonsense



**Figure 1: (a) The example of the video-based commonsense captioning task. It requires the network to not only generate event-wise captions but also provide visually grounded commonsense descriptions about the underlying event in the video. The task seeks to describe the intentions of the agent, effects of the action, and the attributes of the agent’s characteristic. (b) Semantic-level reasoning to interplay with different commonsense clues. (c) Word-level reasoning generates the next word, which is guided by the history sequences as well as using the ground truth words during training.**

semantic coherence by interacting with the clues from various commonsense types. For example, humans can correctly reason about the intention in Figure 1 (b) not only benefit from the video input and event-wise caption but also with the help of accurate predictions of attribute and effect. Such high-level semantic interaction can facilitate multiple commonsense reasoning.

Besides, humans can remain a logic loop when performing commonsense inference. Taking the intention sentence as an example in Figure 1 (a), when humans see the history token sequence like “to practice more and more and eventually get very”, we can easily infer the next word should be positive (e.g., “good”). The “good” can also be used to reason the accuracy of previously predicted words and then correct them for contextual logic consistency (i.e. logic loop). As we can see, such word-level reasoning in the description can benefit from each predicted word. However, current methods [11, 34, 45] follow the Transformer [35] training mechanism to predict the next word only with the ground truth words as context while generating the entire sequence from scratch at

† means corresponding author.

This work was done during the first author’s internship in MSR Asia.

This is an arXiv preprint, to be presented at ACM Multimedia 2021.

The code is publicly available at <https://github.com/yuweijiang/HybridNet>.

inference, which leads to a gap between the previously predicted words and the next word.

To address the above issues, we achieve the semantic-level reasoning in Figure 1 (b) to jointly learn the multiple commonsense types. Then we perform the word-level reasoning in Figure 1 (c) in the network for bridging the gap between the previously predicted words and the next word. In this paper, we propose a novel Hybrid Reasoning Network (HybridNet) in Figure 2 (a) to subtly integrate these two levels of reasoning into a unified framework.

For semantic-level reasoning, we propose multi-commonsense learning as shown in Figure 2 (a) to achieve semantic-level reasoning, in which multiple learning commonsense descriptions are solved at the same time while exploiting commonalities and differences across commonsense semantics. Multi-commonsense learning is a mechanism to inductive transfer that improves generalization by using the domain information contained in the training signals of related commonsense as an inductive bias. It achieves this by reasoning various types of commonsense descriptions in parallel while using a shared representation. What is learned for each commonsense semantic can help other commonsense semantics be learned better. Furthermore, we merge 3D motion features, 1D audio features and 2D appearance features together via the multimodal fusion for enforcing our model to learn to reason commonsense semantics from diverse information.

To utilize the previous context words for predicting the next word, we propose a memory module to record the history information from previous generation processes. Then we introduce a memory-routed multi-head attention (MMHA) to incorporate the history information with the ground truth words into the network for next word prediction during the training. Specifically, our MMHA infers the next word conditioned on history information and ground truth words by learning a merging attention map. In Figure 3 (b), the MMHA first predicts the conditional attention map routed by the previous states from the memory module. Then we design a triangle convolution to learn the conditional attention map and meanwhile prevent foreseeing subsequent positions. Next, we merge the conditional attention map with several attention maps, including the original attention map from the multi-head attention and the previous attention maps of the front blocks. The previous attention maps of different blocks are bridged via our proposed contextual residual connections to learn the long-range attention context. Finally, our MMHA updates the word feature (i.e.  $V$ ) via a merged attention map to achieve word-level reasoning. In this paper, we utilize these two modules to achieve word-level reasoning. Their architectures can be seen in Figure 2 (b) and Figure 3 (b).

**Contributions.** (1) We propose a Hybrid Reasoning Network (HybridNet) to jointly perform semantic-level reasoning and word-level reasoning. (2) A multi-commonsense learning is proposed to achieve the semantic-level reasoning based on the video and caption inputs. (3) A memory-routed multi-head attention is introduced to cooperate with a novel memory module to execute the word-level reasoning. (4) Extensive experiments and abundant analysis have demonstrated the effectiveness and superiority of our proposed modules.

## 2 RELATED WORK

**Visual Comprehension.** More and more researchers pay more attention to the visual comprehension community by targeting visual question answering [1], visual dialogue [9], visual question generation [20], image captioning [38], visual commonsense reasoning [44], visual grounding [27], situation recognition [42] and video captioning [23]. Recently, the visual reasoning tends to cognition-level reasoning by incorporating the commonsense concepts [4, 5, 16, 31–33] from the natural language processing (NLP). For example, early works [26, 39] utilized the prior commonsense knowledge to assist in the prediction of action motivation. Zellers et al. [44] proposed a visual commonsense reasoning task to not only provide question answering but also predict the correct rationale behind the answer based on the question and image. Based on the cause-effect clues, Wang et al. [40] presented a visual commonsense R-CNN on object detection by mining the commonsense knowledge behind the object categories. Yu et al. [43] proposed a heterogeneous graph learning method to seamlessly integrate the intra-graph and inter-graph reasoning in order to bridge multi-modal domains for visual commonsense reasoning. For video-based commonsense captioning, Fang et al. [11] used individual transformers for different commonsense captioning, which lacks of the commonsense interaction. In this work, we propose a hybrid reasoning network to jointly learn multiple commonsense descriptions via semantic-level reasoning and word-level reasoning.

**Video Captioning.** Visual perception and language expression are two key capabilities of human intelligence, and video captioning is an insight example towards learning from humans to bridge vision and language. There are mainly two aspects for video captioning development: datasets and methods. For developing the video captioning, some early works proposed some specific-domain datasets like movie [28, 30] and cooking [10, 29], which is limited and small for deep learning. Some researchers tended to open-domain video captioning datasets such as MSVD [7], MSR-VTT [41] and TGIF [21]. Recently, video captioning tries to connect with the commonsense to explore the commonsense descriptions in the video, which comes up with a dataset named Video-to-Commonsense [11]. In the method sight, current deep-learning-based video captioning often performs sequence-to-sequence learning in an encoder-decoder paradigm. In between, an encoder equipped with powerful deep neural networks is exploited to learn video representation. A decoder of sentence generation is utilized to translate the learned representation into a sentence with more flexible structures. Venugopalan et al. [37] applied the sequence-to-sequence model into the video captioning by end-to-end learning way. To bridge the sentence semantics and visual content, an attention-based LSTM named aLSTMs [12] is proposed to better transfer videos to natural sentences by capturing salient structures of video. For dense video captioning, Zhou et al. [45] used an end-to-end transformer [35] architecture to jointly learning the video encoder, proposal decoder, and captioning decoder. Fang et al. [11] used an encoder-decoder way to individually model the specific commonsense captioning without using the commonsense correlations. In this paper, we propose to generate the commonsense descriptions in the video from semantic-level and word-level inference.

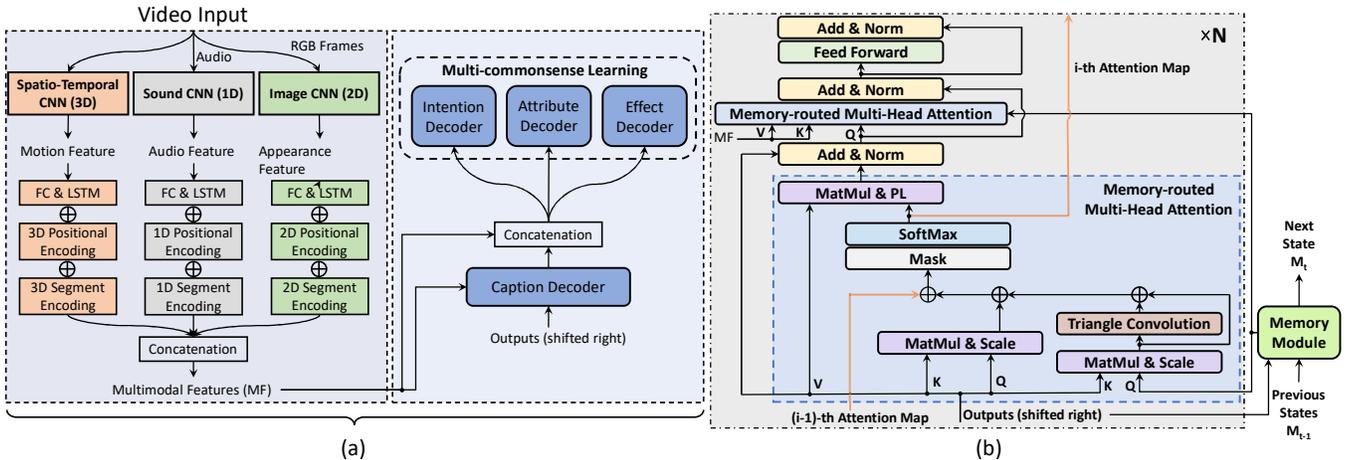


Figure 2: (a) Overview of our HybridNet. Taking the video as input, a multimodal fusion is shown to merge the motion feature, audio feature and appearance feature as multimodal features (MF). Then the MF is fed into the decoder stage for captioning, including the caption decoder and three commonsense decoders. We utilize multi-commonsense learning on three commonsense decoders. (b) The pipeline of the proposed memory-routed multi-head attention (MMHA) in each decoder block. The MMHA is guided by the memory module that is illustrated in Figure 3 (b). The “MatMul & Scale” indicates the scaled dot-product operation and the “PL” means the projection layer. The orange lines with arrows denote the contextual residual connections to bridge contextual attention maps from different blocks. The  $\oplus$  means the addition operation in this work.

### 3 HYBRID REASONING NETWORK

In this section, we explain the architecture and design of our proposed Hybrid Reasoning Network (HybridNet), which can appropriately interact with multiple commonsense semantics and preserve the connection between contextual predicted sequences. As shown in Figure 2 (a), our HybridNet is an encoder-decoder architecture, including a video encoder, a caption decoder and three commonsense decoders. Given a video input, there are three pre-trained models to extract multiple features including motion feature, audio feature and appearance feature. Then a multimodal fusion is utilized to merge the extracted features as multimodal features. The multimodal features are fed into the caption decoder to obtain caption encoding as well as predict event-wise captions. Then we concatenate the multimodal features with caption encoding as the input of commonsense decoder to generate the commonsense descriptions. Note that the multi-commonsense learning is applied on the commonsense decoder. And the decoder blocks are enhanced by our memory-routed multi-head attention (MMHA), memory module and contextual residual connections. Our contribution mainly focuses on the design of the multi-commonsense learning and innovative modules (e.g., MMHA and memory module), which are unveiled and discussed in detail in the following sub-sessions.

There are two settings for the two sub-tasks. **Completion task**: the ground truth caption and video are given to generate the commonsense descriptions. **Generation task**: given the video input, the event-wise captions should be predicted first, and then use both of them to predict the desired commonsense sentences.

#### 3.1 Encoder

Given a video, we use the pre-trained models including ResNet152 [14], SoundNet [2] and I3D [6] to extract the appearance feature, audio

feature and motion feature, respectively. Then we use a multimodal fusion to merge the three types of features. As shown in Figure 2 (a)(left), we use separate linear layers (FC) and LSTM [15] to individually encode the different features, and utilize the last hidden states of the LSTM as the final representations. Finally, the customized positional encoding and segment encoding are added to the final representations, which are concatenated together as the multimodal features. Taking the motion feature as an example

$$E^{3D} = SE^{3D} + PE^{3D} + \text{LSTM}(\text{FC}(V^{3D})), \quad (1)$$

where  $E^{3D}$  is the encoded motion feature and  $V^{3D}$  means motion feature. The  $SE^{3D}$  and  $PE^{3D}$  are 3D segment encoding and 3D positional encoding, respectively. Similarly, we can obtain the encoded audio feature  $E^{1D}$  and encoded appearance feature  $E^{2D}$ . Then we concatenate them together to get the multimodal features (MF).

#### 3.2 Decoder

In our decoder, we propose two main innovations: (1) multiple commonsense learning for semantic-level reasoning; (2) memory-routed multi-head attention cooperated with memory module for word-level reasoning. The first one is to improve the high-level inference ability from various commonsense semantics. The second one aims to mine the low-level reasoning from different words.

#### 3.3 Semantic-level Reasoning

**Multi-commonsense Learning.** Multi-commonsense learning is a training paradigm in which machine learning models are trained with data from multiple types of commonsense descriptions simultaneously, using shared representations to learn the common ideas between a collection of related commonsense. These shared representations increase data efficiency and can potentially yield a faster

learning speed for correlated descriptions. In this paper, we develop the transformer based language model [35] as three commonsense decoders for three particular commonsense domains. There are many different factors to consider when creating a shared architecture for multi-commonsense learning, such as the portion of the model’s parameters that will be shared between commonsense, and how to parameterize and combine commonsense-specific and shared modules. In our HybridNet, we share the parameters of encoder and caption decoder for all commonsense decoders. During the training, the commonsense decoder takes the video encoding  $\mathbf{v}$ , caption encoding  $\hat{\mathbf{s}}$  and ground truth of corresponding commonsense captions (e.g.,  $\mathbf{c}_{\text{att}}$ ,  $\mathbf{c}_{\text{eff}}$ ,  $\mathbf{c}_{\text{int}}$ ) as input to iteratively generate commonsense descriptions, which can be formulated as

$$\hat{\mathbf{c}}_{\text{att}} = \mathbf{D}_{\text{ATT}}(\mathbf{v}, \hat{\mathbf{s}}, \mathbf{c}_{\text{att}}), \quad (2)$$

$$\hat{\mathbf{c}}_{\text{eff}} = \mathbf{D}_{\text{EFF}}(\mathbf{v}, \hat{\mathbf{s}}, \mathbf{c}_{\text{eff}}), \quad (3)$$

$$\hat{\mathbf{c}}_{\text{int}} = \mathbf{D}_{\text{INT}}(\mathbf{v}, \hat{\mathbf{s}}, \mathbf{c}_{\text{int}}), \quad (4)$$

where  $\hat{\mathbf{c}}_{\text{att}}$ ,  $\hat{\mathbf{c}}_{\text{eff}}$ ,  $\hat{\mathbf{c}}_{\text{int}}$  are the generated commonsense sequences decoded by the corresponding commonsense decoders ( $\mathbf{D}_{\text{ATT}}$ ,  $\mathbf{D}_{\text{EFF}}$ ,  $\mathbf{D}_{\text{INT}}$ ). The loss function for the  $\mathbf{D}_{\text{ATT}}$  can be formulated as

$$\mathcal{L}_{\text{att}} = \sum_{t=1}^{N_{\text{att}}} \log p(\mathbf{y}_t | \mathbf{y}_{t-1}, [\mathbf{v}, \hat{\mathbf{s}}]; \Theta_{\text{att}}), \quad (5)$$

where  $\mathbf{y}_t$  denotes the one-hot vector probability of each word at time  $t$ ,  $N_{\text{att}}$  denotes the length of the attribute. The attribute decoder parameters  $\Theta_{\text{att}}$  are trained to maximize the log-likelihood over the training set. Similarly, we can obtain the objection functions of effect decoder and intention decoder like  $\mathcal{L}_{\text{eff}}$  and  $\mathcal{L}_{\text{int}}$ . The caption decoder can be optimized by

$$\mathcal{L}_{\text{cap}} = \sum_{t=1}^{N_{\text{cap}}} \log p(\mathbf{y}_t | \mathbf{y}_{t-1}, \mathbf{v}; \Theta_{\text{cap}}). \quad (6)$$

Finally, we train our HybridNet by using  $\mathcal{L} = \mathcal{L}_{\text{cap}} + \mathcal{L}_{\text{att}} + \mathcal{L}_{\text{eff}} + \mathcal{L}_{\text{int}}$  to jointly optimize our framework.

### 3.4 Word-level Reasoning

**Memory Module.** For any relevant videos, they may share similar patterns in their descriptions that can be used as good references for each other to help the generation process. Besides, the previous sequence can be recorded to guide the next word prediction for contextual consistency. To exploit such characteristics, we propose to use an extra component named memory module to enhance Transformer to learn from the previous word information and facilitate computing the interactions among previous information and the generation process.

As shown in Figure 3(b), our memory module uses a matrix  $\mathbf{M}$  to transfer its states over generation steps, where the states record the important word information with each row (namely, memory slot) representing some word information.<sup>1</sup> During the generation, the matrix is updated step-by-step by incorporating the output from previous steps. Then, at time step  $t$ , the matrix from the previous step,  $\mathbf{M}_{t-1}$ , is functionalized as the query and its concatenations with the previous output serve as the key and value

<sup>1</sup>Note that the rows (memory slots) and word states do not follow one-to-one mapping, where the entire matrix serves as a whole unit to deliver the word information.

to feed into the multi-head attention module. Given  $H$  heads used in Transformer, there are  $H$  sets of queries, keys and values via three linear transformations, respectively. For each head, we obtain the query, key and value in the memory module through  $\mathbf{Q} = \mathbf{M}_{t-1} \mathbf{W}_{\mathbf{q}}$ ,  $\mathbf{K} = [\mathbf{M}_{t-1}; \mathbf{y}_{t-1}] \mathbf{W}_{\mathbf{k}}$  and  $\mathbf{V} = [\mathbf{M}_{t-1}; \mathbf{y}_{t-1}] \mathbf{W}_{\mathbf{v}}$ , respectively, where  $\mathbf{y}_{t-1}$  is the embedding of the last output (at step  $t-1$ );  $[\mathbf{M}_{t-1}; \mathbf{y}_{t-1}]$  is the row-wise concatenation of  $\mathbf{M}_{t-1}$  and  $\mathbf{y}_{t-1}$ . The  $\mathbf{W}_{\mathbf{q}}$ ,  $\mathbf{W}_{\mathbf{k}}$  and  $\mathbf{W}_{\mathbf{v}}$  are the trainable weights of linear transformation of the query, key and value, respectively. Multi-head attention is used to model  $\mathbf{Q}$ ,  $\mathbf{K}$  and  $\mathbf{V}$  so as to depict relations of different patterns. As a result,

$$\mathbf{Z} = \text{softmax}(\mathbf{Q}\mathbf{K}^{\top} / \sqrt{d_k})\mathbf{V}, \quad (7)$$

where  $d_k$  is the dimension of  $\mathbf{K}$ , and  $\mathbf{Z}$  the output of the multi-head attention module. Consider that the memory module is performed in a recurrent manner along with the decoding process, it potentially suffers from gradient vanishing and exploding. Therefore, we introduce residual connections and a series of gate operations. The former is formulated as

$$\mathbf{M}'_t = f_{\text{mlp}}(\mathbf{Z} + \mathbf{M}_{t-1}) + \mathbf{Z} + \mathbf{M}_{t-1}, \quad (8)$$

where  $f_{\text{mlp}}(\cdot)$  refers to the multi-layer perceptron (MLP). Moreover, we apply the forget and input gates to balance the inputs from  $\mathbf{M}_{t-1}$  and  $\mathbf{y}_{t-1}$ , respectively. To ensure that  $\mathbf{y}_{t-1}$  can be used for computation with  $\mathbf{M}_{t-1}$ , it is extended to a matrix  $\mathbf{Y}_{t-1}$  by duplicating it to multiple rows. Therefore, the forget gate  $\mathbf{G}_t^f$  and input gate  $\mathbf{G}_t^i$  are formalized as

$$\mathbf{G}_t^f = \mathbf{Y}_{t-1} \mathbf{W}^f + \tanh(\mathbf{M}_{t-1}) \mathbf{U}^f, \quad (9)$$

$$\mathbf{G}_t^i = \mathbf{Y}_{t-1} \mathbf{W}^i + \tanh(\mathbf{M}_{t-1}) \mathbf{U}^i, \quad (10)$$

where  $\mathbf{W}^f$  and  $\mathbf{W}^i$  are trainable weights for  $\mathbf{Y}_{t-1}$  in each gate; similarly,  $\mathbf{U}^f$  and  $\mathbf{U}^i$  are the trainable weights for  $\mathbf{M}_{t-1}$  in each gate. The final output of the gate mechanism is formalized as

$$\mathbf{M}_t = \sigma(\mathbf{G}_t^f) \odot \mathbf{M}_{t-1} + \sigma(\mathbf{G}_t^i) \odot \tanh(\mathbf{M}'_t), \quad (11)$$

where  $\odot$  refers to the Hadamard product and  $\sigma$  is the sigmoid function. The  $\mathbf{M}_t$  is the output of the entire memory module at step  $t$ , which is fed into the MMHA for routing the decoder.

**Memory-routed Multi-Head Attention (MMHA).** We design the MMHA in each decoder block to bridge the previous word states and the next predicting word. We think the previous words can cooperate with the multimodal features to better reason the generation processes. As shown in Figure 2 (b), given the input representation  $\mathbf{X}$ , we get the query  $\mathbf{Q}$ , key  $\mathbf{K}$  and value  $\mathbf{V}$  through three linear layers. Then a multi-head scale dot-product operation is utilized to generate the original attention map

$$\mathbf{A}_o = \text{Attention}(\mathbf{X}) = \mathbf{Q}\mathbf{K}^{\top} / \sqrt{d_k}, \quad (12)$$

where  $\mathbf{A}_o$  denotes the original attention map and  $d_k$  is the hidden dimension size of  $\mathbf{K}$ . Because of the guidance of the memory module, we can obtain the memory state  $\mathbf{M}$  from the memory module that records the previous sequence state. The  $\mathbf{M}$  is regarded as another  $\mathbf{Q}$  in the MMHA. Then another multi-head scale dot-product operation is applied to getting the conditional attention map  $\mathbf{A}_c$ . Assume there are  $K$  heads in each layer, then we get  $K$  conditional attention maps. They construct a tensor  $\mathbf{A}_c \in \mathbb{R}^{N \times N \times K}$  ( $N$  is the sequence length), which can be viewed as a  $N \times N$  image with  $K$  channels. Taking this

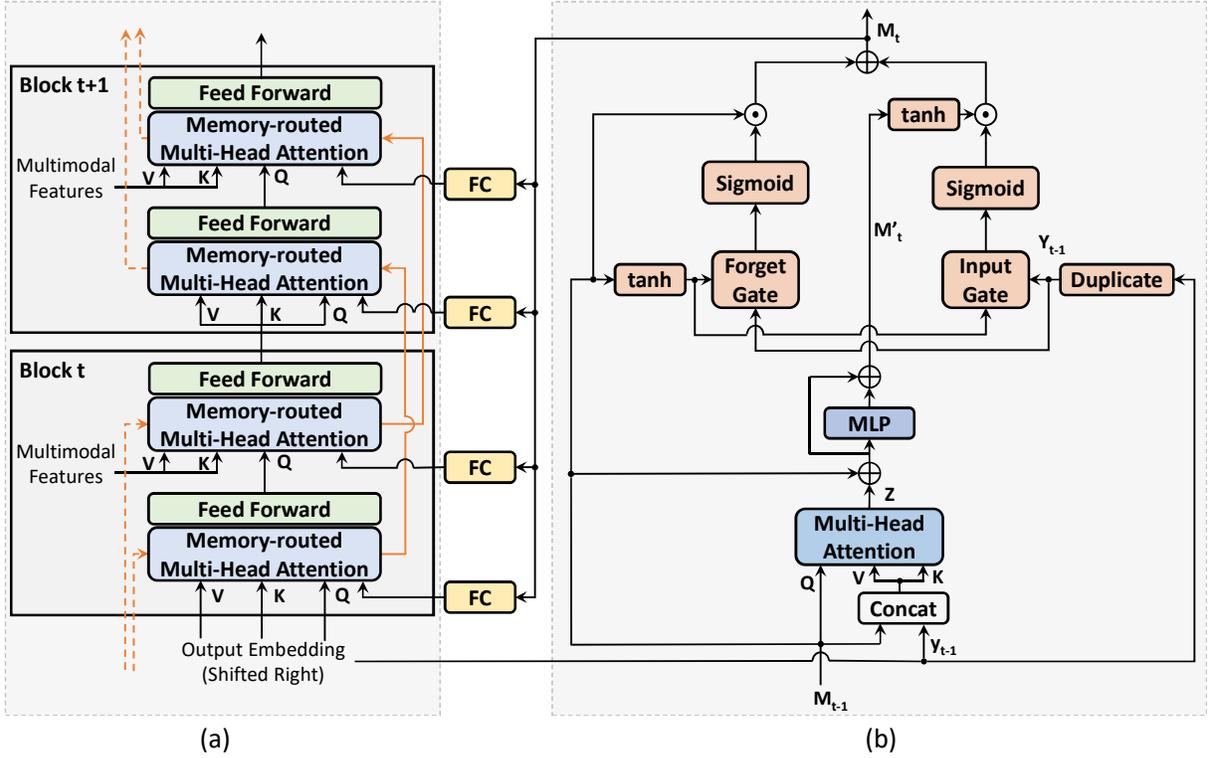


Figure 3: (a) Our residual connections that are marked as orange propagate the contextual information between different decoder blocks. Note that we omit additions and norm in the figure for brevity. (b) Our memory module records the information from previous generation processes.  $\odot$  denotes the Hadamard product.

as input, we adopt one 2D-convolutional layer with  $3 \times 3$  kernels to capture the evolution of attention patterns, as this inductive bias emphasizes local details and produces more precise attention maps by reasoning on previous ones. The output channel is also set to be  $K$ , so the attention maps of all heads can be generated jointly.

To prevent foreseeing subsequent positions, we improve the 2D-convolutional layer by proposing a triangle convolution. The triangle convolution can be implemented as follows: (1) executing standard  $3 \times 3$  convolution with masks in the upper-right corner; (2) after convolution, shifting the entire attention matrix by one pixel to the bottom and one pixel to the right. We apply a ReLU activation after each 2D-convolution layer to provide non-linearity and sparsity. After the triangle convolution, the result attention map  $A_{triangle}$  is combined with the input conditional attention map  $A_c$ , original attention map  $A_o$  and the attention maps from previous blocks  $A_{previous}$ . Mathematically,

$$\begin{aligned} A_{condition} &= \alpha \cdot A_{triangle} + (1 - \alpha) \cdot A_c, \\ A_{guide} &= \beta \cdot A_{condition} + (1 - \beta) \cdot A_o, \\ A_{merge} &= \gamma \cdot A_{previous} + (1 - \gamma) \cdot A_{guide}, \end{aligned} \quad (13)$$

where  $\alpha, \beta, \gamma \in [0, 1]$  are hyper-parameters for linear combination. In our experiments, the values of  $\alpha, \beta$  and  $\gamma$  are chosen empirically on the validation set for each task. Note that the  $A_{previous}$  is the residual attention map of previous blocks via contextual residual connections, which can be seen in Figure 3 (a). These contextual

residual connections aim to bridge the gap of attention maps between different blocks for improving the quality of attention maps globally. Finally, the output of the MMHA can be obtained by

$$A_{out} = \text{softmax}(\text{MASK}(A_{merge})), \quad (14)$$

$$X_{out} = \text{FC}(A_{out} V^T), \quad (15)$$

where  $A_{out}$  is the attention map of current MMHA and the  $X_{out}$  means the output representation of the MMHA.

## 4 EXPERIMENTS

### 4.1 Datasets and Evaluation

We evaluate our proposed HybridNet and compare it with other state-of-the-art methods on Video-to-Commonsense (V2C) [11] benchmark, a representative large-scale video-based commonsense captioning dataset containing a total of 121,618 captions derived from 9,721 video scenes. The dataset is officially split into a training set consisting of 6819 videos with 85,100 captions, a test set containing 2903 videos with 36,518 captions. We follow this data partition in all experiments. We have done some statistics on the V2C and found that there are 5 candidate descriptions per video for intention, effect, and attribute respectively. And the number of candidate captions for each video is not fixed. For the training set, there are 766 samples contain 1~4 captions per video. The number of the video containing 5~14 captions is 3082. The test set also has the same

**Table 1: Evaluation of V2C completion task and generation task in terms of the attribute, effect and intention by using CIDER, BLEU, Rouge, and Meteor metrics. We use only BLEU-1 to evaluate the attribute generation on the completion task since the average length of the ground truth is just less than 2. “Attribute+C” means the attribute descriptions and the predicted event-wise captions on the generation task. The best performing results are marked in red.**

	Relation	Model	CIDER	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L
Completion Task	Attribute	S2VT [37]	-	35.9	-	-	-	-	-
		Attention-Enc-Dec [12]	-	38.3	-	-	-	-	-
		Dense Captioner [45]	-	46.0	-	-	-	-	-
		Video CMS Transformer [11]	-	47.3	-	-	-	-	-
		Our HybridNet	-	<b>58.7<sup>+11.4</sup></b>	-	-	-	-	-
	Effect	S2VT [37]	28.3	24.9	18.6	16.2	14.3	15.4	22.1
		Attention-Enc-Dec [12]	29.5	26.5	19.4	18.8	15.1	17.5	23.9
		Dense Captioner [45]	36.9	33.7	24.8	21.0	20.2	20.0	29.9
		Video CMS Transformer [11]	37.3	34.8	25.9	22.5	20.4	20.8	30.6
		Our HybridNet	<b>66.2<sup>+28.9</sup></b>	<b>49.0<sup>+14.2</sup></b>	<b>42.9<sup>+17.0</sup></b>	<b>40.3<sup>+17.8</sup></b>	<b>38.8<sup>+18.4</sup></b>	<b>30.0<sup>+9.2</sup></b>	<b>41.5<sup>+10.9</sup></b>
	Intention	S2VT [37]	51.8	48.4	39.9	34.3	26.4	23.3	44.3
		Attention-Enc-Dec [12]	52.1	51.1	42.6	35.5	28.2	24.3	48.0
		Dense Captioner [45]	60.3	59.3	47.0	37.3	31.5	28.0	53.1
		Video CMS Transformer [11]	62.0	60.8	48.4	39.1	34.1	28.5	54.6
		Our HybridNet	<b>92.6<sup>+30.6</sup></b>	<b>69.4<sup>+8.6</sup></b>	<b>60.5<sup>+12.1</sup></b>	<b>55.4<sup>+16.3</sup></b>	<b>53.1<sup>+19.0</sup></b>	<b>35.8<sup>+7.3</sup></b>	<b>60.1<sup>+5.5</sup></b>
Generation Task	Attribute+C	S2VT [37]	38.5	69.1	53.6	42.0	32.3	23.9	59.1
		Attention-Enc-Dec [12]	34.0	67.0	51.7	40.7	31.4	23.3	58.0
		Dense Captioner [45]	36.8	68.4	52.1	39.8	30.0	24.1	57.7
		Video CMS Transformer [11]	40.2	70.2	54.8	42.7	32.6	24.7	59.0
		Our HybridNet	<b>41.6<sup>+1.4</sup></b>	<b>71.3<sup>+1.1</sup></b>	<b>57.0<sup>+2.2</sup></b>	<b>45.6<sup>+2.9</sup></b>	<b>35.7<sup>+3.1</sup></b>	<b>25.5<sup>+0.8</sup></b>	<b>60.4<sup>+1.4</sup></b>
	Effect+C	S2VT [37]	29.9	69.8	54.1	39.9	29.1	21.9	55.3
		Attention-Enc-Dec [12]	26.1	70.2	51.8	38.6	28.7	22.5	53.5
		Dense Captioner [45]	30.6	72.1	54.5	42.3	33.2	25.2	56.0
		Video CMS Transformer [11]	32.1	72.5	56.1	44.3	35.2	25.6	57.4
		Our HybridNet	<b>34.2<sup>+2.1</sup></b>	<b>73.2<sup>+0.7</sup></b>	<b>57.4<sup>+1.3</sup></b>	<b>46.3<sup>+2.0</sup></b>	<b>37.2<sup>+2.0</sup></b>	<b>26.3<sup>+0.7</sup></b>	<b>58.3<sup>+0.9</sup></b>
	Intention+C	S2VT [37]	35.4	71.3	53.9	41.3	31.2	21.6	58.6
		Attention-Enc-Dec [12]	33.2	75.4	59.4	45.1	33.5	24.6	59.6
		Dense Captioner [45]	37.0	76.1	60.2	46.7	35.9	26.5	60.9
		Video CMS Transformer [11]	37.8	76.2	61.2	48.1	37.3	26.9	61.9
		Our HybridNet	<b>40.4<sup>+2.6</sup></b>	<b>77.5<sup>+1.3</sup></b>	<b>62.9<sup>+1.7</sup></b>	<b>50.4<sup>+2.3</sup></b>	<b>40.2<sup>+2.9</sup></b>	<b>27.8<sup>+0.9</sup></b>	<b>62.9<sup>+1.0</sup></b>

challenge. Hence, the V2C is challenging because of the complex and diverse language, multiple scenes and hard inference types as mentioned in [11]. Followed other works [11, 12, 37, 45], for two sub-tasks, we measure the performance of our proposed method via Meteor [3], Rouge [22], CIDEr [36] and BLEU (n=1-4) [24].

## 4.2 Implementation Details

We conduct all experiments by using a single NVIDIA 3090 card on a single server. We implement our proposed HybridNet and re-implement other state-of-the-art methods via PyTorch [25] and python3.8 to train and test. The Nvidia CUDA of 11.1 and cuDNN of 8.0 are utilized for acceleration. In our HybridNet, each decoder consists of 6 transformer blocks with 8 attention heads. Note that the traditional multi-head attention is replaced with our memory-routed multi-head attention in the decoder block of our HybridNet. Unless otherwise noted, settings are the same for all experiments. During the training, we set the batch size to be 128 for one GPU and use the Adam [17] optimizer with 5000 warm-up steps, and learning rate initialized at  $1e-4$ , and a dropout probability of 0.1

after the residual layer. It takes about 10 hours for the training set. The number of epoch is 800. We set the hyper-parameters for linear combination of attention maps in our experiments, including  $\alpha=0.1$ ,  $\beta=0.4$  and  $\gamma=0.1$ . During the test time, we validate learning outcomes after each learning epoch and select the model weights with the best CIDEr as our final results.

## 4.3 Results and Comparisons

**Quantitative Results.** We report our state-of-the-art results of the test on V2C [11] dataset for two tasks in Table 1. The previous works individually generate different types of commonsense captions by using separate networks. They can not predict all commonsense descriptions in a unified way. However, our HybridNet can generate all commonsense captions and benefit from their interaction.

On completion task, we used the reported results [11] for a fair comparison. As we can see, our HybridNet achieves the best performance on all metrics compared with the state-of-the-art methods [11, 12, 37, 45]. On the attribute part, our HybridNet performs



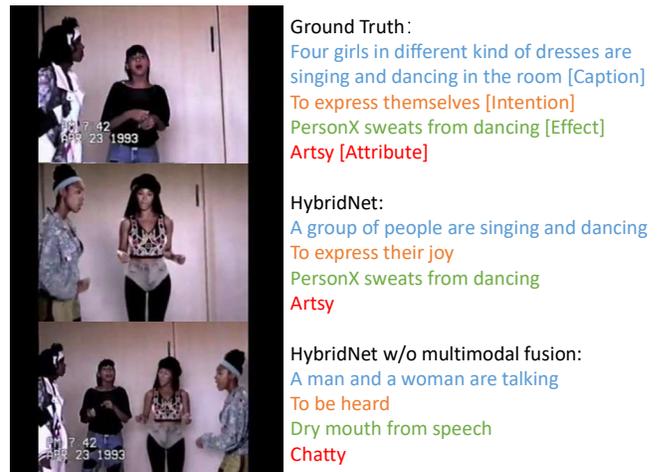
**Figure 4: Qualitative visual results on completion task (left) and generation task (right). The (0)–(4) denote the prediction results of our HybridNet, Video CMS Transformer [11], Dense Captioner [45], S2VT [37] and Attention-Enc-Dec [12] respectively. Compared with the other state-of-the-art methods, our HybridNet can generate more precise and logical descriptions.**

11.4% improvement on BLEU-1 better than the video CMS transformer [11]. On the effect and intention part, our network tends to bring better gains on the metrics that is applied for evaluating long and logical sentences, like the improvements on BLEU-4 (e.g., +19.0%, +18.4%) are better than the BLEU-1 (e.g., +8.6%, +14.2%). It can validate that our network can capture well the long-range context to make the semantics of the generated long sentences more reasonable and logical. The possible reason is our MMHA can learn the correlation between history sequence and next word prediction, which is cooperated with our memory module and contextual residual connections for word-level reasoning.

On the generation task, we re-implement and train the previous methods by using their official codes. Then we use some evaluation metrics that are widely used and accepted for language generation task to objectively estimate the performance of different methods.<sup>2</sup> In Table 1, our HybridNet achieves the best performances of 34.2%, 73.2%, 57.4%, 37.2%, 26.3% and 58.3% with respect to 7 metrics (i.e. CIDEr, BLEU-1 to BLEU-4, Meteor and Rouge-L) compared with 4 state-of-the-art methods on effect+C part. Our method also gets the best scores on other evaluation metrics. It can further prove that our network not only can perform well on the completion task but also can handle the more challenging task like the generation task.

**Qualitative Results.** Figure 4 shows the comparison results by different methods on completion task on the left and generation task on the right. On the completion task, our HybridNet can predict more precise intention results such as “to score a basket” compared with other methods. As we can see, other methods mainly focus on the vague intention expression (e.g., play a game) rather than the specific sports. On the generation task, our network still performs the best to jointly predict all the correct results. Although some methods can generate the correct caption (e.g., running around a race) and attribute (e.g., athletic), they still fail to provide the right intention and effect. On the contrary, thanks to the correct

<sup>2</sup>The previous methods only use the human evaluation for the generation task, which is too subjective.



**Figure 5: Visualization to compare the difference between HybridNet and HybridNet without multimodal fusion. The multimodal signal like audio in this example is necessary to distinguish between singing and talking.**

parsing into the caption and attribute, our method can also correctly provide reasonable and logical effect result (e.g., wins the race) and locate at the accurate intention (e.g., compete against others). Such semantic-level interaction can be achieved by our multi-commonsense learning, which can further demonstrate that our model can successfully learn the constraint and interaction of the multiple commonsense semantics.

#### 4.4 Ablation Studies

**Effect of Multimodal Features.** In Table 2, the advantage of multimodal fusion increases BLEU-1 by 0.1% (58.6% vs 58.7%) on attribute completion, 0.7% on effect completion and 0.7% on intention completion. In Figure 5, we display the effectiveness and importance

Table 2: Ablation studies on two sub-tasks. The “multimodal”, “multi-cms”, “MMHA” and “CRC” means the multimodal fusion, multi-commonsense learning, memory-routed multi-head attention and contextual residual connections, respectively.

Relation		multimodal	multi-cms	MMHA	CRC	CIDER	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L	
Completion Task	Attribute					-	47.3	-	-	-	-	-	
		✓				-	56.5	-	-	-	-	-	
		✓	✓			-	57.6	-	-	-	-	-	
		✓	✓	✓		-	57.9	-	-	-	-	-	
		✓	✓	✓	✓	-	58.6	-	-	-	-	-	
						-	<b>58.7</b>	-	-	-	-	-	
	Effect						37.3	34.8	25.9	22.5	20.4	20.8	30.6
		✓					60.6	45.6	39.0	35.9	34.8	28.0	38.9
		✓	✓				63.1	45.8	39.2	36.0	33.9	28.4	40.0
		✓	✓	✓			64.8	47.2	40.9	38.0	36.2	29.2	40.7
		✓	✓	✓	✓	✓	65.8	48.3	42.3	39.7	38.0	29.4	41.0
							<b>66.2</b>	<b>49.0</b>	<b>42.9</b>	<b>40.3</b>	<b>38.8</b>	<b>30.0</b>	<b>41.5</b>
	Intention						62.0	60.8	48.4	39.1	34.1	28.5	54.6
		✓					84.0	66.7	57.7	51.5	49.1	33.9	58.1
		✓	✓				86.4	67.3	57.7	51.9	49.2	34.4	58.6
✓		✓	✓			91.8	68.7	59.3	53.7	51.1	35.5	60.1	
✓		✓	✓	✓	✓	92.0	68.7	59.5	53.8	51.1	35.6	60.1	
						<b>92.6</b>	<b>69.4</b>	<b>60.5</b>	<b>55.4</b>	<b>53.1</b>	<b>35.8</b>	<b>60.1</b>	
Generation Task	Attribute+C					40.2	70.2	54.8	42.7	32.6	24.7	59.0	
		✓				40.2	70.0	55.7	44.7	35.2	25.0	59.8	
		✓	✓			40.2	70.0	56.5	45.5	35.4	25.1	60.2	
		✓	✓	✓		41.0	71.1	56.9	45.6	35.6	25.2	60.3	
		✓	✓	✓	✓	40.6	69.6	55.4	44.2	34.6	24.9	59.7	
							<b>41.6</b>	<b>71.3</b>	<b>57.0</b>	<b>45.6</b>	<b>35.7</b>	<b>25.5</b>	<b>60.4</b>
	Effect+C						32.1	72.5	56.1	44.3	35.2	25.6	57.4
		✓					32.3	72.4	56.7	45.8	36.0	25.9	57.7
		✓	✓				32.4	72.8	57.0	46.1	36.4	26.1	58.1
		✓	✓	✓			33.1	72.9	57.2	46.2	36.8	26.3	58.2
		✓	✓	✓	✓	✓	32.7	71.9	56.4	45.6	37.1	26.0	57.8
							<b>34.2</b>	<b>73.2</b>	<b>57.4</b>	<b>46.3</b>	<b>37.2</b>	<b>26.3</b>	<b>58.3</b>
	Intention+C						37.8	76.2	61.2	48.1	37.3	26.9	61.9
		✓					37.7	76.7	62.0	49.7	39.8	27.2	62.0
		✓	✓				38.2	77.1	62.2	50.1	40.1	27.4	62.1
✓		✓	✓			39.6	77.2	62.4	50.3	40.2	27.6	62.9	
✓		✓	✓	✓	✓	38.9	76.5	62.3	49.8	39.8	27.2	62.9	
						<b>40.4</b>	<b>77.5</b>	<b>62.9</b>	<b>50.4</b>	<b>40.2</b>	<b>27.8</b>	<b>62.9</b>	

of multimodal fusion in the video-based commonsense captioning task. As we can see, the model guided with the audio feature and motion feature can better recognize the singing and dancing rather than talking. Based on the assistance of the multimodal fusion, our network can easily infer the correct commonsense descriptions.

**Effect of Multi-Commonsense Learning.** The effect of our multi-commonsense learning is apparently to boost BLEU-1 by around 9.2% (47.3%→56.5%), 10.8% (34.8%→45.6%) and 5.9% (60.8%→66.7%) from the baseline on the completion task w.r.t. attribute, effect and intention in Table 2. The effectiveness of our multi-commonsense learning can also be generalized to the generation task. It gets great scores on BLEU-3 and BLEU-4 compared with the baseline (e.g., 42.7%→44.7%, 44.3%→45.8% and 48.1%→49.7% on the attribute, effect and intention in terms of BLEU-3). These significant improvements can support the advantage and effectiveness of our multi-commonsense learning for semantic-level reasoning.

**Effect of Memory-routed Multi-Head Attention (MMHA).** In Table 2, we observe that our MMHA can bring improvements on

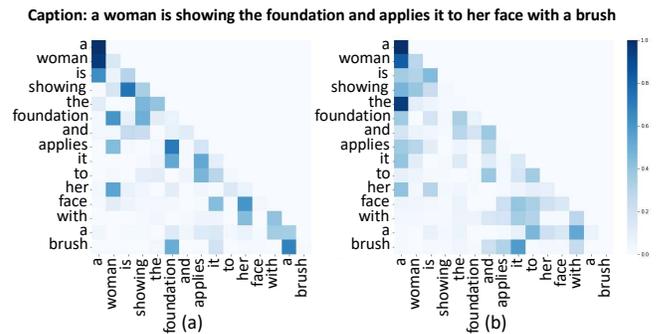


Figure 6: (a) Visualization of the attention map from our MMHA. (b) Visualization of the attention map from the traditional multi-head attention.

all metrics in terms of the completion and generation tasks. For example, our MMHA can promote the network on the intention

generation from 38.2% to 39.6% on CIDEr, and from 62.1% to 62.9% on Rouge-L. And it can reach a higher improvement of 5.4% CIDEr on the intention completion. In Figure 6 (a), the model with the MMHA can successfully reason the correlation between the pronoun “it” and the noun “foundation” since it remembers the history sequence that is applied to infer the next word. However, in Figure 6 (b), the model without MMHA is failed to reason the word-level correlation between “it” and “foundation”. Besides, when predicting the pronoun “her”, our MMHA infers the logical connection with the “woman”, while the model without MMHA ignores the history memory and wrongly connects “her” with “a”. The quantitative and qualitative results demonstrate the effectiveness and availability of our MMHA in video-based commonsense reasoning.

**Effect of Contextual Residual Connection.** Our contextual residual connection (CRC) aims to enhance the attention maps with a global perspective between different layers. In Table 2, our CRC collaborated with other modules is evaluated on two tasks and increases 0.8% CIDEr (39.6%→40.4%) and 2.6% BLEU-4 (36.2%→38.8%) on intention generation and effect completion, respectively. It can support the feasibility of incorporating the contextual residual connection into the decoder of our HybridNet.

## 5 CONCLUSION

In this paper, we propose a novel Hybrid Reasoning Network (HybridNet) for video-based commonsense captioning, which is jointly optimized by semantic-level reasoning and word-level reasoning. Multi-commonsense learning is built to achieve the semantic-level reasoning by jointly training different commonsense types in a unified network. To perform the word-level reasoning, a memory-routed multi-head attention (MMHA) is presented to inject the previous relational states from the memory module into the decoder for updating the attention maps. Then the updated attention maps are merged with other attention maps from previous decoder blocks by using contextual residual connections. The final obtained attention maps can be used to evolve representation. Our HybridNet achieves a new state-of-the-art on the large-scale Video-to-Commonsense benchmark and provided abundant analysis to demonstrate the effectiveness of our proposed modules.

## 6 ACKNOWLEDGEMENTS

This work was supported in part by National Natural Science Foundation of China under Grant No.U1811461, in part by the National Key Research and Development Project under Grant No.2020AAA0106600, in part by Natural Science Foundation of Guangdong Province, China under Grant No.2018B030312002, and in part by the Major Program of Guangdong Basic and Applied Research No.2019B030302002.

## REFERENCES

- [1] Stanislaw Antol, Aishwarya Agrawal, Jiaseen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*. 2425–2433.
- [2] Yusuf Aytar, Carl Vondrick, and Antonio Torralba. 2016. Soundnet: Learning sound representations from unlabeled video. In *Advances in Neural Information Processing Systems*.
- [3] Satyanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the acl workshop on intrinsic and extrinsic evaluation measures for machine translation and/or summarization*. 65–72.
- [4] Chandra Bhagavatula, Ronan Le Bras, Chaitanya Malaviya, Keisuke Sakaguchi, Ari Holtzman, Hannah Rashkin, Doug Downey, Wen-tau Yih, and Yejin Choi. 2020. Abductive Commonsense Reasoning. In *ICLR*.
- [5] Yonatan Bisk, Rowan Zellers, Ronan LeBras, Jianfeng Gao, and Yejin Choi. 2020. PIQA: Reasoning about Physical Commonsense in Natural Language. In *AAAI*. 7432–7439.
- [6] Joao Carreira and Andrew Zisserman. 2017. Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*.
- [7] David Chen and William B Dolan. 2011. Collecting highly parallel data for paraphrase evaluation. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*. 190–200.
- [8] Thilini Cooray, Ngai-Man Cheung, and Wei Lu. 2020. Attention-Based Context Aware Reasoning for Situation Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 4736–4745.
- [9] Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, José MF Moura, Devi Parikh, and Dhruv Batra. 2017. Visual dialog. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 326–335.
- [10] Pradipto Das, Chenliang Xu, Richard F Doell, and Jason J Corso. 2013. A thousand frames in just a few words: Lingual description of videos through latent topics and sparse object stitching. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2634–2641.
- [11] Zhiyuan Fang, Tejas Gokhale, Pratay Banerjee, Chitta Baral, and Yezhou Yang. 2020. Video2Commonsense: Generating Commonsense Descriptions to Enrich Video Captioning. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics, Online. 840–860. <https://doi.org/10.18653/v1/2020.emnlp-main.61>
- [12] Lianli Gao, Zhao Guo, Hanwang Zhang, Xing Xu, and Heng Tao Shen. 2017. Video captioning with attention-based LSTM and semantic consistency. *IEEE Transactions on Multimedia* 19, 9 (2017), 2045–2055.
- [13] Abhinav Gupta, Praveen Srinivasan, Jianbo Shi, and Larry S Davis. 2009. Understanding videos, constructing plots learning a visually grounded storyline model from annotated videos. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2012–2019.
- [14] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [15] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [16] Jena D Hwang, Chandra Bhagavatula, Ronan Le Bras, Jeff Da, Keisuke Sakaguchi, Antoine Bosselut, and Yejin Choi. 2020. Comet-atomic 2020: On symbolic and neural commonsense knowledge graphs. *arXiv preprint arXiv:2010.05953* (2020).
- [17] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [18] Jie Lei, Licheng Yu, Tamara L Berg, and Mohit Bansal. 2019. TVQA+: Spatio-temporal grounding for video question answering. In *ACL*.
- [19] Ruiyu Li, Makarand Tapaswi, Renjie Liao, Jiaya Jia, Raquel Urtasun, and Sanja Fidler. 2017. Situation recognition with graph neural networks. In *Proceedings of the IEEE International Conference on Computer Vision*. 4173–4182.
- [20] Yikang Li, Nan Duan, Bolei Zhou, Xiao Chu, Wanli Ouyang, Xiaogang Wang, and Ming Zhou. 2018. Visual question generation as dual task of visual question answering. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6116–6124.
- [21] Yuncheng Li, Yale Song, Liangliang Cao, Joel Tetreault, Larry Goldberg, Alejandro Jaimes, and Jiebo Luo. 2016. TGIF: A new dataset and benchmark on animated GIF description. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 4641–4650.
- [22] Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*. 74–81.
- [23] Pingbo Pan, Zhongwen Xu, Yi Yang, Fei Wu, and Yueting Zhuang. 2016. Hierarchical recurrent neural encoder for video representation with application to captioning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 1029–1038.
- [24] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. BLEU: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*. Association for Computational Linguistics, 311–318.
- [25] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. 2017. Automatic Differentiation in PyTorch. In *NIPS 2017 Workshop on Autodiff* (Long Beach, California, USA). <https://openreview.net/forum?id=BjJsrnFCZ>
- [26] Hamed Pirsiavash, Carl Vondrick, and Antonio Torralba. 2014. Inferring the why in images. *Massachusetts Inst of Tech Cambridge* (2014).
- [27] Bryan A. Plummer, Liwei Wang, Chris M. Cervantes, Juan C. Caicedo, Julia Hockenmaier, and Svetlana Lazebnik. 2015. Flickr30k Entities: Collecting Region-to-Phrase Correspondences for Richer Image-to-Sentence Models. In *ICCV*.
- [28] Michaela Regneri, Marcus Rohrbach, Dominikus Wetzels, Stefan Thater, Bernd Schiele, and Manfred Pinkal. 2013. Grounding action descriptions in videos. *Transactions of the Association for Computational Linguistics* 1 (2013), 25–36.

- [29] Anna Rohrbach, Marcus Rohrbach, Wei Qiu, Annemarie Friedrich, Manfred Pinkal, and Bernt Schiele. 2014. Coherent multi-sentence video description with variable level of detail. In *German conference on pattern recognition*. Springer, 184–195.
- [30] Anna Rohrbach, Marcus Rohrbach, Niket Tandon, and Bernt Schiele. 2015. A dataset for movie description. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 3202–3212.
- [31] Maarten Sap, Ronan LeBras, Emily Allaway, Chandra Bhagavatula, Nicholas Lourie, Hannah Rashkin, Brendan Roof, Noah A Smith, and Yejin Choi. 2018. Atomic: An atlas of machine commonsense for if-then reasoning. In *AAAI*.
- [32] Maarten Sap, Hannah Rashkin, Derek Chen, Ronan Le Bras, and Yejin Choi. 2019. Social IQa: Commonsense Reasoning about Social Interactions. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 4453–4463.
- [33] Alon Talmor, Jonathan Herzig, Nicholas Lourie, and Jonathan Berant. 2019. CommonsenseQA: A Question Answering Challenge Targeting Commonsense Knowledge. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. 4149–4158.
- [34] Hao Tan and Mohit Bansal. 2019. LXMERT: Learning Cross-Modality Encoder Representations from Transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*. 5103–5114.
- [35] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in neural information processing systems*. 5998–6008.
- [36] Ramakrishna Vedantam, C. Lawrence Zitnick, and Devi Parikh. 2015. CIDEr: Consensus-based image description evaluation. In *CVPR*. IEEE Computer Society, 4566–4575. <http://dblp.uni-trier.de/db/conf/cvpr/cvpr2015.html#VedantamZP15>
- [37] Subhashini Venugopalan, Marcus Rohrbach, Jeffrey Donahue, Raymond Mooney, Trevor Darrell, and Kate Saenko. 2015. Sequence to sequence-video to text. In *Proceedings of the IEEE international conference on computer vision*. 4534–4542.
- [38] Oriol Vinyals, Alexander Toshev, Samy Bengio, and Dumitru Erhan. 2016. Show and tell: Lessons learned from the 2015 mscoco image captioning challenge. *IEEE transactions on pattern analysis and machine intelligence* 39, 4 (2016), 652–663.
- [39] Carl Vondrick, Deniz Oktay, Hamed Pirsiavash, and Antonio Torralba. 2016. Predicting motivations of actions by leveraging text. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2997–3005.
- [40] Tan Wang, Jianqiang Huang, Hanwang Zhang, and Qianru Sun. 2020. Visual commonsense r-cnn. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10760–10770.
- [41] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. Msr-vtt: A large video description dataset for bridging video and language. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5288–5296.
- [42] Mark Yatskar, Luke Zettlemoyer, and Ali Farhadi. 2016. Situation recognition: Visual semantic role labeling for image understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 5534–5542.
- [43] Weijiang Yu, Jingwen Zhou, Weihao Yu, Xiaodan Liang, and Nong Xiao. 2019. Heterogeneous Graph Learning for Visual Commonsense Reasoning. In *Advances in Neural Information Processing Systems*. 2769–2779.
- [44] Rowan Zellers, Yonatan Bisk, Ali Farhadi, and Yejin Choi. 2019. From recognition to cognition: Visual commonsense reasoning. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 6720–6731.
- [45] Luowei Zhou, Yingbo Zhou, Jason J Corso, Richard Socher, and Caiming Xiong. 2018. End-to-end dense video captioning with masked transformer. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8739–8748.
- [46] Yukun Zhu, Ryan Kiros, Rich Zemel, Ruslan Salakhutdinov, Raquel Urtasun, Antonio Torralba, and Sanja Fidler. 2015. Aligning books and movies: Towards story-like visual explanations by watching movies and reading books. In *Proceedings of the IEEE international conference on computer vision*. 19–27.

Table 3: Ablation studies on different fusion for multimodal information on completion task.

	MLP Fusion	Concat Fusion	CIDER	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L
Attribute	✓		-	58.1	-	-	-	-	-
		✓	-	58.7	-	-	-	-	-
Effect	✓		65.8	47.6	41.9	38.7	37.2	29.5	40.6
		✓	66.2	49.0	42.9	40.3	38.8	30.0	41.5
Intention	✓		91.8	68.9	59.8	54.7	51.9	35.6	60.1
		✓	92.6	69.4	60.5	55.4	53.1	35.8	60.1

Table 4: Quantitative results of the example in figure 5

	Multimodal	CIDER	BLEU-1	BLEU-2	BLEU-3	BLEU-4	METEOR	ROUGE-L
Attribute	✓	40.7	69.2	56.9	44.9	35.2	25.2	58.4
		39.5	68.3	55.6	43.8	34.6	24.8	57.3
Effect	✓	33.8	72.1	56.2	45.8	36.7	25.4	57.9
		32.6	72.0	55.8	44.7	35.9	24.6	56.8
Intention	✓	40.2	76.8	61.5	49.8	40.1	27.9	61.6
		39.8	76.2	61.3	49.4	39.5	27.2	61.2

## A APPENDIX

### A.1 Speed and Parameter Comparison

We have tested the inference time of [11] and our models by using the FLOPs metric. The CMS Transformer [11] needs 4.55 GFLOPs and our model only needs 2.93 GFLOPs. It can demonstrate the [11] is more time-consuming compared with our model. Besides, the parameters of our model are 103.4M, which is smaller than CMS (159.1M).

### A.2 Quantitative Analysis of the example in Figure 5

We tested the quantitative results of the example in Figure 5 between our model and [11]. The results are reported in Table 4. The example comes from the generation task. On attribute part, the [11] achieves 68.3 BLUE-1 and our model performs 69.2. Moreover, our model outperforms [11] by 1.1% Rough-L (57.9 vs 56.8) on effect part, and by 0.7% Meteor on intention part.

### A.3 The Effect of Different Fusion for Multimodal Inputs

For the encoder part, we mainly focus on the multimodal inputs. Hence the three pretrained encoders aim to extract multimodal features (e.g., 1D, 2D, 3D). In addition, we have done experiments to compare our fusion method in the encoder part with other fusion way. We replace our concatenation with MLP layers and train the whole network on the completion task. As shown in Table 3, the MLP fusion method achieves 65.8 Cider on effect and 91.8 on intention. Our concatenation performs 66.2 Cider on effect and 92.6 on intention. We believe that simply stacking parameters to increase the fusion block (MLP fusion) will not necessarily improve the video captioning. In other words, the multimodal inputs are the key to improve the final performance.

### A.4 Human Evaluation

We follow the human evaluation setting in [11] and hire 5 students for estimate the results for each model. We estimate our model on generation task with human evaluations and also compare it with the gold annotations from [11]. Our proposed model gets 66.45/78.96/71.78 on Effect/Attribute/Intention by using human evaluation. The Gold Annotations of human evaluation in [11] are 75.19/83.03/80.11 based on the ground truth.