

Watching the News: Towards VideoQA Models that can Read

Soumya Jahagirdar[†] Minesh Mathew[†] Dimosthenis Karatzas[‡] C. V. Jawahar[†]

{soumya.jahagirdar, minesh.mathew}@research.iiit.ac.in dimos@cvc.uab.es jawahar@iiit.ac.in

[†] CVIT, IIIT Hyderabad, India [‡] Computer Vision Center, UAB, Spain

Abstract

Video Question Answering methods focus on common-sense reasoning and visual cognition of objects or persons and their interactions over time. Current VideoQA approaches ignore the textual information present in the video. Instead, we argue that textual information is complementary to the action and provides essential contextualisation cues to the reasoning process. To this end, we propose a novel VideoQA task that requires reading and understanding the text in the video. To explore this direction, we focus on news videos and require QA systems to comprehend and answer questions about the topics presented by combining visual and textual cues in the video. We introduce the “NewsVideoQA” dataset that comprises more than 8,600 QA pairs on 3,000+ news videos obtained from diverse news channels from around the world. We demonstrate the limitations of current Scene Text VQA and VideoQA methods and propose ways to incorporate scene text information into VideoQA methods.

1. Introduction

Visual Question Answering has evolved in numerous directions over the past few years. Two promising directions are, on one hand, the attempt to apply VQA on more dynamic scenarios, namely on video inputs and on the other hand, the introduction of scene text as an extra modality in the VQA process.

The reasoning processes required to tackle these challenges are not trivial to incorporate into a model. Taking into account the temporal dimension of an unfolding event requires reasoning over the evolution of certain actions, retrieving information from a specific time in the sequence, or a combination of the two. At the same time, recognizing the fact that world around us is littered with textual information that often carries important semantics necessary to interpret the scene has spawned a new direction in VQA. Introducing the scene text modality in the process requires incorporating error-prone reading systems, and connecting scene text semantics and literal transcriptions with the answer space.

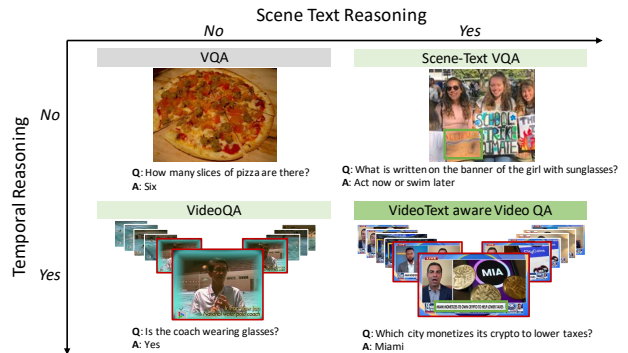


Figure 1: We address the task of text based Video Question Answering, incorporating VideoText (VideoText is the textual content embedded in the videos) information (bottom right). We propose a new dataset of News Videos along with QA annotations grounded on video text, and explore VQA models that jointly reason over temporal and text based information.

In this work, we attempt for the first time to join these two lines of research, and introduce the VideoText (VideoText is the textual content embedded in the videos) modality into Video Visual Question Answering.

Various attempts to apply VQA to the video setting have been proposed [8, 17, 28, 34, 40]. Such VideoQA methods have put forward datasets and methods focusing on recognizing actions, emotions, activities, and reasoning over temporal, causal correspondences, and knowledge graphs. However, they fall short in reasoning over the text appearing in the videos.

Scene Text VQA [3, 27], on the other hand, focuses on methods that allow VQA systems to incorporate scene text in the reasoning process. On one hand, this entails extracting semantics from noisy textual input, and on the other hand it requires dynamically expanding the answer space to incorporate new answer tokens afforded by the scene text [3, 21, 22, 27, 30]. Nevertheless, all scene text VQA methods are limited to processing a single image and cannot be readily extended to a multi-frame video input.

In this work, we attempt to combine multi-frame based, VideoQA architectures with the scene text modality

(Fig. 1). To explore this novel research direction, we define a new task and associated dataset: NewsVideoQA. Motivated by the prominent function of scene text in news video snippets, and the complementary information it carries to the visual modality, we consider that Visual Question Answering over News Videos is an adequate task to advance in models that jointly reason over temporal and scene text based information.

We present and thoroughly analyse the NewsVideoQA dataset, indicating key statistics and theoretical upper bound performance in various scenarios. We subsequently explore various baseline methods and demonstrate the limitations of both VideoQA and Scene Text VQA methods. We show that Scene Text VQA methods only yield top performance when they are applied on the video frame corresponding to the question (that includes the information needed to answer), but there is no trivial way for such methods to automatically retrieve the right frame. On the other hand, we show that VideoQA methods that do not consider the scene text, result in very low performance on the NewsVideoQA dataset. Finally, we repurpose a recently proposed VideoQA method to incorporate scene text information and show that it yields top results on the NewsVideoQA dataset, combining the benefits of both VideoQA and Scene Text VQA genres. The dataset is available at <http://cvit.iit.ac.in/research/projects/cvit-projects/videoqa>

The contributions of our work are the following:

- We introduce a new task of text based Video Question Answering, in which models must have the ability to read and reason about textual content in the videos (multi-frame input) to answer questions.
- We propose a new dataset: NewsVideoQA to explore the proposed task. This dataset comprises questions defined over the textual content in news videos and requires models to read and reason over it to obtain an answer.
- We evaluate various baselines on the NewsVideoQA dataset. These baselines include simple heuristic methods, text-only (machine comprehension) models, Scene Text VQA and VideoQA models.
- We repurpose the SINGULARITY [15] VideoQA model to the NewsVideoQA task and yield acceptable results compared to the original model.

2. Related Work

In this section, we briefly discuss some essential works in this space that is relevant to our work.

Video Question Answering. One of the early attempts at VideoQA is a retrieval-based approach for factoid QA proposed by Yang et al. [37]. Their system relies on speech

transcripts and external knowledge to answer the questions. One or more sentences from the transcript are returned as the output of the QA system, and the output is considered correct if the target answer is contained within the retrieved sentences. For QA evaluation, they used a private dataset containing only 40 QA pairs. Contrary to this work, our NewsVideoQA focuses particularly on the text appearing in the news videos, and is defined over a much larger dataset.

More recent works in VideoQA [17, 28, 34, 40] require models to reason about the events taking place in videos, but disregard any textual information in the videos. Tapaswi et al. [28] introduced a dataset that aims to study story comprehension using video and subtitles. Zhou et al. [40] introduced a large-scale VideoQA dataset that consists of videos of different activities. A method that gradually refines attention over the appearance and motion features is proposed in [34], along with an automatically generated dataset for VideoQA using subtitles. Yang et al. [36] and Maharaj et al. [20] focused on automatic generation of the VideoQA datasets. As the questions in [36] are automatically generated using captions, they are largely based on the visual appearance of objects and actions. Gupta et al. [8] explore knowledge-based question answering on news videos by proposing a new dataset. Questions in this dataset are primarily concerned with people seen in the videos, and the proposed models primarily rely on transcripts and an external knowledge base to find the answer. Questions in the above-mentioned works primarily require visual content and the transcripts of the videos to answer questions. Recently works such as [15, 16, 18, 19] have introduced transformer-based models with different pretraining strategies and yield state-of-the-art performance on existing VideoQA datasets.

Table 1 summarises existing works on VideoQA. It can be seen that majority of models focus on the visual content, transcripts and external knowledge to answer the questions. The text seen in the videos is an important source of information critical to understanding the content of news videos and videos shot outdoors. However, existing works on VideoQA largely disregard text on the videos. This motivates the community to have a publicly available video question answering dataset in which the questions require understanding of the textual content in the videos to obtain the answers.

Scene Text Aware Visual Question Answering (VQA). Early VQA datasets for natural images mainly included questions that seek information present in the visual content of the images [33]. However, realizing the importance of reading scene text to understanding natural images, researchers have recently started working on VQA tasks for natural images where questions based on textual information in the images are prioritized. This VQA branch is referred to as Scene Text VQA. Two popular benchmarks

Table 1: **A comparative overview of VideoQA datasets.** Datasets prior to our work, consider video, video + subtitles, video + knowledge base as input. Our work introduces a new line of research where the questions in proposed dataset are framed based on textual content in the news videos. The column Synthetic Gen. indicates the dataset which are synthetically/automatically generated.

Dataset	Subtitles	Text in video	Type of videos	Synthetic Gen.	Free-form	#Video	#QA
VideoQA [41]	✗	✗	Cooking, movies	✓	✗	109K	390K
MSVD-QA [34]	✗	✗	YouTube	✓	✓	1.9K	50K
ActivityNet-QA [40]	✗	✗	YouTube	✗	✓	5.8K	58K
MSRVTT-QA [34]	✗	✗	YouTube	✓	✓	10K	243K
MoviesQA [28]	✓	✗	Movies	✗	✗	6.7K	6.4K
TVQA [17]	✓	✗	TV shows	✗	✗	21K	152K
HowtoVQA69M [36]	✓	✗	TV shows	✓	✗	69M	69M
QA News Videos [37]	✗	✓	Web videos	-	-	-	40
NewsKVQA [8]	✓	✗	News videos	✓	✗	5.8K	58K
NewsVideoQA (Ours)	✓	✓	News videos	✗	✓	3.0K	8.6K

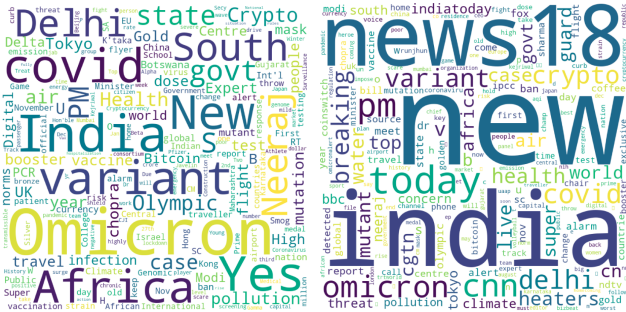


Figure 2: Word clouds of words in answers (left) and word clouds of words in OCR tokens (right).

for English scene text VQA are Scene Text VQA [3] and TextVQA [27]. Wang et al. [32] extended scene text VQA to a bi-lingual setting by introducing a new dataset that contains images with English and Chinese scene text. For scene text VQA, Singh et al. [27] proposed a model called LoRRA that uses top-down and bottom-up attention on scene text tokens and visual features to select an answer either from the OCR tokens or from a fixed vocabulary. M4C [10] uses a multimodal transformer-based model for Scene text VQA and Text VQA. This model, unlike the LoRRA can generate answers of any length by combining tokens from a fixed vocabulary or the scene text tokens found on the image. The current state-of-the-art models for scene text VQA typically use a Transformer-based architecture that is trained in two stages; a pretraining stage and a finetuning stage [2, 38]. The pretraining stage in these works is designed to learn multimodal interactions. In TAP [38], Yang et al. propose to pretrain an M4C-like architecture using pretraining tasks suitable for alignment between scene text and visual objects. TAP uses visual features corresponding to visual objects detected on the images using a pretrained object de-

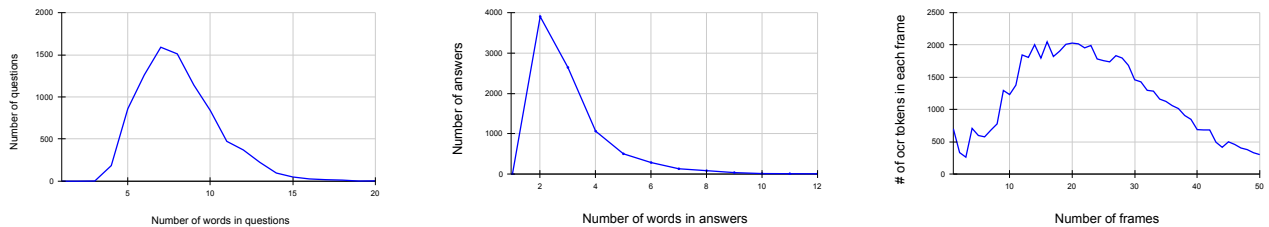
tection model, as done in most of the previous VQA works like LoRRA and M4C. Unlike TAP, which uses scene text, positional information, and visual features for pretraining on natural images, LaTr [2] uses document images for pre-training and uses only text and layout information. In the finetuning stage, LaTr uses visual features extracted using a pre-trained vision transformer.

In addition to scene text VQA, many specialized VQA tasks require reading and reasoning text on the images. There are multiple datasets for VQA on charts where text on the charts is critical to answer the questions [11, 12]. Mishra et al. [23] introduced a VQA dataset where all images are book covers, and the questions in the dataset are synthetically created using metadata associated with the books and question templates. Since the questions are created using information such as author names, titles and names of the publisher, questions purely depend on the text on the book covers and need little visual information. DocVQA [21] extends VQA to text-rich document images. This dataset has questions grounded on various document elements such as unstructured text in the form of paragraphs, forms, tables, and figures.

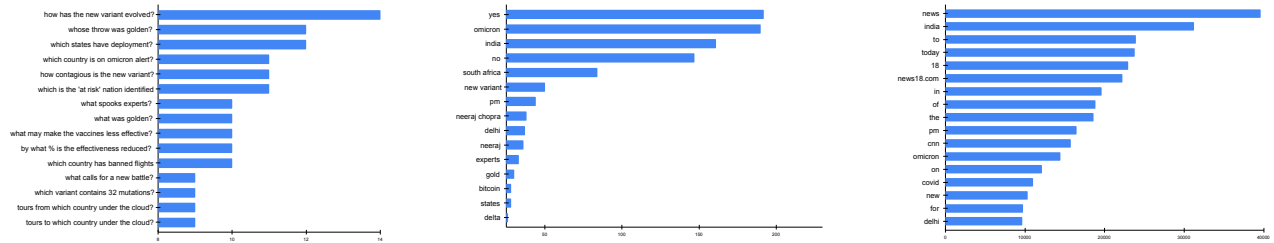
Similar to these existing VQA works that involve text on still images, we propose a VQA task that requires reading and understanding the text on videos. While the exact context—a single image—is given directly for the VQA problems, the proposed NewsVideoQA with multiple frames requires models to automatically find the right frames that informs the answer and reason over the textual content in those frames.

3. NewsVideoQA Dataset

In this section, we explain the data collection and annotation process. Also, we share statistics and analysis of the proposed NewsVideoQA dataset.



(a) **Questions with particular length.** The average length of questions in the dataset is 6.79 words. (b) **Answers with particular length.** The average number of words in the answers is 2.02 words. (c) **OCR tokens with particular length.** Average number of OCR tokens per frame is 26.14 tokens.



(d) Top 15 most occurring questions in the dataset. (e) Top 15 most occurring answers in the dataset. (f) Top 15 most occurring OCR tokens in the dataset.

Figure 3: Statistics for question, answer and OCR tokens in **NewsVideoQA** dataset.

3.1. Data Collection

News Videos: We collect news videos from English news channels around the world. We obtain videos from the following YouTube channels like BBC, ABC Australia, India Today, TRT World, AL Jazeera, CNN, NHK World Japan, Fox News, WION, NDTV, ABC News, CNN-News18, CTV News, CGTN, and IPCC. While collecting the news videos, we manually ensure that the videos are text-rich because the proposed task relies on video question answering, which requires reading text. The collected videos are split into 10 seconds of non-overlapping clips. The proposed dataset contains 3,083 videos, with at least 20 videos from each channel. The average number of questions per video is 2.96. The maximum number of questions defined for a video in the dataset is 20. The minimum number of questions defined for a video is 1.

Questions and Answers: The annotation process was organized into two stages. In stage 1, the annotators were instructed to define question-answer pairs based on textual information present in the news videos. Specifically, they were provided with the following instruction: ‘*Ensure that answering the questions generated requires reading of the text present in the news videos and should be related to the topic of that video*’. Annotators were asked to frame factoid questions that can be answered by reading the text present in the news videos. They were also instructed to add a timestamp: the time (with up to 1 second precision) of the video when the question was framed.

A second stage of verification was introduced to check the correctness of the data. Here, the annotators were asked to verify the data collected in the first stage. The annotators were shown the video-question pair for a video clip, and were asked to enter the answer and the timestamp and check the correctness of the question-answer pair based on its relevance to the textual content of the news video. They were asked to reject the questions with any grammatical mistakes in the questions and answers. During this stage, if the annotator finds a question-answer pair irrelevant to the topic or if the question was framed on the audio of the news videos, then such question-answer pairs were rejected from the dataset. A total of 1,200 QA pairs were rejected after the verification step. An extra stage was also added where the authors reviewed randomly picked question-answer pairs and their correctness and relevance to the task proposed.

3.2. Statistics and Analysis

The NewsVideoQA dataset comprises 8,672 questions framed on 3,083 news videos. The data is split randomly in 80-10-10 ratio to train, validation and test split. The train split has 6,994 questions over 2,407 videos, the validation split has 714 questions over 330 video clips, and the test split has 964 questions over 346 video clips.

Fig. 3a shows the distribution of question lengths for the questions in NewsVideoQA dataset. The average question length is 7.04 words. Among the 8,672 questions 7,008

(80.81%) are unique. Higher diversity in questions is reflective of the fact that questions are based on textual content. Fig. 3d shows the top 15 most frequent questions and their frequencies. Fig. 4 shows a sunburst plot of the first three words of the questions. It can be observed from Fig. 4 that there is variability in the question types like, questions starting with "What" which account for the questions related to the answer being directly present in the text of news videos. We provide subtitles of the news videos using a publicly available speech-to-text tool [29]. A total of 1,388 (17.36%) questions can be answered with sub-titles of the videos. This low percentage is observed due to two reasons, (a) smaller duration of the videos (10 seconds), resulting in incomplete sentences in the subtitles, and (b) most of the questions are based on textual content of the news videos. In total, there are 4,150 (47.85%) unique answers. Word cloud on the right in Fig. 2 shows the most common words in the answers. The answer space is broad and involves names of countries, events, games, people, etc. The distribution of answer lengths is shown in Fig. 3b. The average answer length is 2.02. The top 15 answers in the dataset are shown in Fig. 3e. We obtain OCR tokens using Google OCR. We uniformly sample the video at 2 frames per second and also retain the first frame of the video. Fig. 2 on the left shows the word cloud of OCR tokens. In Fig. 3f we show the top 15 OCR tokens present in the dataset. An average of 26.14 OCR tokens per frame is observed, and an average of 532.55 OCR tokens per video clip are observed in the dataset.

4. Baseline Methods

We evaluate three different methods as strong baselines for the newly introduced task of scene-text aware VQA on NewsVideoQA dataset. In this section, we briefly discuss the original methods and explain how these methods are adapted for the new task.

4.1. Heuristic methods and Upper Bounds

Inspired by heuristic baselines evaluated on scene text VQA [3, 27] and DocVQA [21] datasets, we evaluate the following heuristic baselines and upper bounds: **(i) Majority answer:** measures the performance when the most frequent answer in the train split is considered as the answer for all the questions in the test set. **(ii) Biggest OCR token:** measures the performance when the OCR token that occupies the largest area in the video is considered as the answer.

We compute upper-bound (UB) for the following cases: **(i) Vocabulary UB:** measures the maximum performance obtainable on the test set, if an answer is picked from a vocabulary of most common answers in the train split. **(ii) OCR Substring of single frame UB:** this measures the performance that can be obtained when we restrict our vo-

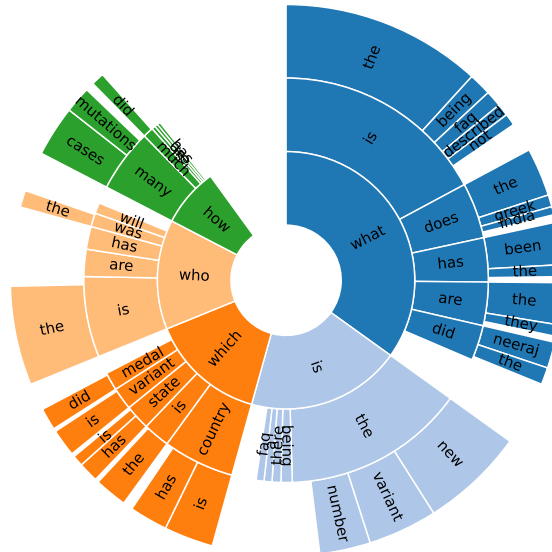


Figure 4: Distribution of questions by their starting 3-grams. Note that there is a diverse range of types of questions in the dataset. The question type "What" has a maximum count with questions such as "What is the ...?", "What does the ...?" and so on.

cabulary to list of OCR tokens of the frame on which the question was defined. **(iii) OCR Substring of all frames UB:** measures the performance we can obtain if the answer in the test split is a substring in the concatenated list of OCR tokens from uniformly sampled frames of the video.

4.2. Reading comprehension model

As observed in section 3, by design, almost all of the questions in NewsVideoQA are grounded on the text in the videos. For this reason, we evaluate a QA baseline that only considers the text in the videos to answer the questions. Specifically, we evaluate the BERT [5] QA model that is originally developed for extractive text-only QA. Extractive QA is a task of extracting a short snippet from the document/context on which the question is asked. The answer snippet is called a 'span' and the span is defined in terms of its start and end tokens. BERT is a transformer encoder-based method of pretraining language representations from unlabelled text. These pretrained models can be used later for downstream tasks with addition of output suitable for the task at hand. For the task of extractive QA, the additional layer, is an output layer that predicts start and end tokens of the span of the answer. For NewVideoQA, we concatenate the OCR tokens in a frame (assuming we know the correct frame) or the whole video—in our experiments we try out both settings—in the default reading order (i.e., top-left to bottom right order) and use this sequence as the context for the BERT QA model.

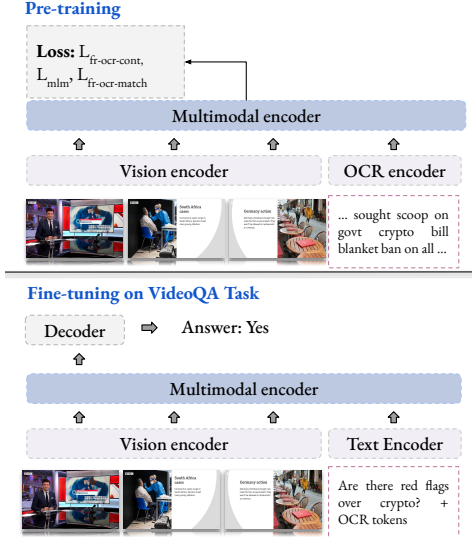


Figure 5: **OCR-aware SINGULARITY**. We extend SINGULARITY [15] for the task of text-based video question answering by incorporating OCR information by pretraining and finetuning on proposed NewsVideoQA dataset.

4.3. VQA Model

To evaluate the performance of current VQA models on NewsVideoQA dataset, we use M4C [10] model which takes into account the text present in the frames of the news videos. We pair each question with the frame corresponding to the timestamp of the question defined and consider it as input to M4C. M4C uses a multimodal transformer and an iterative answer prediction module. The tokens in the questions are embedded using a BERT model [5]. Each frame is represented using the following features: (i) appearance features of the objects detected using a Faster-RCNN pretrained on Visual Genome [14] and (ii) location information - bounding box coordinates of the detected objects.

Each OCR token recognized from the frame is represented using the following features: (i) a pretrained word embedding, which is FastText [4], (ii) appearance feature of the token’s bounding box from Faster-RCNN [25] (iii) PHOC [1] representation of the token and (iv) bounding box coordinates of the token. The representations of the entities mentioned, i.e., question tokens, objects and OCR tokens are projected to a common, learned embedding space. Later, a stack of transformer [31] layers is applied over these features in the common embedding space. The multi-head self-attention in transformers enables both inter-entity and intra-entity attention. In the end, answers are predicted through iterative decoding in an auto-regressive manner. At each step in the decoding, the decoded word is either an OCR token from the considered frame or a word from the fixed vocabulary of the common answer words.

Table 2: **Heuristics and Upper bound baseline results**. It can be seen that answers are substrings for more than 50 % of the serialized OCR tokens of a single frame corresponding to the timestamp of the question.

Heuristic Baselines	Acc. (%)
Majority answer	3.00
Biggest OCR token	1.03
Vocab Upper Bound	76.58
Substring single frame UB	53.05
Substring all frames UB	74.43

4.4. VideoQA Model

In addition to the text-only QA model and the text-based VQA models, we evaluate the performance of NewsVideoQA on a recently proposed transformer-based Retrieval and VideoQA method called SINGULARITY [15]. This method studies the importance of temporal relations to answer questions. SINGULARITY is a vision-language model pretrained on many video and image captioning datasets [9, 13, 34, 35, 39, 40]. It consists of three components, a vision encoder [6], a language encoder [5] and a multi-modal encoder [31]. For pretraining, each video/image is paired with its corresponding caption. The multi-modal encoder applies cross-attention to collect information from visual representations using the text as the key. Three pretraining objectives are defined: (i) Vision-Text Contrastive: a contrastive loss that aligns vision and text representations, (ii) Masked Language Modeling (MLM): predicts the masked visual and text contexts, and (iii) Vision-Text Matching: predicts the matching score of a vision-text pair with multi-modal encoder. For QA task, a multi-modal decoder is initialized from pretrained multi-modal encoder, which takes the outputs of multi-modal encoder as input. This generates an answer text with "[CLS]" as start token.

We extend the original SINGULARITY model [15], Fig. 5 and propose a new **OCR-aware VideoQA** version that can read the text in the videos and thereby answer questions based on the text in the videos. To this end, we include the OCR tokens in the videos as additional input during pretraining and finetuning stages. At the time of pretraining, unlike the original model that uses image/video + caption pairs, we use image/video + OCR tokens pairs. Similar to the original model, the following three pretraining objectives are employed.

(i) **Vision-OCR Contrastive loss**: aligns the visual features and OCR tokens, (ii) **Masked Language Modeling**: follows the formulation in BERT [5] to predict a randomly masked OCR token and (iii) **Vision-OCR Matching**: similar to Vision-OCR Contrastive loss, this allows the models

Table 3: **Comparison of all baselines on test set:** It can be seen that models such as BERT-QA [5] have poor performance when input of 12 frames followed by voting is provided at the test time. SINGULARITY [15] without any OCR information performs very poor as it does not consider OCR tokens as input. OCR-aware SINGULARITY performs better than all the baselines.

Model	#Frames for training	#Frames for testing	Acc. (%)	ANLS
BERT-QA [5]	1	1 (1 frame from the video)	28.70	34.21
BERT-QA [5]	1	1 (1 frame on which question was defined)	46.55	56.81
M4C [10]	1	1	28.49	32.17
BERT-QA [5]	1	2 (2 frames from the video)	15.03	17.65
BERT-QA [5]	1	2 (2 frames on which question was defined)	56.36	67.11
M4C [10]	1	2	27.87	31.54
BERT-QA [5]	1	12 (OCR tokens from 12 random frames)	53.86	65.27
M4C [10]	1	12	30.68	34.90
SINGULARITY [15]	1	12	4.82	5.78
OCR-aware SINGULARITY	1	12 (OCR tokens from a single frame)	33.57	37.52
OCR-aware SINGULARITY	1	12 (OCR tokens from 12 random frames)	32.47	35.56

to improve the alignment between paired vision and OCR inputs by using output of [CLS] token from multimodal encoder for binary classification. In essence, it says whether or not the input frame and OCR tokens pair match. Similar to the original model, we add multimodal decoder that has same architecture as that of multimodal encoder. This decoder uses multimodal encoder outputs as its cross-attention inputs. It decodes the answer with [CLS] as the start token.

5. Experiments

In this section, we explain evaluation metrics, and experimental settings and report the results. In all the experiments, we use the validation split of the dataset to save the best-performing checkpoints.

5.1. Evaluation Metrics

We use two evaluation metrics—Accuracy (Acc.) and Average Normalized Levenshtein Similarity (ANLS). Accuracy is the percentage of questions for which the predicted answer matches exactly with the target answer. The accuracy metric awards a zero score even when the prediction is a little different from the target answer. ANLS is a Levenshtein Similarity-based metric that acts softly on minor answer mismatches that might stem from an error in recognizing text on the images (i.e., OCR errors). Since all the answers in our dataset are derived from text seen in the videos, we found ANLS to be a suitable metric for NewsVideoQA.

5.2. Experimental setup

We run a commercial OCR engine to obtain OCR tokens for the evenly sampled frames.

BERT-QA. In the case of NewsVideoQA, we use the OCR tokens of the sampled video frames as context for BERT-QA. For each question, we obtain the OCR tokens of the frame on which the question is defined. We use the default OCR token ordering from the OCR system: top-left to bottom-right. To convert the NewsVideoQA dataset in SQuAD format, we find the first substring of the answer in the context, which is an approximation of the answer span as followed in [21]. We finetune the BERT QA checkpoint that is already pretrained and finetuned for QA on SQuAD dataset [24]. Specifically, we use the ‘bert-large-uncased-whole-word-masking-finetuned-squad’ checkpoint [7]. We train the BERT QA model starting from this checkpoint on NewsVideoQA dataset for ten epochs with a batch size of 32 and a learning rate of $2e - 05$.

M4C. For M4C, we use the official implementation along with default hyperparameters [26]. The fixed vocabulary used for answer generation is 3,751 words from answers in the train split of NewsVideoQA. Since M4C is a model for VQA on images, we train it using video frame + question pairs in the train split of NewsVideoQA. Similar to how BERT-QA was trained on NewsVideoQA, for each question the corresponding matching frame is found using the time-stamp information for each question that was collected during annotation.

SINGULARITY [15]. We use the pretrained model of SINGULARITY and finetune it on NewsVideoQA. We finetune it for 20 epochs and all the hyperparameters and training settings are kept the same as in the official implementation. SINGULARITY uses a single frame while training, and 12 randomly sampled frames while testing.

OCR-aware SINGULARITY. We continue pretraining the original SINGULARITY on our NewsVideoQA dataset



Figure 6: **Qualitative results** for different baselines on the proposed task. Results for baselines are shown in green for the correct predictions and in red for the incorrect predictions.

for 10 epochs. The vision encoder and the multimodal encoder are initialized similar to the original work. Following the pretraining on NewsVideoQA, we finetune the pre-trained model for 20 epochs with a batch size of 4 and learning rate of $1e-5$. The only difference compared to the original model is that we append OCR tokens to the question tokens. We keep the hyperparameters and pretraining settings the same as the SINGULARITY [15]. More details on experimental settings for above mentioned baselines can be found in supplementary material. In order to maintain the constant setting throughout all the baselines, (multiframe at the time of testing), for BERT-QA and M4C we perform additional experiments where these models are trained on single frame and are tested on multiple randomly sampled frames followed by a majority answer voting to obtain the final answer. Similar to SINGULARITY, we fix the number of frames used at the time of testing to be equal to 12.

5.3. Results

In Table. 2, we show the results of heuristics and upper-bound baselines. 3.0% of the questions can be answered by predicting “yes” which is the most common answer in the train split. Vocab Upper Bound of 76.58% shows that many answers in the train split repeat in the test split as well. In Table. 3, we show the comparative results for all the baselines. From the first four rows in the Tab. 3, it can be seen that BERT-QA (text only model), and M4C (Text-based single image VQA model) perform good when they are tested on one frame and two frames settings (frames based on the timestamp of the question). It can be seen that the performance of M4C reduces significantly when it is tested on 12 frames. For BERT-QA, the first row shows the performance when OCR tokens a single frame (not necessarily the frame on which the question was defined). This is followed by testing BERT-QA on OCR tokens of the frame on which the question was defined. From the table,

it can be seen that M4C performs poorly when the correct information required to answer the questions is not given as input to these models. SINGULARITY (without fine-tuning on NewsVideoQA) has poor performance compared to other baselines as the majority of the questions framed are based on textual content in the videos (Note: Results for BERT-QA in the first version were on OCR tokens not from the frames on which question was defined which led to poor performance of BERT-QA). We perform several experiments on the baselines which are present in the supplementary material. In Fig. 6, we show qualitative results from our experiments. The left example shows the predictions of baselines. As the frame contains less textual information all the baselines predict the correct answer. Whereas in the center and right example, the number of OCR instances increases thereby increasing the difficulty to obtain the correct answer.

6. Conclusion

We introduce and explore the problem of text based Video Question Answering, in which the models are encouraged to read and reason about the textual content in the videos. Towards this, we propose a new dataset, NewsVideoQA, which contains questions defined over textual content in news videos. We adopt existing baselines for text based video question answering on NewsVideoQA. Furthermore, we redesign the existing VideoQA method by incorporating OCR tokens to yield better results compared to the original method. Our exhaustive analysis and findings encourage the concurrent use of visual and textual cues for better video understanding systems. Our work will encourage researchers to develop better text-based Video Question Answering models and better insights into well-designed multimodal machine understanding models.

Acknowledgements This work is supported by MeitY, Government of India.

References

- [1] Jon Almazán, Albert Gordo, Alicia Fornés, and Ernest Valveny. Word spotting and recognition with embedded attributes. *IEEE Trans. Pattern Anal. Mach. Intell.*, 36(12):2552–2566, 2014.
- [2] Ali Furkan Biten, Ron Litman, Yusheng Xie, Srikar Apalparaju, and R. Manmatha. Latr: Layout-aware transformer for scene-text vqa. In *CVPR*, 2022.
- [3] Ali Furkan Biten, Rubèn Tito, Andrés Mafla, Lluís Gómez i Bigorda, Marçal Rusiñol, C. V. Jawahar, Ernest Valveny, and Dimosthenis Karatzas. Scene text visual question answering. In *ICCV*, pages 4290–4300. IEEE, 2019.
- [4] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomáš Mikolov. Enriching word vectors with subword information. *Trans. Assoc. Comput. Linguistics*, 5:135–146, 2017.
- [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *NAACL-HLT (1)*, pages 4171–4186. Association for Computational Linguistics, 2019.
- [6] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*. OpenReview.net, 2021.
- [7] Hugging Face. Huggingface models. <https://huggingface.co/models>. Accessed on 27 August 2022.
- [8] Pranay Gupta and Manish Gupta. Newskvqa: Knowledge-aware news video question answering. In *PAKDD (3)*, volume 13282 of *Lecture Notes in Computer Science*, pages 3–15. Springer, 2022.
- [9] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan C. Russell. Localizing moments in video with natural language. In *ICCV*, pages 5804–5813. IEEE Computer Society, 2017.
- [10] Ronghang Hu, Amanpreet Singh, Trevor Darrell, and Marcus Rohrbach. Iterative answer prediction with pointer-augmented multimodal transformers for textvqa. In *CVPR*, pages 9989–9999. Computer Vision Foundation / IEEE, 2020.
- [11] Kushal Kafle, Brian Price, Scott Cohen, and Christopher Kanan. DVQA: Understanding data visualizations via question answering. In *CVPR*, 2018.
- [12] Samira Ebrahimi Kahou, Vincent Michalski, Adam Atkinson, Ákos Kádár, Adam Trischler, and Yoshua Bengio. Figureqa: An annotated figure dataset for visual reasoning. In *ICLR (Workshop)*. OpenReview.net, 2018.
- [13] Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. Dense-captioning events in videos. In *ICCV*, pages 706–715. IEEE Computer Society, 2017.
- [14] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson, Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalantidis, Li-Jia Li, David A. Shamma, Michael S. Bernstein, and Li Fei-Fei. Visual genome: Connecting language and vision using crowdsourced dense image annotations. *Int. J. Comput. Vis.*, 123(1):32–73, 2017.
- [15] Jie Lei, Tamara L. Berg, and Mohit Bansal. Revealing single frame bias for video-and-language learning, 2022.
- [16] Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L. Berg, Mohit Bansal, and Jingjing Liu. Less is more: Clipbert for video-and-language learning via sparse sampling. In *CVPR*, pages 7331–7341. Computer Vision Foundation / IEEE, 2021.
- [17] Jie Lei, Licheng Yu, Mohit Bansal, and Tamara L. Berg. TVQA: localized, compositional video question answering. In *EMNLP*, pages 1369–1379. Association for Computational Linguistics, 2018.
- [18] Junnan Li, Dongxu Li, Caiming Xiong, and Steven C. H. Hoi. BLIP: bootstrapping language-image pre-training for unified vision-language understanding and generation. In *ICML*, volume 162 of *Proceedings of Machine Learning Research*, pages 12888–12900. PMLR, 2022.
- [19] Linjie Li, Yen-Chun Chen, Yu Cheng, Zhe Gan, Licheng Yu, and Jingjing Liu. HERO: hierarchical encoder for video+language omni-representation pre-training. In *EMNLP (1)*, pages 2046–2065. Association for Computational Linguistics, 2020.
- [20] Tegan Maharaj, Nicolas Ballas, Anna Rohrbach, Aaron C. Courville, and Christopher Joseph Pal. A dataset and exploration of models for understanding video data through fill-in-the-blank question-answering. In *CVPR*, pages 7359–7368. IEEE Computer Society, 2017.
- [21] Minesh Mathew, Dimosthenis Karatzas, and C. V. Jawahar. Docvqa: A dataset for vqa on document images. In *2021 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 2199–2208, 2021.
- [22] Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. Ocr-vqa: Visual question answering by reading text in images. In *ICDAR*, 2019.
- [23] Anand Mishra, Shashank Shekhar, Ajeet Kumar Singh, and Anirban Chakraborty. OCR-VQA: Visual question answering by reading text in images. In *ICDAR*, 2019.
- [24] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100, 000+ questions for machine comprehension of text. In *EMNLP*, pages 2383–2392. The Association for Computational Linguistics, 2016.
- [25] Shaoqing Ren, Kaiming He, Ross B. Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. In *NIPS*, pages 91–99, 2015.
- [26] Amanpreet Singh, Vedanuj Goswami, Vivek Natarajan, Yu Jiang, Xinlei Chen, Meet Shah, Marcus Rohrbach, Dhruv Batra, and Devi Parikh. Mmf: A multimodal framework for vision and language research. <https://github.com/facebookresearch/mmf>, 2020.
- [27] Amanpreet Singh, Vivek Natarajan, Meet Shah, Yu Jiang, Xinlei Chen, Devi Parikh, and Marcus Rohrbach. Towards vqa models that can read. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8317–8326, 2019.
- [28] Makarand Tapaswi, Yukun Zhu, Rainer Stiefelwagen, Antonio Torralba, Raquel Urtasun, and Sanja Fidler.

- Movieqa: Understanding stories in movies through question-answering. In *CVPR*, pages 4631–4640. IEEE Computer Society, 2016.
- [29] Silero Team. Silero models: pre-trained enterprise-grade stt / tts models and benchmarks. <https://github.com/snakers4/silero-models>, 2021.
- [30] Rubèn Tito, Dimosthenis Karatzas, and Ernest Valveny. Document collection visual question answering. In *ICDAR (2)*, volume 12822 of *Lecture Notes in Computer Science*, pages 778–792. Springer, 2021.
- [31] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NeurIPS*, pages 5998–6008, 2017.
- [32] Xinyu Wang, Yuliang Liu, Chunhua Shen, Chun Chet Ng, Canjie Luo, Lianwen Jin, Chee Seng Chan, Anton van den Hengel, and Liangwei Wang. On the general value of evidence, and bilingual scene-text visual question answering. In *CVPR*, 2020.
- [33] Qi Wu, Damien Teney, Peng Wang, Chunhua Shen, Anthony Dick, and Anton van den Hengel. Visual question answering: A survey of methods and datasets. *CVIU*, 163, 2017.
- [34] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang Zhang, Xiangnan He, and Yueting Zhuang. Video question answering via gradually refined attention over appearance and motion. In *ACM Multimedia*, 2017.
- [35] Jun Xu, Tao Mei, Ting Yao, and Yong Rui. MSR-VTT: A large video description dataset for bridging video and language. In *CVPR*, pages 5288–5296. IEEE Computer Society, 2016.
- [36] Antoine Yang, Antoine Miech, Josef Sivic, Ivan Laptev, and Cordelia Schmid. Just ask: Learning to answer questions from millions of narrated videos. In *ICCV*, pages 1666–1677. IEEE, 2021.
- [37] Hui Yang, Lekha Chaisorn, Yunlong Zhao, Shi-Yong Neo, and Tat-Seng Chua. Videoqa: Question answering on news video. pages 632–641, 01 2003.
- [38] Zhengyuan Yang, Yijuan Lu, Jianfeng Wang, Xi Yin, Dinei Florencio, Lijuan Wang, Cha Zhang, Lei Zhang, and Jiebo Luo. Tap: Text-aware pre-training for text-vqa and text-caption. In *CVPR*, 2021.
- [39] Youngjae Yu, Jongseok Kim, and Gunhee Kim. A joint sequence fusion model for video question answering and retrieval. In *ECCV (7)*, volume 11211 of *Lecture Notes in Computer Science*, pages 487–503. Springer, 2018.
- [40] Zhou Yu, Dejing Xu, Jun Yu, Ting Yu, Zhou Zhao, Yueting Zhuang, and Dacheng Tao. Activitynet-qa: A dataset for understanding complex web videos via question answering. In *AAAI*, pages 9127–9134, 2019.
- [41] Linchao Zhu, Zhongwen Xu, Yi Yang, and Alexander G. Hauptmann. Uncovering the temporal context for video question answering. *Int. J. Comput. Vis.*, 124(3), 2017.