# Holistic Multi-modal Memory Network for Movie Question Answering

Anran Wang, Anh Tuan Luu, Chuan-Sheng Foo, Hongyuan Zhu, Yi Tay, Vijay Chandrasekhar

*Abstract*—Answering questions according to multi-modal context is a challenging problem as it requires a deep integration of different data sources. Existing approaches only employ partial interactions among data sources in one attention hop. In this paper, we present the Holistic Multi-modal Memory Network (HMMN) framework which fully considers the interactions between different input sources (multi-modal context, question) in each hop. In addition, it takes answer choices into consideration during the context retrieval stage. Therefore, the proposed framework effectively integrates multi-modal context, question, and answer information, which leads to more informative context retrieved for question answering. Our HMMN framework achieves state-of-the-art accuracy on MovieQA dataset. Extensive ablation studies show the importance of holistic reasoning and contributions of different attention strategies.

*Index Terms*—Question answering, multi-modal learning, MovieQA.

## I. INTRODUCTION

With recent tremendous progress in computer vision and natural language processing, increasing attention has been drawn to the joint understanding of the visual and textual semantics. Related topics such as image-text retrieval [17], [7], [29], image/video captioning [33], [21], [22], [20], [35], [34], visual question answering (VQA) [1], [4], [11], [9] have been intensely explored. In particular, VQA remains a relatively challenging task, and it requires understanding of a given image to answer the question.

Due to the multi-modal nature of the world, developing question answering (QA) systems that are able to attain an understanding of the world based on multiple sources of information is a natural next step. Several QA datasets incorporating multiple data modalities have recently been developed towards this end [12], [26], [13]. In this work, we focus on Movie question answering (MovieQA) [26], which requires systems to demonstrate story comprehension by successfully answering multiple choice questions relating to videos and subtitles taken from movies.

A key challenge in multi-modal QA is to integrate information from different data sources. In the context of MovieQA, both query-to-context attention and inter-modal attention between videos and subtitles should be considered. Recently

Anran Wang, Anh Tuan Luu, Chuan-Sheng Foo, Hongyuan Zhu, and Vijay Chandrasekhar are with Institute for Infocomm Research, A*STAR, Singapore 138632 (e-mail: wang_anran@i2r.a-star.edu.sg; at.luu@i2r.a-star.edu.sg; foo_chuan_sheng@i2r.a-star.edu.sg; zhuh@i2r.a-star.edu.sg; vijay@i2r.a-star.edu.sg).

Yi Tay is with the School of Computer Science and Engineering, Nanyang Technological University, Singapore 639798 (e-mail: YTAY017@e.ntu.edu.sg).

developed methods have adopted the classic strategies of early-fusion [19] and late-fusion [26], both of which have their limitations. Early-fusion of different modalities may limit the ability to pick up meaningful semantic correlations due to the increased noise at the feature level, while late-fusion does not allow for cross-referencing between modalities to define the higher level semantic features. Wang *et al.* [28] proposed to utilize inter-modal attention for MovieQA. However, their method does not fully integrate the input data, where different attention stage considers a different subset of interactions between the question, videos, subtitles for context retrieval.

Moreover, answer choices are only considered at the final step of the system where they are matched against an integrated representation of the input data, hence the useful contexts among answer choices are not effectively utilized to determine the relevant parts of the input data.

To address these limitations, we propose the Holistic Multi-modal Memory Network (HMMN) framework. Firstly, our framework employs both inter-modal and query-to-context attention mechanisms for effective data integration in each hop. Specifically, our attention mechanism holistically investigates videos, subtitles, question, answer choices to obtain a summarized context in each attention hop, which is different from existing methods that only consider a subset of interactions in each hop. Hence, query-to-context relationship is jointly considered while modeling the multi-modal relationship between context. Secondly, our framework considers answer choices during not only answer prediction stage but also context retrieval stage. Utilizing answer choices to hone in on relevant information is a common heuristic used by students when taking multiple-choice tests. Analogously, we thought this would help on the QA task by restricting the set of inputs considered by the model thus helping it sieve out signal from the noise. Specifically, the retrieved answer-aware context should match the answer choice for correct answers. Otherwise, the resultant context may convey different semantic meaning from the answer choice.

Given this holistically modeling, our HMMN framework achieves state-of-the-art performance on both validation and test sets for video-based movie question answering. Ablation studies confirm the utility of incorporating answer choices during the context retrieval stage, which show that incorporating the answer information contributes to the context retrieval process. In addition, we include analyses of different attention mechanisms (query-to-context attention, inter-modal attention, intra-modal self attention) and confirm the importance of holistic reasoning.

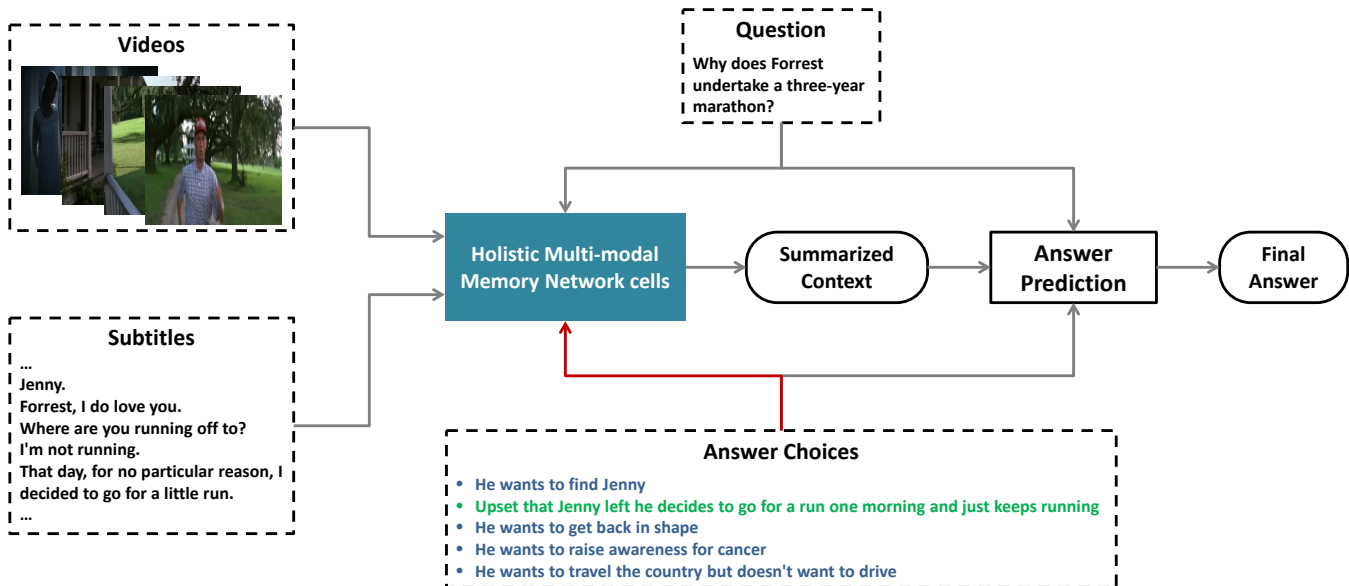The rest of this paper is organized as follows. Section II

Fig. 1. Illustration of our proposed multi-modal feature learning framework. Our Holistic Multi-modal Memory Network cells holistically fuse multi-modal context (videos, subtitles), the question, as well as the answer choices. This framework jointly takes inter-modal and query-to-context attentions into consideration in each attention hop, and incorporates answer choices in both context retrieval and answer prediction stages.

introduces related works on multi-modal question answering and visual question answering as well. Section III describes the HMMN method in detail. Section IV presents the experimental results together with analyses of different attention mechanisms. Section V concludes the paper.

## II. RELATED WORK

In this section, methods of visual question answering and multi-modal question answering which are closely related to our method are introduced.

### A. Visual Question Answering

Besides the significant progress from computer vision and natural language processing, the emergence of large visual question answering (VQA) datasets [2], [5], [11] also leads to the popularity of VQA task. Early works about VQA task use the holistic full-image feature to represent the visual context. For example, Malinowski *et al.* [18] proposed to feed the Convolutional Neural Network (CNN) image features and the question features together into a long short-term memory (LSTM) network and train an end-to-end network. Later, quite a few works have used attention mechanism to pay attention to certain parts of the image, where the alignment between image patches [32], [36] or region proposals [23] with the words in the question have been explored. Several attention mechanisms of connecting the visual context and the question have been proposed [31], [16]. For example, Lu *et al.* [16] presented a co-attention framework which considers visual attention and question attention jointly. Wu *et al.* [30] proposed to incorporate high-level concepts and external knowledge for image captioning and visual question answering. Answer

attention was investigated for the grounded question answering task [6]. Grounded question answering is a special type of VQA, which is to retrieve an image bounding box from the candidate pool to answer the textual question. This method models the interaction between the answer candidates and the question and learns the answer-aware summarization of the question, while our method models the interaction between the answer choices and the context to retrieve more informative context.

### B. Multi-modal Question Answering

In contrast to VQA, which only involves visual context, multi-modal question answering takes multiple modalities as context, and has attracted great interest. Kembhavi *et al.* [12] presented the Textbook Question Answering (TextbookQA) dataset that consists of lessons from middle school science curricula with both textual and diagrammatic context. In [13], PororoQA dataset was introduced, which is constructed from children cartoon Pororo with video, dialogue, and description. Tapaswi *et al.* [26] introduced the movie question answering (MovieQA) dataset which aims to evaluate the story understanding from both video and subtitle modalities. In this paper, we focus on the MovieQA dataset, and related approaches are discussed as follows.

Most methods proposed for the MovieQA dataset are based on the End-to-end Memory Network [25], which was originally proposed for the textual question answering task. In [26], they proposed a straightforward extension of the End-to-end Memory Network [25] to multi-modal scenario on MovieQA. In particular, answer is predicted based on each modality separately, and late fusion is performed to combine the answer prediction scores from two modalities.

Na *et al.* [19] proposed another framework based on the End-to-end Memory Network [25]. Their framework has read and write networks that are implemented by convolutional layers to model sequential memory slots as chunks. Context from textual and visual modalities are early fused as the input to the write network. Specifically, compact bilinear pooling [3] is utilized to obtain the joint embeddings with subshot and sentence features. Deep embedded memory networks (DEMN) model was introduced by [13], where they aim to reconstruct stories from a joint stream of scene and dialogue. However, their framework is not end-to-end trainable, as their method makes a hard context selection. Liang *et al.* [14] presented a focal visual-text attention network which captures correlation between visual and textual sequences for the personal photos and descriptions in the MemexQA dataset [10] and applied this method to the MovieQA dataset. Wang *et al.* [28] proposed a layered memory network with two-level correspondences. Specifically, the static word memory module corresponds words with regions inside frames, and the dynamic subtitle memory module corresponds sentences with frames. However, visual modality is used to attend to the textual modality which is dynamically updated by different strategies. Interactions between the question, videos, subtitles are not holistically considered in each attention stage.

## III. METHODOLOGY

We will first introduce the notations. Then, we will introduce the End-to-end Memory Network (E2EMN) [25]. After that, Holistic Multi-modal Memory Network (HMMN) cell will be introduced together with the prediction framework.

### A. Notations

We assume the same feature dimension for subtitle sentences and frames. In MovieQA dataset, each question is aligned with several relevant video clips. We obtain features for frames and sentences following [28].

Let $S \in \mathbb{R}^{d \times m}$ denote the subtitle modality, where $d$ is the dimension of the feature vectors and $m$ is number of subtitle sentences. For the subtitle modality, we not only gather the subtitle sentences within the relevant video clips, but also incorporate nearby (in time) subtitle sentences to make use of the contextual information. The dimension of the word2vec features for words in subtitles is $d_w$. The word2vec representation of each word is projected to $d$-dim with a projection matrix $W_1 \in \mathbb{R}^{d_w \times d}$. Then, a mean-pooling is performed among all words in each sentence to get the sentence representation.

Similarly, $V \in \mathbb{R}^{d \times n}$ represents the video modality, where $d$ is the dimension of the feature vectors and $n$ is number of frames. We select a fixed number of frames from the relevant video clips for each question. Frame-level representations are generated by investigating attention between regional features and word representations in the vocabulary, where $W_2 \in \mathbb{R}^{d_r \times d_w}$ is utilized to project the regional VGG [24] features to $d_w$-dim to match the dimension of word representations. Here $d_r$ is the dimension of the regional features. With regional features represented by the vocabulary word features, frame-level representations are generated with average pooling followed by a projection with $W_1$. We refer the reader to the original paper [28] or their released code for more details.

The question and answer choices are represented in the same way as subtitle sentences. The question is represented as a vector $q \in \mathbb{R}^d$. Answer choices for each question are represented as $A = [a_0, a_1, a_2, a_3, a_4] \in \mathbb{R}^{d \times 5}$, where each answer choice is encoded as $a_k \in \mathbb{R}^d$. In the whole framework, only $W_1$ and $W_2$ are learnable. The structures to generate representations for the subtitle and video modalities are shown in Fig. 2(a).

### B. End-to-end Memory Network

The End-to-end Memory Network (E2EMN) [25] is originally proposed for a question answering task where the aim is to pick the most likely word from the vocabulary as answer according the textual context. In [26], E2EMN is adapted to multi-modal question answering with multi-choice answers. In particular, scores from two modalities are late-fused to make the final prediction. As E2EMN is designed for textual question answering, this method only deals with context from a single modality. Here we use the subtitle modality $S$ for explanation.

In E2EMN, input features of context $S$ are treated as memory slots. With both memory slots and query (question is used as query here) as input, a summary of context is derived according to the relevance between the query and memory slots. In particular, the match between query $q$ and each memory slot is calculated with the inner product followed by a softmax:

$$\alpha_i = softmax(q^T S_{:i}) \qquad (1)$$

where $\alpha_i$ indicates the importance of $i$-th subtitle sentence to the query. The summarized context $u$ is computed as the weighted sum of subtitle sentence features based on $\alpha_i$:

$$u = \sum_{i=1}^{m} \alpha_i S_{:i} \qquad (2)$$

Then, the answer prediction is made by comparing the answer choice $a_i$ with the sum of query representation $q$ and the summarized context $u$:

$$p = softmax((q+u)^T A) \qquad (3)$$

where $p \in \mathbb{R}^5$ is the confidence vector. Here 5 is the number of answer choices for MovieQA. The process to derive the summarized context can be performed in multiple hops, where the output of one layer can be used as part of the query of the next layer.

### C. Holistic Multi-modal Memory Network (HMMN)

Different from E2EMN, our HMMN framework takes multi-modal context as input. HMMN framework investigates interactions between multi-modal context and the question jointly. By doing this, query-to-context relationship is jointly considered while modeling the multi-modal relationship between context. In addition, it not only exploits answer choices for answer prediction but also in the process of summarizing the context from multiple modalities.
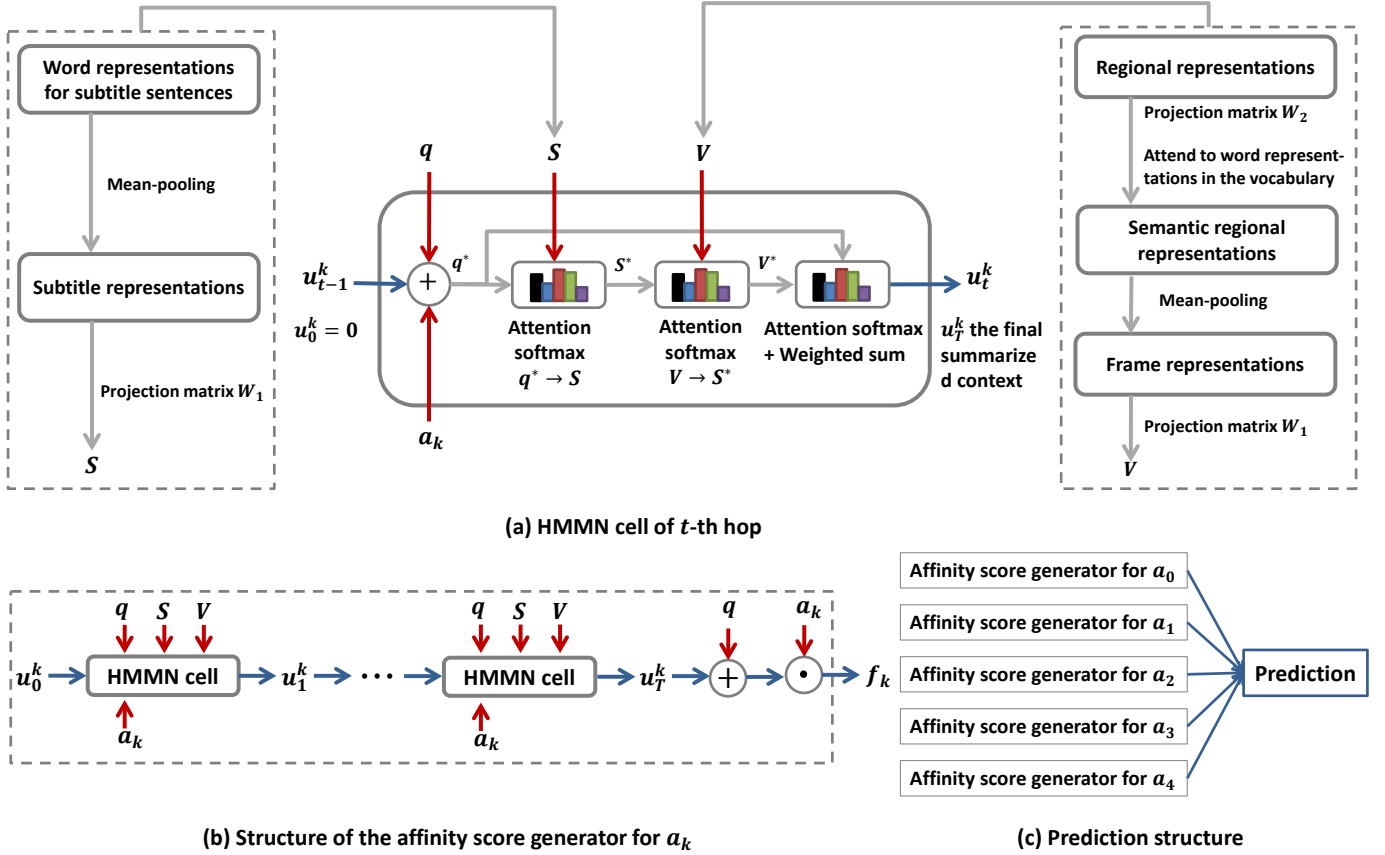
Fig. 2. Illustration of (a) the HMMN cell of $t$-th hop, (b) affinity score generator for answer choice $a_k$, and (c) the prediction structure. In (b), $T$ is the number of hops, which also denotes the number of stacked HMMN cells. $\oplus$ denotes the element-wise addition and $\odot$ denotes the element-wise multiplication.

The inference process is performed by stacking small building blocks, called HMMN cell. The structure of HMMN cell is shown in Fig. 2(a). The process to generate $S$ and $V$ is only for illustration, while our main contribution lies in the HMMN structure. Each HMMN cell takes as input the question, one answer choice, context from videos and subtitles, and derives the answer-aware summarized context. We call this process as one hop of reasoning. Let $u_t^k$ be the output of the $t$-th reasoning hop with respect to answer choice $k$. The output of the $t$-th hop will be utilized as the input of the $(t+1)$-th hop.

*1) Involving answers in context retrieval:* The HMMN cell incorporates the answer choice as a part of the query in the context retrieval stage. The query involving $k$-th answer choice for the $t$-th hop is calculated by combining the output of previous hop $u_{t-1}^k$, the question $q$, answer choice $a_k$:

$$q^* = u_{t-1}^k + a_k + \lambda q \qquad (4)$$

where $\lambda$ is a tradeoff parameter between the question and the rest of the query.

The intuition of incorporating answer choices in the context retrieval stage is to mimic behaviors of students who take a reading test with multi-choice questions. When the context is long and complicated, a quick and effective way to answer the question is to locate relevant information with respect to

each answer choice. For one answer choice, if the retrieved answer-aware context conveys similar ideas with the answer choice, it tends to be the correct answer. Alternatively, if the retrieved context has a different semantic meaning, the answer is likely to be wrong.

*2) Holistically considering different attention mechanisms in each hop:* Instead of only taking a subset of interactions between the query and multi-modal context, our framework jointly considers inter-modal and query-to-context attention strategies in each hop.

The HMMN cell takes the query $q^*$ to gather descriptive information from multi-modal context, where interactions between the question, answer choices, videos, subtitles are exploited holistically. In particular, we utilize the updated query to highlight the relevant subtitle sentences in $S$ by performing the query-to-context attention (denoted as $(q^* \to S)$). The resulted re-weighted subtitle modality is represented as $S^*$:

$$\delta_i = softmax(q^{*T} S_{:i})$$
$$S_{:i}^* = \delta_i S_{:i} \qquad (5)$$

where more relevant subtitle sentences are associated with large weights.

Inter-modal attention reasoning is applied by using the video modality $V$ to attend to the subtitle modality $S$ (denoted as $(V \rightarrow S^*)$), which aims to generate the subtitle-aware representations for frames as $V^*$. Each frame is represented with the weighted sum of all the subtitle sentence features according to the relevance:

$$\varepsilon_{ij} = V_{:i}^T S_{:j}^*$$
$$V_{:i}^* = \sum_{j=1}^{m} \varepsilon_{ij} S_{:j}^* \qquad (6)$$

The resulted $V^*$ can be summarized with respect to the query $q^*$ as the hop output. In particular, the $t$-hop summarized context with respect to the $k$-th answer choice is caculated as $u_t^k$:

$$\zeta_i = softmax(q^{*T} V_{:i}^*)$$
$$u_t^k = \sum_{i=1}^{n} \zeta_i V_{:i}^* \qquad (7)$$

In each reasoning hop, the output of previous hop, the answer choice, the question, multi-modal context are holistically integrated. The reason of using $V$ to attend to $S$ (not using $S$ to attend to $V$) is that, the subtitle modality is more informative than the video modality for the MovieQA task. Typically, the subtitle modality includes descriptions of the story such as character relationships, story development. By attending to $S$, the feature representations in $S$ will be used to form the summarized context. The re-weighted $S$ acts as an information screening step to derive more informative representations for the subtitle modality.

*3) Predicting answer with affinity scores:* It is shown in the original E2EMN that multiple hops setting yields improved results. We stack the HMMN cells to do $T$ hops of reasoning. Given the final summarized context $u_T^k$ with respect to answer choice $a_k$, an affinity score $f_k$ for $a_k$ is generated. This score is derived by comparing the sum of the question and answer-aware summarized context with the answer choice as:

$$f_k = (q + u_T^k)^T a_k \qquad (8)$$

This score indicates whether the retrieved context has the consistent semantic meaning with the answer choice. The structure of generating the affinity score is shown in Fig. 2(b). Then the affinity scores for all the answer choices will be passed to a softmax function to get the final answer prediction as shown in Fig. 2(c), and the cross-entropy loss is minimized with the standard stochastic gradient descent. This indicates that if one answer choice matches with the answer-aware summarized context, it is likely to be the correct answer.

## IV. EXPERIMENTS

### A. Dataset and Setting

The MovieQA dataset [26] consists of 408 movies and 14,944 questions. Diverse sources of information are collected including video clips, plots, subtitles, scripts, and Descriptive Video Service (DVS). Plot synopses from Wikipedia are utilized to generate questions. For multi-modal question answering task with videos and subtitles, there are 6,462 questions with both videos clips and subtitles from 140 movies. We

TABLE I
PERFORMANCES OF HMMN STRUCTURES W/ AND W/O ANSWER ATTENTION WITH DIFFERENT NUMBERS OF LAYERS ON THE VALIDATION SET.

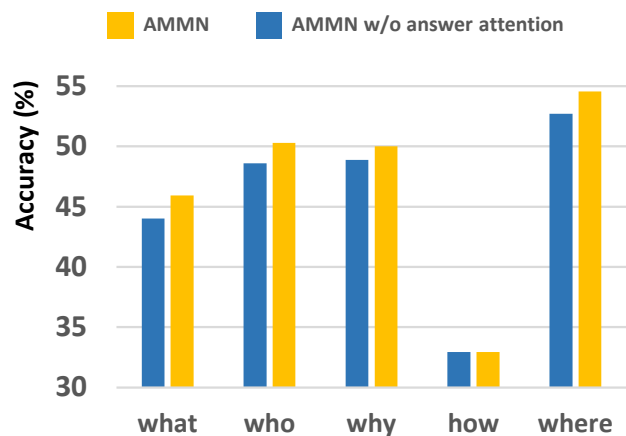| Method | Accuracy (%) |
|---|---|
| HMMN (1 layer)  w/o answer attention | 43.35 |
| HMMN (2 layers) w/o answer attention | **44.47** |
| HMMN (3 layers) w/o answer attention | 44.24 |
| HMMN (1 layer) | 45.71 |
| HMMN (2 layers) | **46.28** |
| HMMN (3 layers) | 44.13 |



Fig. 3. Comparison between HMMN with and without answer attention for different question types.

follow the public available train, validation, test split. The accuracy of multi-choice questions is measured.

### B. Implementation Details

For the subtitle modality, we consider the subtitle sentences which fall into the time interval derived by extending the starting and ending time points of video clips by 300 seconds. For the video modality, 32 frames are selected from the relevant video clips following [28]. We use the word2vec representations provided by [26]. The dimension of the word2vec representation $d_w$ is 300. The dimension of the regional features from 'pool5' of VGG-16 $d_r$ is 512. 10% of the training examples are kept as the development set. The batch size is 8. The learning rate is set to 0.005. The tradeoff parameter $\lambda$ is set to be 0.45. The dimension of features $d$ is set to 300. Our model is trained up to 50 epochs, and early stopping is performed.

### C. Quantitative Analysis

Table I presents results for HMMN structures with and without answer attention with different numbers of layers, where "with answer attention" means considering answer choices when retrieving the context. We can see that by

TABLE II
COMPARISON OF STATE-OF-THE-ART METHODS ON VALIDATION AND TEST SETS.

| Method | Accuracy on Val (%) | Accuracy on Test (%) |
|---|---|---|
| Tapaswi *et al.* [26] | 34.20 | - |
| Na *et al.* [19] | 38.67 | 36.25 |
| Kim *et al.* [13] | 44.7 | 34.74 |
| Liang *et al.* [14] | 41.0 | 37.3 |
| Wang *et al.* [28] | 42.5 | 39.03 |
| Our HMMN framework w/o answer attention | 44.47 | 41.65 |
| Our HMMN framework | **46.28** | **43.08** |

TABLE III
PERFORMANCE OF BASELINES WITH DIFFERENT ATTENTION STRATEGIES ON THE VALIDATION SET.

| Method | Accuracy (%) |
|---|---|
| $V$ | 37.69 |
| $S$ | 39.62 |
| $V'$ $(q \to V)$ Query-to-context Attention | 37.92 |
| $S'$ $(q \to S)$ Query-to-context Attention | 40.86 |
| $\bar{V}$ $(V \to S)$ Inter-modal Attention | 42.73 |
| $\bar{S}$ $(S \to V)$ Inter-modal Attention | 35.10 |
| $\hat{V}$ $(V \to V)$ Intra-modal Self Attention | 37.92 |
| $\hat{S}$ $(S \to S)$ Intra-modal Self Attention | 40.29 |

TABLE IV
PERFORMANCE OF BASELINES EXPLORING HIGHER-LEVEL INTER-MODAL ATTENTION ON THE VALIDATION SET.

| Method | Accuracy (%) | Method | Accuracy (%) |
|---|---|---|---|
| $V \to S$ | 42.73 | $S \to V$ | 35.10 |
| $V \to S'$ | **43.35** | $S' \to V$ | 35.21 |
| $V \to \bar{S}$ | 37.47 | $\bar{S} \to V$ | 35.10 |
| $V \to \hat{S}$ | 41.08 | $\hat{S} \to V$ | 35.10 |
| $V' \to S$ | 43.12 | $S \to V'$ | 35.21 |
| $V' \to S'$ | **43.35** | $S' \to V'$ | 35.44 |
| $V' \to \bar{S}$ | 38.14 | $\bar{S} \to V'$ | 35.32 |
| $V' \to \hat{S}$ | 37.02 | $\hat{S} \to V'$ | 35.44 |
| $\bar{V} \to S$ | 41.76 | $S \to \bar{V}$ | 40.29 |
| $\bar{V} \to S'$ | 41.08 | $S' \to \bar{V}$ | 40.18 |
| $\bar{V} \to \bar{S}$ | 39.84 | $\bar{S} \to \bar{V}$ | 40.85 |
| $\bar{V} \to \hat{S}$ | 37.47 | $\hat{S} \to \bar{V}$ | 38.60 |
| $\hat{V} \to S$ | 43.12 | $S \to \hat{V}$ | 34.55 |
| $\hat{V} \to S'$ | **43.35** | $S' \to \hat{V}$ | 35.77 |
| $\hat{V} \to \bar{S}$ | 37.58 | $\bar{S} \to \hat{V}$ | 35.10 |
| $\hat{V} \to \hat{S}$ | 38.15 | $\hat{S} \to \hat{V}$ | 34.98 |

incorporating the answer choices in the context retrieval stage, the performance is significantly improved. And 2-layer structures achieve the best performance, thus the number of layers which is also the number of hops $T$ is set to 2. Fig. 3 shows the the comparison of HMMN w/ and w/o answer attention for different question types, with the starting word as 'what', 'who', 'why', 'who', 'where'. It can be seen that HMMN framework performs consistently better than the HMMN framework w/o answer attention.

Table II shows the comparison with state-of-the-art methods. We compare with [26], [19], [13], [14], [28]. It can be seen that our proposed method significantly outperforms all four state-of-the-art methods on both validation and test sets.

In particular, [26] performs a late fusion. [19] conducts an early fusion. Neither late fusion nor early fusion can well exploit the relationship between modalities, which results in suboptimal results. [13] is not end-to-end trainable. [14] assumes sequences of the visual and textual representations have the same length, which is not true for MovieQA. Similar with [19], [14] cuts off visual information of those frames without accompanying subtitles. [28] takes the multi-modal relationship into consideration, however, in each attention stage, a different subset of interactions between the question, videos, subtitles are considered. In comparison, our HMMN framework w/o answer attention holistically incorporates the output of previous hop, the question, videos, subtitles in each hop, which leads to superior performance.

*1) Ablation study:* In this paper, we also explore different attention mechanisms. Effects of different attention strategies are investigated.

**(i) Query-to-context Attention** $(q \to S)$ $S'$

Query-to-context attention indicates which memory slots are more relevant to the query. Here $S$ of subtitle modality is used as the context for illustration. With the calculated similarity between the query and each memory slot in Eq. 1, more relevant subtitle sentences can be highlighted with:

$$S'_{:i} = \alpha_i S_{:i} \tag{9}$$

This process is denoted as $q \to S$, and the re-weighted memory slots are represented as $S'$.

**(ii) Inter-modal Attention** $(S \to V)$ $\bar{S}$

Inter-modal attention from $S$ to $V$ indicates that, for each subtitle sentence, we intent to find the most relevant frames. The retrieved frame features will be fused to represent the subtitle sentence. The output can be interpreted as the video-aware sentence representation.

First, the coattention matrix between frames and subtitle sentences can be defined as:

$$\beta_{ij} = S_{:i}^T V_{:j} \tag{10}$$

where $\beta_{ij}$ indicates the relevance between $j$-th frame and $i$-th subtitle sentence. Then $i$-th subtitle sentence can be represented by the weighted sum of all frames based on $\beta_{ij}$:

$$\bar{S}_{:i} = \sum_{j=1}^{n} \beta_{ij} V_{:j} \tag{11}$$

The resulted representation for subtitle modality $\bar{S}$ is of the same size as $S$. Similarly, the result for $(V \to S)$ is $\bar{V}$.

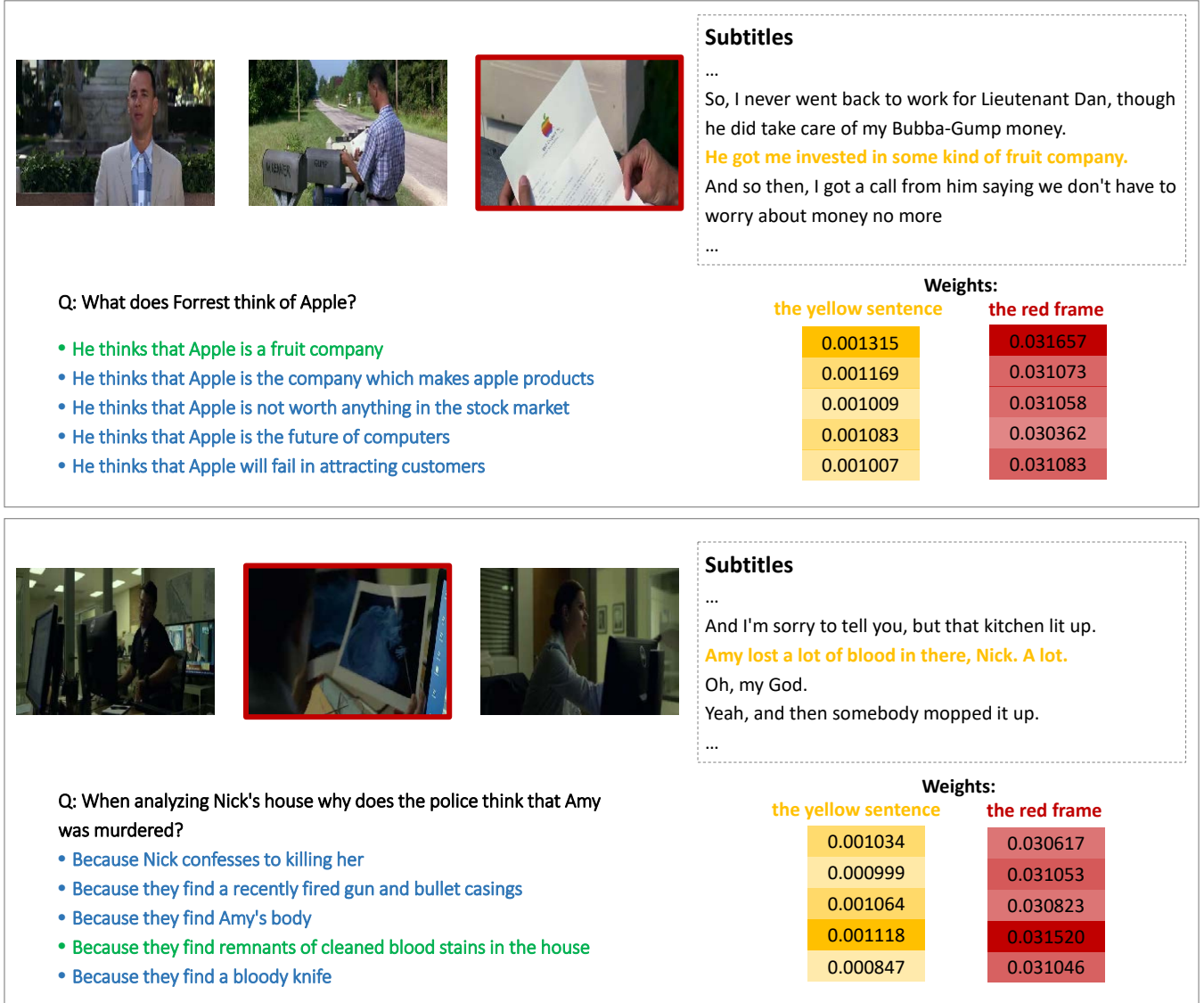**(iii) Intra-modal Self Attention** $(S \to S)$ $\hat{S}$

Fig. 4. Visualization of weights of two success cases in the HMMN framework. The correct answer choice is colored in green. Weights of relevant frame and subtitle sentence (highlighted in red and yellow) are shown. When using the correct answer choice to search for the relevant context, the relevant frame and sentence are associated with higher weights than those for other answer choices. The correct answer choice is in green.

Self attention has shown its power in tasks such as question answering and machine translation [8], [27]. The intuition is that contextual information among other memory slots can be exploited. Similar with the inter-modal attention, the coattention between different memory slots in the same modality is calculated as:
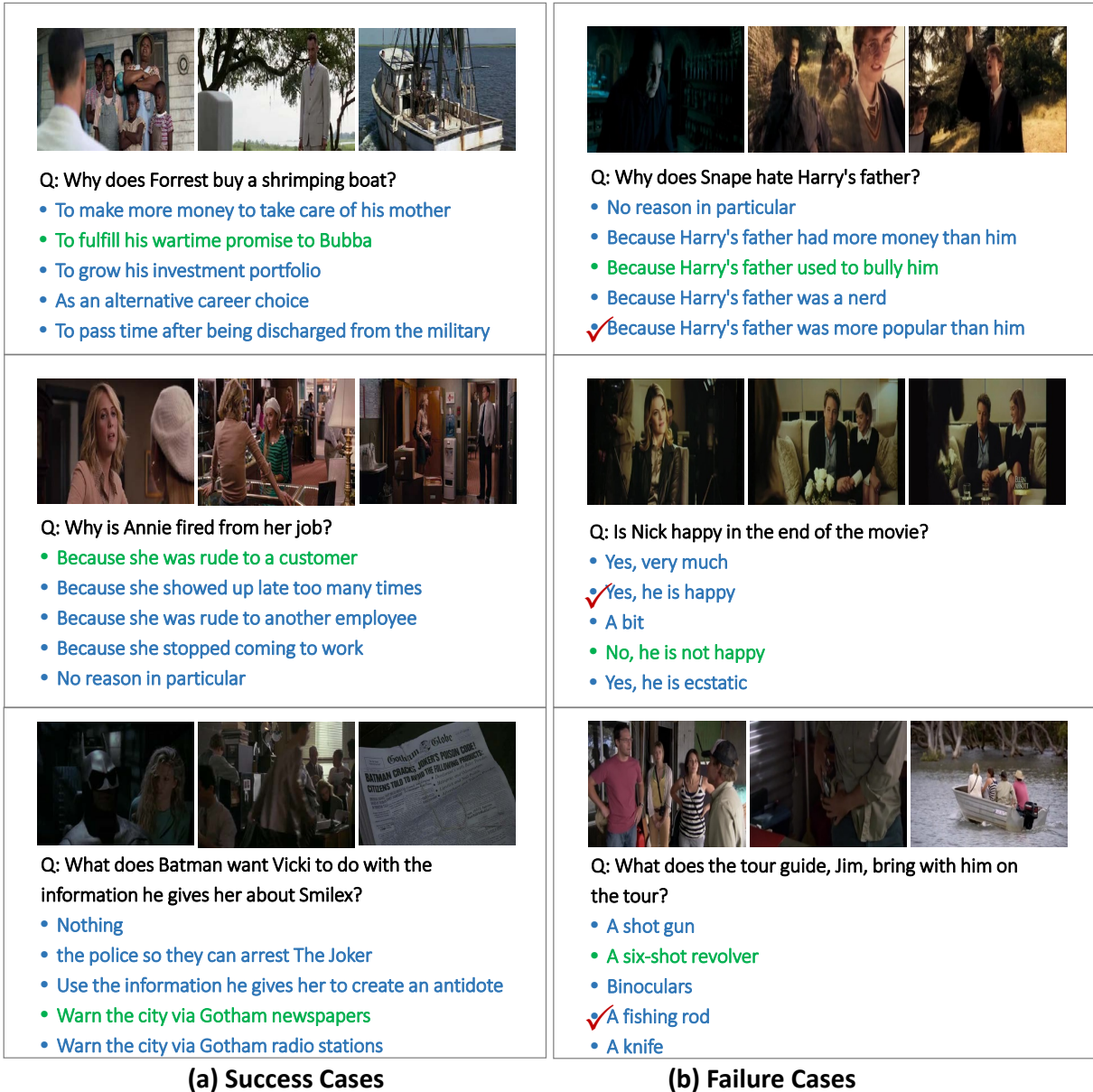
$$\gamma_{ij} = I(i \neq j) S_{:i}^T S_{:j} \qquad (12)$$

Noted that the correlation between one sentence with itself is set to zero. The resulted representation with self-attention is $\hat{S}$. Each subtitle sentence will be represented by the weighted sum of features of all the subtitle sentences based on $\gamma_{ij}$:

$$\hat{S}_{:i} = \sum_{j=1}^{m} \gamma_{ij} S_{:j} \qquad (13)$$

With derived $S$, $S'$, $\bar{S}$, $\hat{S}$, $V$, $V'$, $\bar{V}$, $\hat{V}$, each of them can be treated as memory slots in the original E2EMN. Similarly, Eq. 1, Eq. 2 and Eq. 3 can be applied to predict the answer.

Table III shows results of baselines with different attention strategies. For $S$, $S'$, $\bar{S}$, $\hat{S}$, $V$, $V'$, $\bar{V}$, $\hat{V}$, each baseline treats one of them as memory slots in the original E2EMN to predict the answer. The baseline $S$ performs better than $V$, as the subtitle modality contains more informative descriptions of character relationships and story development. By using each attention strategy to obtain the enhanced representations, the performance improvement is achieved. Particularly, the inter-modal attention ($V \rightarrow S$) brings a significant improvement.

To explore higher-level inter-modal attention, we use $V$, $V'$, $\bar{V}$, $\hat{V}$ to attend to $S$, $S'$, $\bar{S}$, $\hat{S}$, and vise versa with Eq. 10 and Eq. 11. The results are shown in Table IV. Typically, the baselines on the left side perform better than ones on

Q: Why does Forrest buy a shrimping boat?
• To make more money to take care of his mother
• To fulfill his wartime promise to Bubba
• To grow his investment portfolio
• As an alternative career choice
• To pass time after being discharged from the military

Q: Why is Annie fired from her job?
• Because she was rude to a customer
• Because she showed up late too many times
• Because she was rude to another employee
• Because she stopped coming to work
• No reason in particular

Q: What does Batman want Vicki to do with the information he gives her about Smilex?
• Nothing
• the police so they can arrest The Joker
• Use the information he gives her to create an antidote
• Warn the city via Gotham newspapers
• Warn the city via Gotham radio stations

**(a) Success Cases**

Q: Why does Snape hate Harry's father?
• No reason in particular
• Because Harry's father had more money than him
• Because Harry's father used to bully him
• Because Harry's father was a nerd
✓Because Harry's father was more popular than him

Q: Is Nick happy in the end of the movie?
• Yes, very much
✓Yes, he is happy
• A bit
• No, he is not happy
• Yes, he is ecstatic

Q: What does the tour guide, Jim, bring with him on the tour?
• A shot gun
• A six-shot revolver
• Binoculars
✓A fishing rod
• A knife

**(b) Failure Cases**

Fig. 5. Success and failure cases. The correct answer choice is in green and the mistake made by our framework is marked in red. Failure cases indicate the necessity of exploiting common sense knowledge and object detection results.

the right side. As mentioned in the methodology section, it is because when we use $V \to S$ attention, the video modality will be represented by the subtitle features, which are more descriptive. According to this observation, when designing the structure of attention mechanism for any general multi-modal tasks, using the less discriminative modalities as clues to attend to the more discriminative modalities tends to achieve better performance. We can observe that the intra-modal self attention does not bring much improvement to this task. Baselines of using the video modality to attend to re-weighted subtitle modality ($V \to S'$) perform considerably well, and our 1-layer HMMN framework w/o answer attention degenerates to the $V \to S'$ baseline.

*D. Qualitative Analysis*

To demonstrate that relevant context can be well captured by the HMMN framework, we visualize the attention weights of frames and subtitles of two success cases in Fig. 4(a). The observed relevant frame and subtitle are highlighted in red and yellow respectively. For the question "What does Forrest think of Apple?", the first answer choice is correct. Following attention weights of the second layer are visualized: 1) the attention weight of the subtitle sentence with respect to the answer choice and question; 2) the attention weight of the frame with respect to the answer choice and question. By using the first answer choice to retrieve the context, the relevant frame and subtitle sentence are associated with larger weights compared to those of other choices. Thus the key information can be captured and a high affinity score will be generated. On

the other hand, the HMMN framework w/o answer attention picks the second answer choice as the correct one. This is because that the word "Apple" is not mentioned in the subtitle, thus by using the question only to retrieve the information, important context about "Apple" is missed. Similar remarks can be made for the second example question in Fig. 4(a), the relevant frame and subtitle sentence are given high weights.

Fig. 4(b) shows 3 typical failure cases. In the first example, although both video and subtitles contain the information that Harry's father made fun of Snape, it is difficult to associate them with the word "bully" that has a high-level semantic meaning, where common sense knowledge reasoning is required. The second example is from "Gone girl". Although the husband is not happy with the main character Amy, the actors act as having a happy ending. This question also requires common sense to answer. In the third example, the revolver appears in the frames, but is not mentioned in the subtitles. Although we use conventional CNN features to generate frame-level representations, the associations between visual patterns and object labels are not enforced during training. More success and failure cases can be found in Fig. 5.

**Future work**: Failure cases indicate the necessity of exploiting common sense knowledge and object detection results, which we leave for future work. High-level semantic information can be injected to the network by leveraging well-built knowledge graph, e.g. ConceptNet [15]. Object detector can capture more descriptive visual information than the current grid-based regions, and labels generated along with detected regions further enforce the correspondence of visual and textual information.

## V. CONCLUSION

We presented a Holistic Multi-modal Memory Network framework that learns to answer questions with context from multi-modal data. In the proposed HMMN framework, we investigate both inter-modal and query-to-context attention mechanism to jointly model the interactions between multi-modal context and the question. In addition, we explore the answer attention by incorporating the answer choices in the context retrieval stage. Our HMMN framework achieves state-of-the-art results on the MovieQA dataset. We also presented a detailed ablation study for different attention mechanisms, which could provide guidance for future model design.

## REFERENCES

[1] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. Lawrence Zitnick, and D. Parikh, "Vqa: Visual question answering," in *ICCV*, 2015.

[2] S. Antol, A. Agrawal, J. Lu, M. Mitchell, D. Batra, C. L. Zitnick, and D. Parikh, "VQA: Visual Question Answering," in *ICCV*, 2015.

[3] Y. Gao, O. Beijbom, N. Zhang, and T. Darrell, "Compact bilinear pooling," in *CVPR*, 2016.

[4] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh, "Making the v in vqa matter: Elevating the role of image understanding in visual question answering," in *CVPR*, 2017.

[5] Y. Goyal, T. Khot, D. Summers-Stay, D. Batra, and D. Parikh, "Making the V in VQA matter: Elevating the role of image understanding in Visual Question Answering," in *CVPR*, 2017.

[6] J. Hu, D. Fan, S. Yao, and J. Oh, "Answer-aware attention on grounded question answering in images," 2017.

[7] Y. Huang, W. Wang, and L. Wang, "Instance-aware image and sentence matching with selective multimodal lstm," in *CVPR*, 2017.

[8] D. A. Hudson and C. D. Manning, "Compositional attention networks for machine reasoning," *arXiv preprint arXiv:1803.03067*, 2018.

[9] U. Jain, S. Lazebnik, and A. Schwing, "Two can play this game: Visual dialog with discriminative question generation and answering," in *CVPR*, 2018.

[10] L. Jiang, J. Liang, L. Cao, Y. Kalantidis, S. Farfade, and A. Hauptmann, "Memexqa: Visual memex question answering," *arXiv preprint arXiv:1708.01336*, 2017.

[11] J. Johnson, B. Hariharan, L. van der Maaten, L. Fei-Fei, C. L. Zitnick, and R. Girshick, "Clevr: A diagnostic dataset for compositional language and elementary visual reasoning," in *CVPR*, 2017.

[12] A. Kembhavi, M. Seo, D. Schwenk, J. Choi, A. Farhadi, and H. Hajishirzi, "Are you smarter than a sixth grader? textbook question answering for multimodal machine comprehension," in *CVPR*, 2017.

[13] K.-M. Kim, M.-O. Heo, S.-H. Choi, and B.-T. Zhang, "Deepstory: video story qa by deep embedded memory networks," in *IJCAI*, 2017.

[14] J. Liang, L. Jiang, L. Cao, L.-J. Li, and A. Hauptmann, "Focal visual-text attention for visual question answering," in *CVPR*, 2018.

[15] H. Liu and P. Singh, "Conceptnet?a practical commonsense reasoning tool-kit," *BT technology journal*, vol. 22, no. 4, pp. 211–226, 2004.

[16] J. Lu, J. Yang, D. Batra, and D. Parikh, "Hierarchical question-image co-attention for visual question answering," in *NIPS*, 2016.

[17] L. Ma, Z. Lu, L. Shang, and H. Li, "Multimodal convolutional neural networks for matching image and sentence," in *ICCV*, 2015.

[18] M. Malinowski, M. Rohrbach, and M. Fritz, "Ask your neurons: A neural-based approach to answering questions about images," in *ICCV*, 2015.

[19] S. Na, S. Lee, J. Kim, and G. Kim, "A read-write memory network for movie story understanding," in *CVPR*, 2017.

[20] F. Nian, T. Li, Y. Wang, X. Wu, B. Ni, and C. Xu, "Learning explicit video attributes from mid-level representation for video captioning," *Computer Vision and Image Understanding*, vol. 163, pp. 126–138, 2017.

[21] P. Pan, Z. Xu, Y. Yang, F. Wu, and Y. Zhuang, "Hierarchical recurrent neural encoder for video representation with application to captioning," in *CVPR*, 2016.

[22] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel, "Self-critical sequence training for image captioning," in *CVPR*, 2017.

[23] K. J. Shih, S. Singh, and D. Hoiem, "Where to look: Focus regions for visual question answering," in *CVPR*, 2016.

[24] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2015.

[25] S. Sukhbaatar, J. Weston, R. Fergus *et al.*, "End-to-end memory networks," in *NIPS*, 2015.

[26] M. Tapaswi, Y. Zhu, R. Stiefelhagen, A. Torralba, R. Urtasun, and S. Fidler, "MovieQA: Understanding Stories in Movies through Question-Answering," in *CVPR*, 2016.

[27] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, and I. Polosukhin, "Attention is all you need," in *NIPS*, 2017.

[28] B. Wang, Y. Xu, Y. Han, and R. Hong, "Movie question answering: Remembering the textual cues for layered visual contents," in *AAAI*, 2018.

[29] L. Wang, Y. Li, J. Huang, and S. Lazebnik, "Learning two-branch neural networks for image-text matching tasks," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018.

[30] Q. Wu, C. Shen, P. Wang, A. Dick, and A. van den Hengel, "Image captioning and visual question answering based on attributes and external knowledge," *IEEE transactions on pattern analysis and machine intelligence*, vol. 40, no. 6, pp. 1367–1381, 2018.

[31] C. Xiong, S. Merity, and R. Socher, "Dynamic memory networks for visual and textual question answering," in *ICML*, 2016.

[32] H. Xu and K. Saenko, "Ask, attend and answer: Exploring question-guided spatial attention for visual question answering," in *ECCV*, 2016.

[33] Q. You, H. Jin, Z. Wang, C. Fang, and J. Luo, "Image captioning with semantic attention," in *CVPR*, 2016.

[34] M. Zhang, Y. Yang, H. Zhang, Y. Ji, H. T. Shen, and T.-S. Chua, "More is better: Precise and detailed image captioning using online positive recall and missing concepts mining," *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 32–44, 2019.

[35] L. Zhou, Y. Zhou, J. J. Corso, R. Socher, and C. Xiong, "End-to-end dense video captioning with masked transformer," in *CVPR*, 2018.

[36] Y. Zhu, O. Groth, M. Bernstein, and L. Fei-Fei, "Visual7w: Grounded question answering in images," in *CVPR*, 2016.