

# TutorialVQA: Question Answering Dataset for Tutorial Videos

Anthony Colas\*, Seokhwan Kim†, Franck Deroncourt†,

Siddhesh Gupte\*, Daisy Zhe Wang\*, Doo Soon Kim†

\*University of Florida, † Adobe Research

\* 432 Newell Dr, Gainesville, FL 32611, † 345 Park Avenue San Jose, CA 95110-2704

\*{acolas1, daisyw, siddhesh.gupte}@ufl.edu

†{seokim, deronco,dkim}@adobe.com

## Abstract

Despite the number of currently available datasets on video-question answering, there still remains a need for a dataset involving multi-step and non-factoid answers. Moreover, relying on video transcripts remains an under-explored topic. To adequately address this, we propose a new question answering task on instructional videos, because of their verbose and narrative nature. While previous studies on video question answering have focused on generating a short text as an answer, given a question and video clip, our task aims to identify a span of a video segment as an answer which contains instructional details with various granularities. This work focuses on screencast tutorial videos pertaining to an image editing program. We introduce a dataset, TutorialVQA, consisting of about 6,000 manually collected triples of (video, question, answer span). We also provide experimental results with several baseline algorithms using the video transcripts. The results indicate that the task is challenging and call for the investigation of new algorithms.

**Keywords:** question answering, resources, evaluation, corpora

## 1. Introduction

Video is the fastest growing medium to create and deliver information today. Consequentially, videos have been increasingly used as main data sources in many question answering problems (Yang et al., 2003; Tapaswi et al., 2016; Jang et al., 2017; Maharaj et al., 2017; Kim et al., 2017; Jang et al., 2017; Zeng et al., 2017). These previous studies have mostly focused on factoid questions, each of which can be answered in a few words or phrases generated by understanding multimodal contents in a short video clip.

However, this problem definition of video question answering causes some practical limitations for the following reasons. First, factoid questions are just a small part of what people actually want to ask on video contents. Especially if a short video is given to users, most fragmentary facts within the scope of previous tasks can be easily perceived by themselves even before asking questions. Thus, video question answering is expected to provide answers to more complicated non-factoid questions beyond the simple facts. For example, users may have “*how to*” type questions where the answer involves many fragmented steps to complete the task, as shown in Fig. 2.

Accordingly, the answer format needs to also be improved towards more flexible ways besides multiple choice (Tapaswi et al., 2016; Jang et al., 2017) or fill-in-the-blank questions (Maharaj et al., 2017; Kim et al., 2017). Although open-ended video question answering (Yang et al., 2003; Jang et al., 2017; Zeng et al., 2017) has been explored, it still aims to generate just a short word or phrase-level answer, which is not enough to cover the various granularities of non-factoid question answering.

The other issue is that many videos with sufficient amounts of information, which a user is likely to pose questions on, have much longer lengths than the video clips in the existing datasets. Therefore, in practice, the most relevant part of a whole video needs to be determined prior to each an-

swer generation phase. However, this localization task has been out of scope for previous studies.

In this work, we propose a new question answering problem for non-factoid questions on instructional videos. In accord with the nature of the media created for instructional purposes, we assume that many answers may already exist within the given video contents. Under this assumption, we formulate the problem as a localization task to specify the span of a video segment as the direct answer to a given video and question, as illustrated in Figure 1.

The remainder of this paper is structured as follows: Section 2 goes over some related work. Section 3 introduces TutorialVQA dataset as a case study of our proposed problem. The dataset includes about 6,000 triples, comprised of videos, questions, and answer spans manually collected from screencast tutorial videos with spoken narratives for a photo-editing software. Section 4 presents the baseline models and their experiment details on the sentence-level prediction and video segment retrieval tasks on our dataset. Then, we discuss the experimental results in Section 5 and conclude the paper in Section 6.

## 2. Related Work

Most relevant to our proposed work is the reading comprehension task, which is a question answering task involving a piece of text such as a paragraph or article. Such datasets for the reading comprehension task, such as SQuAD (Rajpurkar et al., 2016) based on Wikipedia, TriviaQA (Joshi et al., 2017) constructed from trivia questions with answer evidence from Wikipedia, or those from Hermann et al. based on CNN and Daily Mail articles (Hermann et al., 2015) are factoid-based, meaning the answers typically involve a single entity. Differing from video transcripts, the structures of these data sources, namely paragraphs from Wikipedia and news sources, are typically straightforward since they are meant to be read. In contrast, video transcripts originate from spoken dialogue, which can be verbose, unstruc-



Figure 1: An illustration of our task, where the red in the timeline indicates where answers can be found in a video.

tured, and disconnected. Furthermore, the answers in instructional video transcripts can be longer, spanning multiple sentences if the process is multi-step or even fragmented into multiple segments throughout the video.

Visual corpora in particular have proven extremely valuable to visual question-answering tasks (Antol et al., 2015), the most similar being MovieQA (Tapaswi et al., 2016) and VideoQA (Yang et al., 2003). Similar to how our data is generated from video tutorials, the MovieQA and VideoQA corpus is generated from movie scripts and news transcripts, respectively. MovieQA’s answers have a shorter span than the answers collected in our corpus, because questions and answer pairs were generated after each paragraph in a movie’s plot synopsis (Tapaswi et al., 2016). The MovieQA dataset also contains specific annotated answers with incorrect examples for each question. In the VideoQA dataset, questions focus on a single entity, contrary to our instructional video dataset. Although not necessarily a visual question-answering task, the work proposed by Gupta et al. involved answering questions over transcript data (Gupta et al., 2018). Contrary to our work, Gupta et al.’s dataset is not publically available and their examples only showcase factoid-style questions involving single entity answers.

Malmaud et al. focus on aligning a set of instructions to a video of someone carrying out those instructions (Malmaud et al., 2015). In their task, they use the video transcript to represent the video, which they later augment with a visual cue detector on food entities. Their task focuses on procedure-based cooking videos, and contrary to our task is primarily a text alignment task. In our task we aim to answer questions—using the transcripts—on instructional-style videos, in which the answer can involve steps not mentioned in the question.

### 3. TutorialVQA Dataset

In this section, we introduce the TutorialVQA dataset and describe the data collection process.<sup>1</sup>

#### 3.1. Overview

Our dataset consists of 76 tutorial videos pertaining to an image editing software. All of the videos include spoken instructions which are transcribed and manually segmented

number of videos	76
number of segments	408
number of QA pairs	6,195
avg. length of answer (sec)	31.39
avg. length of transcript (sentences)	48
avg. length of question (words)	9
avg. length of answer (sentences)	6

Table 1: Statistics of TutorialVQA dataset.

```
{
  "video_id": "19197",
  "question": "how do i output selections used with masks?",
  "answer_start": 44,
  "answer_end": 49
}
```

Figure 2: An example of a QA annotation.

into multiple segments. Specifically, we asked the annotators to manually divide each video into multiple segments such that each of the segments can serve as an answer to any question. For example, Fig. 1 shows example segments marked in red (each which are a complete unit as an answer span). Each sentence is associated with the starting and ending time-stamps, which can be used to access the relevant visual information.

The dataset contains 6,195 non-factoid QA pairs, where the answers are the segments that were manually annotated. Fig. 2 shows an example of the annotations. `video_id` can be used to retrieve the video information such as meta information and the transcripts. `answer_start` and `answer_end` denote the starting and ending sentence indexes of the answer span. Table 1 shows the statistics of our dataset, with each answer segment having on average about 6 sentences, showing that our answers are more verbose than those in previous factoid QA tasks.

#### 3.2. Basis

We chose videos pertaining to an image editing software because of the complexity and variety of tasks involved. In these videos, a narrator is communicating an overall goal by utilizing an example. For example, in 1 the video pertains to combining multiple layers into one image. How-

<sup>1</sup><https://github.com/acolas1/TutorialVQAData>

ever, throughout the videos multiple subtasks are achieved, such as the opening of multiple images, the masking of images, and the placement of two images side-by-side. These subtasks involve multiple steps and are of interest to us in segmenting the videos. Each segment can be seen as a sub-task within a larger video dictating an example. We thus chose these videos because of the amount of procedural information stored in each video for which the user may ask. Though there is only one domain, each video corresponds to a different overall goal.

### 3.3. Data Collection

We downloaded 76 videos from a tutorial website about an image editing program.<sup>2</sup> Each video is pre-processed to provide the transcripts and the time-stamp information for each sentence in the transcript. We then used Amazon Mechanical Turk<sup>3</sup> to collect the question-answer pairs.<sup>4</sup> One naive way of collecting the data is to prepare a question list and then, for each question, ask the workers to find the relevant parts in the video. However, this approach is not feasible and error-prone because the videos are typically long and finding a relevant part from a long video is difficult. Doing so might also cause us to miss questions which were relevant to the video segment. Instead, we took a reversed approach. First, for each video, we manually identified the sentence spans that can serve as answers. These candidates are of various granularity and may overlap. The segments are also complete in that they encompass the beginning and end of a task. In total, we identified 408 segments from the 76 videos. Second we asked AMT workers to provide question annotations for the videos.

Our AMT experiment consisted of two parts. In the first part, we presented the workers with the video content of a segment. For each segment, we asked workers to generate questions that can be answered by the presented segment. We did not limit the number of questions a worker can input to a corresponding segment and encouraged them to input a diverse set of questions which the span can answer. Along with the questions, the workers were also required to provide a justification as to why they made their questions. We manually checked this justification to filter out the questions with poor quality by removing those questions which were unrelated to the video. One initial challenge worth mentioning is that at first some workers input questions they had about the video and not questions which the video could answer. This was solved by providing them with an unrelated example. The second part of the question collection framework consisted of a paraphrasing task. In this task we presented workers with the questions generated by the first task and asked them to write the questions differently while keeping the semantics the same. In this way, we expanded our question dataset. After filtering out the questions with low quality, we collected a total of 6,195 questions.

It is important to note the differences between our data collection process and the the query generation process em-

<sup>2</sup><https://helpx.adobe.com/photoshop/tutorials.html>

<sup>3</sup><https://www.mturk.com/>

<sup>4</sup>We recruited workers who have completed at least 100 AMT tasks and have at least a 95% approval rating.

Example	Video ID
why might some alignment options appear grey?	4051
how would i go about selecting all the layers at once?	4051
what does the crop tool do?	4177
where is the crop tool located?	4177

Table 2: Examples of question variations

ployed in the Search and Hyperlinking Task at MediaEval (Eskevich et al., 2014). In the Search and Hyperlinking Task, 30 users were tasked to first browse the collection of videos, select interesting segments with start and end times, and then asked to conjecture questions that they would use on a search query to find the interesting video segments. This was done in order to emulate their thought process mechanism. While the nature of their task involves queries relating to the overall videos themselves, hence coming from a video’s interestingness, our task involves users already being given a video and formulating questions where the answers themselves come from within a video. By presenting the same video segment to many users, we maintain a consistent set of video segments and extend the possibility to generate a diverse set of question for the same segment.

### 3.4. Dataset Details

Table 2 presents some extracted sample questions from our dataset. The first column corresponds to an AMT generated question, while the second column corresponds to the video ID where the segment can be found. As can be seen in the first two rows, multiple types of questions can be answered within the same video (but different segments). The last two rows display questions which belong to the same segment but correspond to different properties of the same entity, ‘crop tool’. Here we observe different types of questions, such as “why”, “how”, “what”, and “where”, and can see why the answers may involve multiple steps. Some questions that the worked paraphrased were in the “yes/no” style, however our answer segments then provide an explanation to these questions.

Each answer segment was extracted from an image editing tutorial video that involved multiple steps and procedures to produce a final image, which can partially be seen in 1. The average number of sentences per video was approximately 52, with the maximum number of sentences contained in a video being 187. The sub-tasks in the tutorial include segments (and thus answers) on editing parts of images, instructions on using certain tools, possible actions that can be performed on an image, and identifying the locations of tools and features, with the shortest and longest segment having a span of 1 and 37 sentences respectively, demonstrating the heterogeneity of the answer spans.

## 4. Baselines

Our video question answering task is novel and to our knowledge, no model has been designed specifically for this task. As a first step towards solving this problem, we evaluated the performance of state-of-the-art models developed for other QA tasks, including a sentence-level predic-

tion task and two segment retrieval tasks. In this section, we report their results on the TutorialVQA dataset.<sup>5</sup>

#### 4.1. First Baseline: Sentence-level prediction

Given a transcript (a sequence of sentences) and a question, the first baseline predicts (*starting sentence index, ending sentence index*). The model is based on RaSor (Lee et al., 2016), which has been developed for the SQuAD QA task (Rajpurkar et al., 2016). RaSor concatenates the embedding vectors of the starting and the ending words to represent a span. Following this idea, the first baseline represents a span of sentences by concatenating the vectors of the starting and ending sentences. The left diagram in Fig. 3 illustrates the first baseline’s model.

**Model.** The model takes two inputs, a transcript,  $\{s_1, s_2, \dots, s_n\}$  where  $s_i$  are individual sentences and a question,  $q$ . The output is the span scores,  $y$ , the scores over all possible spans. GLoVe (Pennington et al., 2014) is used for the word representations in the transcript and the questions. We use two bi-LSTMs (Schuster and Paliwal, 1997) to encode the transcript.

$$h_i = biLSTM_{last}(s_i) \text{ for } i = 1..n \quad (\text{Sent Encoding})$$

$$p_i = biLSTM_{all}(\{h_1, \dots, h_n\}) \quad (\text{Psg Encoding})$$

where  $n$  is the number of sentences.<sup>6</sup> The output of Passage-level Encoding,  $p$ , is a sequence of vector,  $p_i$ , which represents the latent meaning of each sentence. Then, the model combines each pair of sentence embeddings ( $p_i, p_j$ ) to generate a span embedding.

$$r_{ij} = [p_i, p_j] \text{ for } i, j = 1..n \quad (\text{Span Generation})$$

where  $[\cdot, \cdot]$  indicates the concatenation. Finally, we use a one-layer feed forward network to compute a score between each span and a question.

$$h^q = biLSTM_{last}(q) \quad (\text{Span Scoring})$$

$$y_{ij} = \text{Softmax}(\text{FFN}([r_{ij}, h^q]))$$

In training, we use cross-entropy as an objective function. In testing, the span with the highest score is picked as an answer.

**Metrics.** We use tolerance accuracy (Tsunoo et al., 2017), which measures how far away the predicted span is from the gold standard span, as a metric. The rationale behind the metric is that, in practice, it suffices to recommend a rough span which contains the answer – a difference of a few seconds would not matter much to the user.

Specifically, the predicted span is counted as correct if  $|pred_{start} - gt_{start}| + |pred_{end} - gt_{end}| \leq k$ , where  $pred_{start/end}$  and  $gt_{start/end}$  indicate the indices of the predicted and ground-truth starting and ending sentences, respectively. We then measure the percentage of correctly predicted questions among the entire test questions.

<sup>5</sup>For the baselines, we considered only the transcript information, since it was non-trivial to include the other modalities such as video. In the future we plan to develop a multi-modal approach.

<sup>6</sup> $biLSTM_{last}$  produces only the last hidden vector while  $biLSTM_{all}$  produces all hidden vectors along the sequence.

#### 4.2. Second baseline: Segment retrieval

We also considered a simpler task by casting our problem as a retrieval task. Specifically, in addition to a plain transcript, we also provided the model with the segmentation information which was created during the data collection phrase (See Section. 3). Note that each segments corresponds to a candidate answer. Then, the task is to pick the best segment for given a query. This task is easier than the first baseline’s task in that the segmentation information is provided to the model. Unlike the first baseline, however, it is unable to return an answer span at various granularities. The second baseline is based on the attentive LSTM (Tan et al., 2016), which has been developed for the InsuranceQA task. The right diagram in Fig. 3 illustrates the second baseline’s model.

**Model.** The two inputs,  $s$  and  $q$  represent the segment text and a question. The model first encodes the two inputs.

$$h^s = biLSTM_{all}(s) \quad (\text{Sentence Encoding})$$

$$h^q = biLSTM_{last}(q) \quad (\text{Question Encoding})$$

$h^s$  is then re-weighted using attention weights.

$$a = \text{FFN}([h^s, h^q]) \quad (\text{Attention})$$

$$h'^s = a \odot h^s$$

where  $\odot$  denotes the element-wise multiplication operation. The final score is computed using a one-layer feed-forward network.

$$y = \text{Softmax}(\text{FFN}([h_s, h_q])) \quad (\text{Scoring})$$

During training, the model requires negative samples. For each positive example, (question, ground-truth segment), all the other segments in the same transcript are used as negative samples. Cross entropy is used as an objective function.

**Metrics.** We used accuracy and MRR (Mean Reciprocal Ranking) as metrics. The accuracy is

$$\frac{\# \text{ of questions where the top answer is correct}}{\text{the total \# of questions}}$$

We split the ground-truth dataset to train/dev/test into the ratio of 6/2/2. The resulting size is 3,718 (train), 1,238 (dev) and 1,239 QA pairs (test).

#### 4.3. Third baseline: Pipeline Segment retrieval

We construct a pipelined approach through another segment retrieval task, calculating the cosine similarities between the segment and question embeddings. In this task however, we want to test the accuracy of retrieving the segments given that we first retrieve the correct video from our 76 videos. First, we generate the TF-IDF embeddings for the whole video transcripts and questions. The next step involves retrieving the videos which have the lowest cosine distance between the video transcripts and question. We then filter and store the top ten videos, reducing the number of computations required in the next step. Finally, we calculate the cosine distances between the question and the

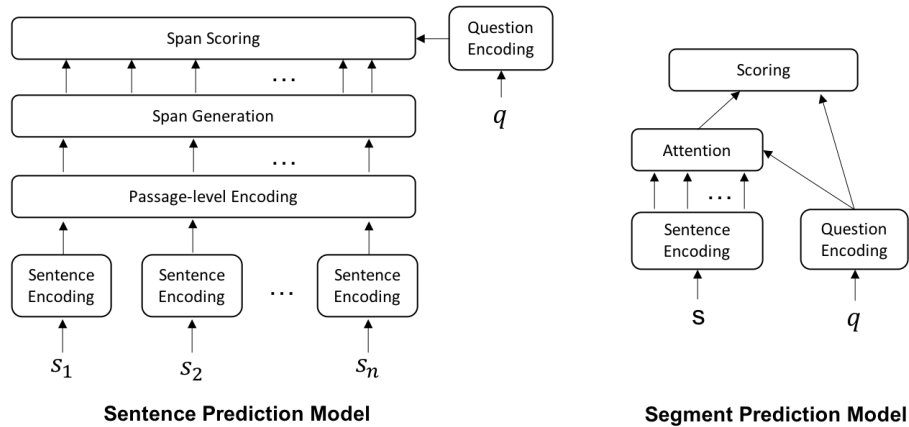


Figure 3: Baseline models for sentence-level prediction and video segment retrieval tasks.

segments which belong to the filtered top 10 videos, marking it as correct if found in these videos. While the task is less computationally expensive than the previous baseline, we do not learn the segment representations, as this task is a simple retrieval task based on TF-IDF embeddings.

**Model.** The first two inputs are the question,  $q$ , and video transcript,  $v$ , encoded by their TF-IDF vectors: (Tata and Patel, 2007):

$$v_{TF-IDF} = TF - IDF(v) \quad (\text{Video Encoding})$$

$$q_{TF-IDF} = TF - IDF(q) \quad (\text{Question Encoding})$$

We then filter the top 10 video transcripts (out of 76) with the minimum cosine distance, and further compute the TF-IDF vectors for their segments,  $S_{\text{top}10}^n$ , where  $n = 10$ . We repeat the process for the corresponding segments:

$$s_{TF-IDF} = TF - IDF(s) \quad (\text{Segment Encoding})$$

selecting the segment with the minimal cosine distance to the query.

**Metrics.** To evaluate our pipeline approach we use overall accuracy after filtering and accuracy given that the segment is in the top 10 videos. While the first metric is similar to 4.2, the second can indicate if initially searching on the video space can be used to improve our selection:

$$\frac{\text{\# of questions where the top answer is correct}}{\text{the total \# of questions answerable by top 10 videos}}$$

#### 4.4. Results

Tables 3, 4, 5 show the results. First, the tables show that the two first baselines under-perform for our task. Even with a tolerance window of 6, the first baseline merely achieves an accuracy of .14. The second baseline, despite being a simpler task, has only an accuracy of .23. Second, while we originally hypothesized that the segment selection task should be easier than the sentence prediction task, Table 4 shows that the task is also challenging. One possible reason is that the segments contained within the same transcript have similar contents, due to the composition of the overall task in each video, and differentiating among them may require a more sophisticated model than just using a

$k$	Train	Dev	Test
0	.0506	.0541	.0523
2	.0825	.0645	.0667
4	.1143	.1129	.1133
6	.1412	.1348	.1443

Table 3: Sentence-level prediction results for the first baseline with different tolerance window sizes  $k$ .

Model	MRR	ACC
Attentive LSTM	.4689	.2341

Table 4: Video segment retrieval results for the second baseline

sequence model for segment representation. Table 5 shows the accuracy of retrieving the correct segment, for baseline both overall and given that the video selected is within the top 10 videos. While the overall accuracy is only .16, by reducing the search space to 10 relevant videos our accuracy increases to .6385. In future iterations, it may then be useful to find better approaches in filtering large paragraphs of text before predicting the correct segment.

## 5. Discussion and Future Work

We performed an error analysis on the first baseline’s results. We first observe that, in 92% of the errors, the predicted span and the ground-truth overlap. Furthermore, in 56% of the errors, the predicted spans are a subset or superset of the ground-truth spans. This indicates that the model finds the rough answer regions but fails to locate the precise boundaries. To address this issue, we plan on exploring the Pointer-network (Vinyals et al., 2015), which finds an answer span by selecting the boundary sentences. Unlike the first baseline which avoids an explicit segmentation step, the Pointer-network can explicitly model which sentences are likely to be a boundary sentence. Moreover, the search space of the spans in the Pointer-network is  $2n$  where  $n$  is the number of sentences, because it selects only two boundary sentences. Note that the search space of the first baseline is  $n^2$ . A much smaller search space might improve the accuracy by making the model consider fewer candidates.

Overall	Given the video is in top 10
.1579	.6385

Table 5: Segment level prediction for the third baseline, both overall and given that the video is in the top 10.

In future work, we also plan to use multi-modal information. While our baselines only used the transcript, complementing the narratives with the visual information may improve the performance, similarly to the text alignment task in (Malmaud et al., 2015).

## 6. Conclusion

We have described the collection, analysis, and baseline results of TutorialVQA, a new type of dataset used to find answer spans in tutorial videos. Our data collection method for question-answer pairs on instructional video can be further adopted to other domains where the answers involve multiple steps and are part of an overall goal, such as cooking or educational videos. We have shown that current baseline models for finding the answer spans are not sufficient for achieving high accuracy and hope that by releasing this new dataset and task, more appropriate question answering models can be developed for question answering on instructional videos.

## 7. Bibliographical References

- Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Lawrence Zitnick, C., and Parikh, D. (2015). Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*, pages 2425–2433.
- Eskevich, M., Aly, R., Racca, D. N., Ordelman, R. J., Chen, S., and Jones, G. J. (2014). The search and hyperlinking task at mediaeval 2014. In *CEUR workshop proceedings*, volume 1263. CEUR-WS. org.
- Gupta, A., Mehrotra, R., and Gupta, M. (2018). Neural attention reader for video comprehension. In *ACM KDD 2018 Deep Learning Day*.
- Hermann, K. M., Kocisky, T., Grefenstette, E., Espeholt, L., Kay, W., Suleyman, M., and Blunsom, P. (2015). Teaching machines to read and comprehend. In *Advances in neural information processing systems*, pages 1693–1701.
- Jang, Y., Song, Y., Yu, Y., Kim, Y., and Kim, G. (2017). Tgif-qa: Toward spatio-temporal reasoning in visual question answering. In *Computer Vision and Pattern Recognition (CVPR), 2017 IEEE Conference on*, pages 1359–1367. IEEE.
- Joshi, M., Choi, E., Weld, D. S., and Zettlemoyer, L. (2017). Triviaqa: A large scale distantly supervised challenge dataset for reading comprehension. *arXiv preprint arXiv:1705.03551*.
- Kim, K.-M., Heo, M.-O., Choi, S.-H., and Zhang, B.-T. (2017). Deepstory: Video story qa by deep embedded memory networks. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence (IJCAI-17)*, pages 2016–2022.
- Lee, K., Salant, S., Kwiatkowski, T., Parikh, A., Das, D., and Berant, J. (2016). Learning recurrent span representations for extractive question answering. *arXiv preprint arXiv:1611.01436*.
- Maharaj, T., Ballas, N., Rohrbach, A., Courville, A. C., and Pal, C. J. (2017). A dataset and exploration of models for understanding video data through fill-in-the-blank question-answering. In *CVPR*, pages 7359–7368.
- Malmaud, J., Huang, J., Rathod, V., Johnston, N., Rabinovich, A., and Murphy, K. (2015). What’s cookin’? interpreting cooking videos using text, speech and vision. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 143–152.
- Pennington, J., Socher, R., and Manning, C. (2014). Glove: Global vectors for word representation. In *Proceedings of the 2014 conference on empirical methods in natural language processing (EMNLP)*, pages 1532–1543.
- Rajpurkar, P., Zhang, J., Lopyrev, K., and Liang, P. (2016). Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*.
- Schuster, M. and Paliwal, K. (1997). Bidirectional recurrent neural networks. *Trans. Sig. Proc.*, 45(11):2673–2681, November.
- Tan, M., Dos Santos, C., Xiang, B., and Zhou, B. (2016). Improved representation learning for question answer matching. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, volume 1, pages 464–473.
- Tapaswi, M., Zhu, Y., Stiefelshagen, R., Torralba, A., Urtasun, R., and Fidler, S. (2016). Movieqa: Understanding stories in movies through question-answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4631–4640.
- Tata, S. and Patel, J. M. (2007). Estimating the selectivity of tf-idf based cosine similarity predicates. *ACM Sigmod Record*, 36(2):7–12.
- Tsunoo, E., Bell, P., and Renals, S. (2017). Hierarchical recurrent neural network for story segmentation. In *INTERSPEECH*, pages 2919–2923.
- Vinyals, O., Fortunato, M., and Jaitly, N. (2015). Pointer networks. In *Advances in Neural Information Processing Systems*, pages 2692–2700.
- Yang, H., Chaisorn, L., Zhao, Y., Neo, S.-Y., and Chua, T.-S. (2003). Videoqa: question answering on news video. In *Proceedings of the eleventh ACM international conference on Multimedia*, pages 632–641. ACM.
- Zeng, K.-H., Chen, T.-H., Chuang, C.-Y., Liao, Y.-H., Niebles, J. C., and Sun, M. (2017). Leveraging video descriptions to learn video question answering. In *AAAI*, pages 4334–4340.