

Contents

1	Introduction	2
1.1	Motivation	2
1.2	Problem Description	3
1.3	Objectives	3
2	Background	4
2.1	Deep Learning	4
2.1.1	Convolutional Neural Networks	4
2.2	Image Processing	5
2.2.1	Object Detection	5
2.2.2	Commonly Used Metrics	6
2.2.3	Multiple Object Tracking	7
2.3	Knowledge Representation and Reasoning	8
2.3.1	Answer Set Programming	8
2.3.2	The Action Language \mathcal{AL}	9
2.3.3	Symbolic Rule Learning	10
3	Related Work	10
3.1	Datasets for VideoQA	10
3.2	VideoQA Implementations	13
3.3	Hybrid Models for VQA	14
4	Project Plan	14
	References	15
	Appendix A ASP encoding of \mathcal{AL}	20

1 Introduction

Writing algorithms which can answer questions on pictures or videos with a high level of accuracy and generality has been a goal of researchers in the AI community for many years. Recently, a lot of progress has been made in this area; with advances in neural network models and the production of larger datasets allowing reseachers to significantly improve accuracy on question answering models.

Formally, Visual Question Answering (VQA) [3] is a task where, given an image and a question posed in natural language about the image, a model is required to produce an open-ended answer to the question. Video Question Answering (VideoQA) is a related task where a model is given a video (mulitple images in sequence) and a question. These questions can be related to a single frame of the video, effectiively making VideoQA a superset of the VQA task.

This project will attempt to produce a hybrid model for VideoQA - one which makes use of both neural networks and knowledge representation and reasoning methods based on first-order logic.

1.1 Motivation

VQA and VideoQA tasks attract attention because of their difficulty; both problems are considered “AI-complete” - they require knowledge from multiple modalities beyond a single sub-domain [2]. Building systems which have a deep understanding of the world would be a significant achievement for AI research; allowing many tasks which require significant human time and effort to be automated. Solving the VideoQA problem, which require image understanding, natural language understanding and commensense reasoning to be deployed, could be a major step towards this.

Furthermore, subproblems of VideoQA have already been shown to have applications in real-world tasks, for example event recognition has been used for identifier attacks on computer networks [9], detecting credit card fraud [42] and recognising cardiac arrhythmias [37].

Finally, the use of a hydrid model for VideoQA brings with it a nmuber of advantages. Firstly, representing knowledge in logical form allows the injection of commensense or background knowledge which can significantly improve accuracy in question answering tasks [34]. Secondly, there has recently been a significant increase in research related to explainable AI methods - machine learning techniques that enable human users to understand, trust and manage emerging artificially intelligent partners [4]. Extracting the knowledge from a neural network into logical form could be an important step toward explaining and understanding their behaviour.

1.2 Problem Description

As mentioned above, the VideoQA problem can be defined as building a model which, when presented with a short video and an open-ended, natural language question about the video, can produce a natural language answer to the given question. In our case we are looking to design a hybrid model for VideoQA. More specifically, convolutional neural networks (CNNs) will be used to extract knowledge from each frame of the video. This knowledge, along with the question to be answered, will be represented in a fashion that is amenable to searching for the answer to the question. This framing of the problem leads us to outline the following sub-problems, which will need to be solved in order to produce a satisfactory VideoQA model.

1. **Object Detection.** Given a frame, we need a model which can produce a rough estimate (a bounding box, for example) of the location of an object in the frame. We will also need a model which can classify each detected object into a set of predefined classes.
2. **Property Extraction.** Given an object, which is the output of the ‘object detection’ model above, we need an algorithm which can produce a set of values for that object for some set of predefined properties. For example, we might need to give a value for the colour, size or shape of an object.
3. **Event Detection.** Given two sequential frames, we need a model which can classify the event(s) which occurred between the two frames into a set of predefined classes (possibly including a catch-all ‘no event’ class). This model will also be required to list objects involved in the event and what their role in the event was. This will require some level of object tracking so that it is clear how objects are related between frames.
4. **Question and Knowledge Representation.** Given the outputs of the models above and the natural language question, we need a way of representing the background knowledge, the question and the knowledge contained in the frames of the video. These must be represented in a manner that allows an answer to the question to be found.

1.3 Objectives

A lot of research has been conducted on building end-to-end neural network models to solve VideoQA tasks (see section 3 for examples), but very little prior work has been done on hybrid models. The main aim of this project, therefore, is to explore the possibility of adding logical representation of knowledge to existing deep learning techniques. The following are the primary objectives of the project:

1. Construct a hybrid VideoQA model which allows injection of background knowledge and helps to increase the explainability of the model.
2. Find a challenging dataset for training the model (or construct a dataset if none are suitable)
3. Compare qualitatively (and quantitatively, if possible) existing approaches to our own hybrid model.
4. Investigate the possibility of learning domain-dependent logical rules, rather than having to have them hard coded for each environment.

2 Background

This section introduces some technical background which will likely be required as part of the project. It includes an introduction to neural networks and CNNs, a comparison of some existing object tracking and object detection algorithms and a discussion on knowledge representation and reasoning and symbolic rule learning.

2.1 Deep Learning

Deep neural networks (DNNs) have emerged as a very successful algorithm for machine learning; deep learning has been used to beat records in tasks such as image recognition, speech recognition and language translation [29]. Many different architectures have been proposed to solve various tasks, these architectures include convolutional neural networks (CNNs), which are designed to process data that come in the form of multiple arrays [29], and recurrent neural networks (RNNs), which are designed to process sequences of arbitrary length [32]. The following section gives a brief introduction to CNNs and describes some of their use cases.

2.1.1 Convolutional Neural Networks

CNNs contain three types of layers: convolution, pooling and fully connected. Units (artificial neurons) in a convolution layer are organised into feature maps. The inputs to each unit in a feature map come from the outputs of the units in a small region of the previous layer, the output of the unit is then calculated by passing the weighted sum of its inputs through an activation function such as ReLU. The set of weights, also known as a filter or kernel, is the part of the layer which is learned through backpropagation. Every unit in a feature map has the same kernel. Each feature map in a layer has its own kernel. Pooling layers reduce the size of the input by merging multiple units into one. A typical pooling operation is max-pooling, which computes the maximum of a local patch of units. Finally, in fully-connected

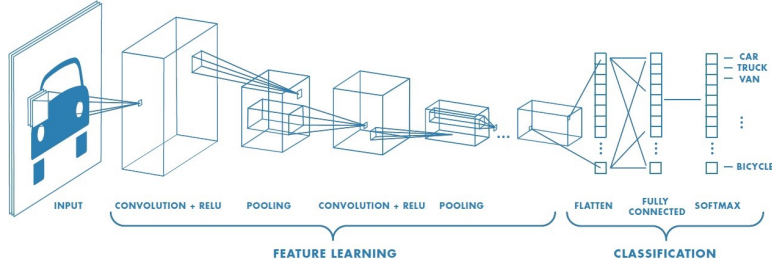


Figure 1: An example of a CNN architecture. The input image is passed through a series of convolution and pooling layers before being flattened into a one-dimensional layer and passed through one final fully connected layer. The softmax classification function is then applied at the output.

layers (which are typically placed at the output of the CNN) every unit in a layer is connected to every unit in the previous layer. An example CNN architecture is shown in figure 1.

CNNs have proven to be adept at a number of tasks involving images, including image classification [22] and object detection [38, 40]. We explore these further in Section 2.2.

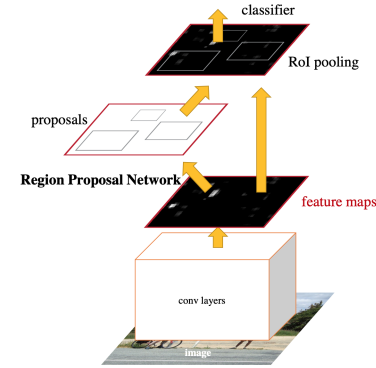
2.2 Image Processing

2.2.1 Object Detection

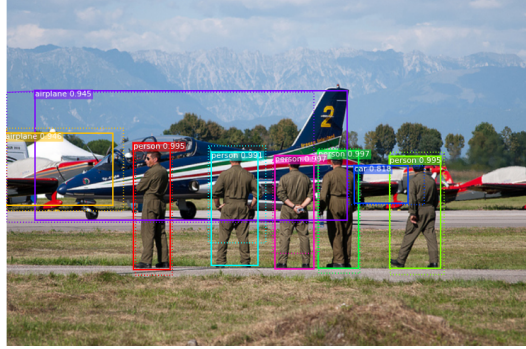
The object detection task could formally be defined as designing a model which, when given an image, can produce a rough localisation of objects of interest in the image (in the form of a bounding box) and classify each of these objects into a set of predefined classes. In this section we introduce two well known object detection algorithms, Faster R-CNN [40] and You Only Look Once (YOLO) [38].

Faster R-CNN is an evolution of previous object detection algorithms, R-CNN [16] and Fast R-CNN [15]. Faster R-CNN builds on its predecessors by adding a region proposal network (RPN) - a neural network which takes an image and produces a set of region of interest (RoI) proposals. This method of region proposal is much faster than previous algorithms (such as those used in [16] and [15]) since it is able to make use of the GPU, as opposed to requiring the CPU. Faster R-CNN then uses a similar classifier and bounding box regressor as Fast R-CNN at the output; this section of the network also receives the feature maps from the final layer of the RPN, in this sense the initial layers of the network are shared between the region proposal section and the classifier/regressor section. A diagram of the Faster R-CNN architecture is shown in figure 2a.

The three object detection algorithms mentioned above all work by first producing region proposals, then producing a more accurate localisation and



(a) Diagram of the Faster R-CNN architecture. Figure from [40].



(b) An example of the bounding boxes and confidence scores produced by an object detection algorithm. Image from [1].

Figure 2

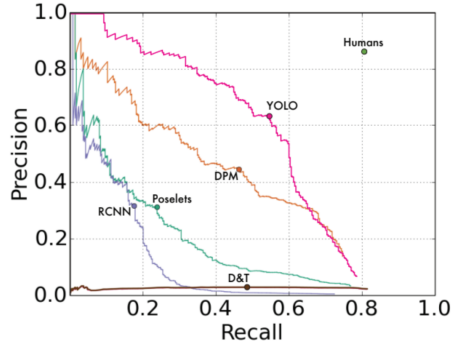
a score for each region and finally removing any low-scoring or redundant regions. This requires the algorithm to ‘look’ at the image multiple times (around 2000 times for R-CNN). You Only Look Once (YOLO) is a significantly more time-efficient algorithm which, as the name suggests, takes a single look at the image. A convolutional neural network is used to simultaneously predict multiple bounding boxes and the class probabilities for each box. As well as being very fast, YOLO makes fewer than half the number of background errors (where the algorithm mistakes background patches for objects) as Fast R-CNN [39]. YOLO is, however, slightly less accurate than some of the slower methods for object detection [38].

2.2.2 Commonly Used Metrics

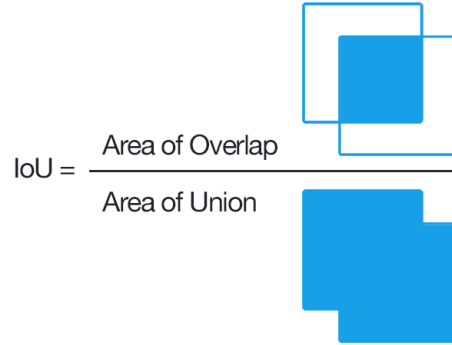
In this section we present some commonly used metrics for classification and object detection tasks. We use TP, TN, FP and FN to mean True Positive, True Negative, False Positive and False Negative, respectively.

Firstly, for classification tasks the following terminology is commonly used:

- $Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$. The accuracy is the ratio of correct predictions to the total number of predictions.
- $Precision = \frac{TP}{TP+FP}$. The precision is the ability of a classifier to not label the negative data as positive.
- $Recall = \frac{TP}{TP+FN}$. The recall is the ability of the classifier to find the positively-labelled data.
- $F_1 = 2 * \frac{precision * recall}{precision + recall}$. The F_1 score is a way of combining the precision and recall scores.



(a) Example precision-recall curves for various object detection models. Image from [39].



(b) Visual explanation of the intersection over union metric. Image from [41].

Figure 3: Precision-recall curves and definition of intersection over union

Each object detector model will output a confidence score for each object classification it makes. We can then set a threshold value which determines what is counted as a classification of an object. Altering this threshold value will give different precision and recall values for the model, which can then be plotted on a precision-recall graph. Example precision-recall curves are shown in figure 3a.

For object detection tasks, where a bounding box is produced as a rough localisation of an object’s position, metrics which measure the accuracy of the localisation of an object are required. One very common metric is the Average Precision (AP), which is roughly defined as the area under the precision-recall curve (estimates of this value are usually used for competition datasets). In order to assess how well a model localises an object in an image the Intersection over Union (IoU) metric is commonly used. Figure 3b gives a visualisation on how IoU is calculated. Each IoU threshold will produce its own precision-recall curve.

2.2.3 Multiple Object Tracking

Multiple Object Tracking (MOT) is a computer vision task that aims to identify and track objects from a sequence of images without any prior knowledge about the appearance or number of targets [7]. Most object tracking methods share a very similar pipeline [7], as follows:

- **Detection.** Each input frame is analysed to identify objects using bounding boxes.
- **Feature Extraction.** Algorithms extract appearance, motion and/or interaction features from objects.

- **Affinity Computation.** Features are used to compute a similarity score between objects.
- **Association.** Similarity scores are used to associate detections and compute the object trajectories.

The Simple Online and Realtime Tracking (SORT) [5] algorithm is one of the best performers [7]. It makes use of Faster R-CNN for object detection, the Kalman filter [20] framework for predicting object motion, IoU as a similarity function and the Hungarian algorithm [23] for association. SORT attempts to track objects using motion prediction, other algorithms, however, focus on extracting features from objects, such as Histogram of Orientated Gradients (HOG) [8] features, and using these to match objects in successor frames. More recent methods tend to use features extracted directly from a CNN.

2.3 Knowledge Representation and Reasoning

Knowledge Representation and Reasoning (KRR) is concerned with how intelligent agents store and manipulate their knowledge. In this section we discuss Answer Set Programming, a logic programming framework, action languages, which can be used to define the behaviour of a system, and methods for learning logical rules inductively.

2.3.1 Answer Set Programming

Answer Set Programming (ASP) is a form of declarative logic programming that can be used to solve difficult search problems. Whereas imperative programs define an algorithm for finding a solution to a problem, logic programs simply define a problem, it is then the job of logic program solvers to find a solution to the problem. ASP also differs from Prolog, another popular logic programming language, in that ASP programs are purely declarative. This means that reordering rules or atoms within rules has no effect on the output of the solver [10]. ASP solvers work by finding the answer sets of the program, where each rule in the program imposes restrictions on what can be an answer set. An answer set can be thought of as a set of ground atoms which satisfies every rule of the program (although the full definition of an answer set is too in-depth for this discussion).

The following templates are a some of the possible forms of rules in an ASP program:

$$a \text{ :- } b_1, \dots, b_k, \text{not } b_{k+1}, \dots, \text{not } b_m. \quad (1)$$

$$\{c_1; \dots; c_n\} \text{ :- } b_1, \dots, b_k, \text{not } b_{k+1}, \dots, \text{not } b_m. \quad (2)$$

$$\text{ :- } b_1, \dots, b_k, \text{not } b_{k+1}, \dots, \text{not } b_m. \quad (3)$$

Where each a , b_i and c_i are atoms in first-order logic. The left hand side of the rule is known as the head, and the right hand side is known as the body. The *not* in the rule body stands for negation-as-failure, which means that *not* b_i will be satisfied when b_i cannot be proved. Each rule requires that when the body is satisfied - that is, when every member of the body is satisfied - the head must be satisfied. Rule 2 is known as a choice rule. The head of a choice rule can be satisfied by any subset of the atoms inside the brackets (including the empty set); in effect, a choice rule creates possible answer sets. Finally, rule 3 is known as a constraint. The body of a constraint must not be satisfied; intuitively, a constraint rules out answer sets.

As well as negation as failure, ASP also has a notion of ‘strong negation’. The strong negation of an atom p is written $\neg p$. Strong negation can be thought of as logical negation, although it does not always have the same properties. In practice, ASP solvers implement strong negation by treating $\neg p$ as an additional atom, and enforce that no answer set can contain both p and $\neg p$.

2.3.2 The Action Language \mathcal{AL}

Action languages are formal models for describing the behaviour of dynamic systems. In this section we present the version of \mathcal{AL} given in [14]. \mathcal{AL} ’s signature contains three special sorts: *statics*, *fluents* and *actions*. Fluents are partitioned into two sorts: *inertial* and *defined*. Statics and fluents are both referred to as ‘domain properties’. A ‘domain literal’ is a domain property or its negation. Statements in \mathcal{AL} can be of the following form:

$$a \text{ \textbf{causes} } l_{in} \text{ \textbf{if} } p_0, \dots, p_m \quad (4)$$

$$l \text{ \textbf{if} } p_0, \dots, p_m \quad (5)$$

$$\text{\textbf{impossible} } a_0, \dots, a_k \text{ \textbf{if} } p_0, \dots, p_m \quad (6)$$

Where a is an action, l and p_0, \dots, p_m are domain literals and l_{in} is a literal formed by an inertial fluent. Statement 4 is known as a *causal law*, 5 as a *state constraint* and 6 as an *executability condition*. A collection of \mathcal{AL} statements is known as a ‘system description’. An \mathcal{AL} system description can be used to model the behaviour of dynamic systems with discrete states; each state can be seen as the set of fluents which are true and transitions between states are caused by actions.

It is also possible to encode a given \mathcal{AL} system description, along with a number of ‘domain-independent’ axioms, in ASP. The method for creating this encoding is given in Appendix A.

2.3.3 Symbolic Rule Learning

Inductive Logic Programming [36] (ILP) is a field of symbolic AI research concerned with learning symbolic rules which, when combined with background knowledge, entail a set of positive examples and do not entail any negative examples. ILASP [28] (Inductive Learning of Answer Set Programs) is an ILP framework for learning ASP programs.

The authors of [24] define the *Learning from Answer Sets* (ILP_{LAS}) task (which is the task solved by the original version of ILASP), by first defining a *partial interpretation*. A partial interpretation E is a pair of sets of atoms E^{inc} and E^{exc} , known as the *inclusions* and *exclusions* of E . An answer set A *extends* E if it contains all of the inclusions ($E^{inc} \subseteq A$) and none of the exclusions ($E^{exc} \cap A = \emptyset$). An ILP_{LAS} task is then defined as the tuple $T = \langle B, S_M, E^+, E^- \rangle$, where B is the background knowledge, S_M is the search space, E^+ and E^- are the partial interpretations for the positive and negative examples, respectively. An hypothesis H is known as an inductive solution of T if and only if all of the following are true:

1. $H \subseteq S_M$
2. $\forall e^+ \in E^+ \exists A \in AS(B \cup H)$ such that A extends e^+
3. $\forall e^- \in E^- \nexists A \in AS(B \cup H)$ such that A extends e^-

Where $AS(P)$ refers to the answer sets of a program P .

Later versions of ILASP are capable of solving more complex tasks, including learning weak constraints [27] (a method for specifying preferences in ASP), learning from context dependent examples [26] and learning from noisy examples [25].

3 Related Work

3.1 Datasets for VideoQA

A number of datasets are available for the VideoQA problem, in this section we discuss each of the available datasets. A comparison of all the datasets discussed in this section is shown in Table 1.

The MovieQA dataset [46] is a VideoQA dataset consisting of 14,944 multiple-choice questions which align with video clips from movies. The clips come from a collection of 408 movies and the Question-Answer (QA) pairs were generated by humans in three steps. Firstly, one annotator would be given a plot summary for a movie and asked to create correct QA pair, additionally this annotator would mark the area of the plot summary to which their QA pair corresponded. Next, a different annotator would be asked to provide four incorrect answers to a QA pair (and was given the

option to correct the original question and answer). Finally, annotators would align each sentence in the plot summary to the video by marking the start and end points, this allowed each QA pair to be aligned with a video clip. The questions and each of the possible answers are written in natural language.

Zeng et al. [52] create a much larger VideoQA dataset by automatically generating QA pairs from videos and associated descriptions collected online. Their dataset consists of 18100 videos as well 151263 and 21352 automatically generated QA pairs in the training and validation sets, respectively. The dataset also contains 2461 human-generated QA pairs to be used for testing. Their questions and answers are free-form natural language, however, a large number of their answers are yes and no (32.5% and 32.5%, respectively).

The TGIF-QA dataset [19] is commonly used for assessing the performance of VideoQA models. The updated version of the TGIF-QA dataset contains human-generated 165,165 QA pairs collected from 71,741 GIFs, sourced from the TGIF dataset [31], which contains a number of GIFs and associated descriptions. There are four possible types of questions in the TGIF-QA dataset, three specific to VideoQA, requiring temporal knowledge to answer, and the last asks questions related to a single frame (VQA). The question types are as follows:

- **Repetition Count.** Counting the number of repetitions of an action. There are 11 possible answers (0 to 10+).
- **Repeating Action.** A multiple-choice question about identifying an action that has been repeated in the video. For example, **Q:** What does the duck do three times? **A:** Shake head.
- **State Transition.** A multiple-choice question about identifying the state before or after another state. For example, **Q:** What does the bear do after sitting? **A:** Stand.
- **FrameQA.** Open-ended questions related to a single frame.

For the VideoQA questions the authors created templates for questions and used a large number of human annotators to speed up the generation process. The FrameQA questions are generated using the descriptions from the TGIF dataset. A number of quality control checks were also included.

Zhu et al. [53] have proposed a VideoQA dataset containing fill-in the blank style questions, with multiple-choice answers. The dataset contains over 100,000 real-world video clips and 400,000 questions. The dataset is generated from three different annotated video sources. On top of questions which ask the model to describe the present (describe the current video), for two of three video sources the authors also introduce two additional question types: infer the past and predict the future. For these two types

of questions the model is asked a question on a part of the video which it is not explicitly given; these questions require the model to use some form of commonsense reasoning to generate a correct answer. One of the advantages of using a multiple-choice dataset such as this is that it is more amenable to quantitative evaluation than datasets with free-form answers, since answers are either right or wrong.

The EgoVQA dataset [11] attempts to address the lack of first-person VideoQA datasets. The dataset contains 581 QA pairs with both multiple-choice questions (with 5 possible answers per question) and open-ended questions. The dataset was generated by manually generating QA pairs from a pre-existing set of 16 first-person videos. The authors also show that existing VideoQA models only marginally outperformed random choice on questions related to the colour of objects. They conjecture that existing models struggle to separate attentions on camera wearers from attentions on third persons.

Xu et al. [48] generate two VideoQA datasets by converting video captions into QA pairs. The first dataset, known as MSVD-QA, is generated from the Microsoft Research Video Description Corpus [6] which is used in many video captioning experiments. MSVD-QA contains 1,970 video clips and 50,505 QA pairs. Similarly, the second dataset, known as MSRVT-T-QA, is generated from the MSR-VTT dataset [49]. The MSRVT-T-QA dataset contains 10,000 video clips and 243,680 QA pairs.

The YouTube2Text-QA dataset [50] is another large dataset for VideoQA generated from a pre-existing video description dataset, in this case the YouTube2Text [17] dataset is used. The YouTube2Text-QA dataset consists of 1,970 videos and 99,421 QA pairs.

The TVQA dataset [30] contains 21,793 video clips and 152,545 QA pairs based on 6 popular TV shows. The QA pairs were annotated manually using *Amazon Mechanical Turk*. Workers were asked to generate questions in the format: [What/How/Where/Why/Who/Other] _____ [when/before/after] _____. The second part of the question localises the relevant video moment within the clip, while the first part contains the question about that moment. The answers to the questions are given in multiple-choice format, with 5 candidate answers for each question.

The PororoQA dataset [21] is created from video clips and subtitles of the children’s cartoon series ‘*Pororo*’. The dataset contains 8,913 multiple-choice QA pairs and 16,066 video clips.

Finally, the MovieFIB dataset [33] is a large-scale fill-in the blank style dataset generated from movie descriptions. The dataset contains 128,085 video clips and 348,998 QA pairs. The questions concern entities, actions and objects; answering these questions therefore implies that a model has some level of visual understanding of the scene, rather than being able to answer based purely on the given sentence. Answers are open-ended (not multiple choice) but each answer is only a single word.

Table 1: Comparison of discussed VideoQA datasets. Each row contains data on: the number of videos/clips, the number of QA pairs, whether the uses multiple-choice questions, whether the dataset uses fill-in the blank questions and the video source.

Dataset	#Videos	#QA pairs	MC	FIB	Source
MovieQA	408 ¹	14,944	Y	N	Movies
Zeng et al.	18,100	175,076	N	N	Online videos
TGIF-QA	71,741	165,165	Y ²	N	Online videos
Zhu et al.	>100,000	400,000	Y	Y	Various
EgoVQA	16	581	Y	N	First-person videos
MSVD-QA	1,970	50,505	N	N	Video desc. corpus
MSRVTT-QA	10,000	243,680	N	N	Video desc. corpus
Youtube2Text-QA	1,970	99,421	Y	N	YouTube videos
TVQA	21,793	152,545	Y	N	TV shows
PororoQA	16,066	8,913	Y	N	Cartoon series
MovieFIB	128,085	348,998	N	Y	Movie description

¹ Full length movies. Some of the QAs come with timestamps, allowing more video clips.

² FrameQA questions (53,083 QA pairs) are not MC.

3.2 VideoQA Implementations

All of the datasets described in the section above (except [11]) were presented along with original neural network models for solving the VideoQA task. This section attempts to summarise some of these approaches, along with a few others [13, 51, 12, 44]. However, since the focus of our approach will not be an end-to-end neural architecture, so we do not give detailed descriptions. An example of a typical VideoQA neural architecture is given in figure 4.

All previous work presented here contains a video encoder for extracting features from frames of the video. These usually include both appearance and motion features. These are extracted from pre-trained networks such as ResNet [18], VGG [43] or GoogLeNet [45], for appearance, and C3D [47], for motion. The extracted features from each frame are then usually passed into LSTM or GRU networks to obtain encodings for the whole video.

Questions are often encoded by generating word embeddings for each word of the sentences (along with basic preprocessing such as tokenisation) and then passing the list of words to a sentence encoder, such as the LSTM or GRU architectures.

Visual attention [35] is used to help neural networks focus on the most relevant areas of an image or video (and decrease the amount of computation required). Applying attention mechanisms to associate a question with its most relevant frames (or areas within frames) has become a key part of more recent VideoQA models [19, 48, 50, 30, 13, 51, 12, 21]. Temporal attention

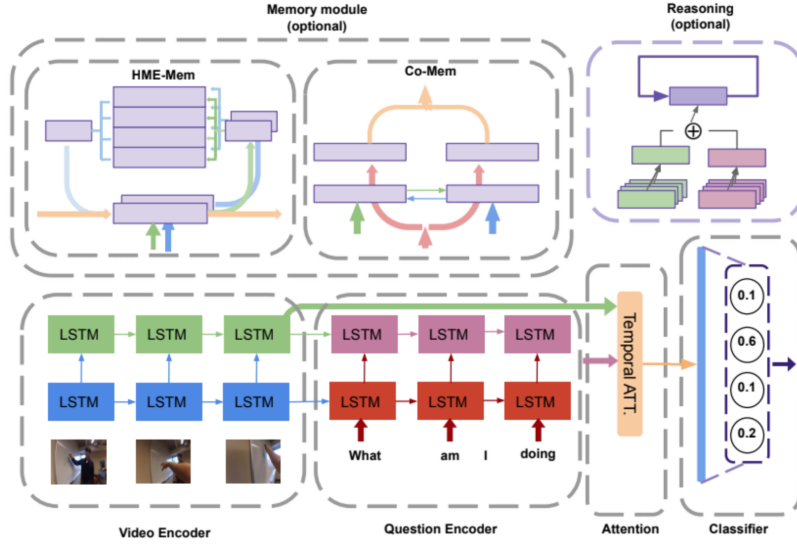


Figure 4: Typical VideoQA architecture. Figure from [11].

is commonly applied to help the model focus on the most salient frames of the video, however, [19] applies both temporal and spatial attention, which also allows the model to attend to the most relevant regions of a frame.

Additional modules can be added to the overall network to improve performance; for example, the architectures from [12, 52, 13] all make use of external memory to allow multiple iterations of inference. Furthermore, [12] adds an additional external memory module and a multimodal fusion layer which can combine features from questions and images.

3.3 Hybrid Models for VQA

4 Project Plan

References

- [1] W. Abdulla. *Mask R-CNN for object detection and instance segmentation on Keras and TensorFlow*. URL: https://github.com/matterport/Mask_RCNN (visited on 01/15/2020).
- [2] Somak Aditya, Yezhou Yang, and Chitta Baral. “Explicit Reasoning over End-to-End Neural Architectures for Visual Question Answering”. In: *CoRR* abs/1803.08896 (2018).
- [3] Aishwarya Agrawal et al. “Vqa: Visual question answering”. In: *International Journal of Computer Vision* 123.1 (2017), pp. 4–31.
- [4] Alejandro Barredo Arrieta et al. *Explainable Artificial Intelligence (XAI): Concepts, Taxonomies, Opportunities and Challenges toward Responsible AI*. 2019. arXiv: 1910.10045 [cs.AI].
- [5] Alex Bewley et al. “Simple online and realtime tracking”. In: *2016 IEEE International Conference on Image Processing (ICIP)*. IEEE. 2016, pp. 3464–3468.
- [6] David L Chen and William B Dolan. “Collecting highly parallel data for paraphrase evaluation”. In: *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*. Association for Computational Linguistics. 2011, pp. 190–200.
- [7] Gioele Ciaparrone et al. “Deep learning in video multi-object tracking: A survey”. In: *Neurocomputing* (2019).
- [8] Navneet Dalal and Bill Triggs. “Histograms of Oriented Gradients for Human Detection”. In: *International Conference on Computer Vision & Pattern Recognition (CVPR '05)*. Vol. 1. IEEE Computer Society, 2005, pp. 886–893.
- [9] Christophe Dousson, Pierre Le Maigat, and France Telecom R&d. “Chronicle Recognition Improvement Using Temporal Focusing and Hierarchization”. In: *IJCAI*. 2007, pp. 324–329.
- [10] Thomas Eiter, Giovambattista Ianni, and Thomas Krennwallner. “Answer Set Programming: A Primer”. In: *Reasoning Web. Semantic Technologies for Information Systems: 5th International Summer School 2009*. Ed. by Sergio Tessaris et al. Springer Berlin Heidelberg, 2009, pp. 40–110.
- [11] Chenyou Fan. “EgoVQA - An Egocentric Video Question Answering Benchmark Dataset”. In: *The IEEE International Conference on Computer Vision (ICCV) Workshops*. 2019.

- [12] Chenyou Fan et al. “Heterogeneous Memory Enhanced Multimodal Attention Model for Video Question Answering”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2019, pp. 1999–2007.
- [13] Jiyang Gao et al. “Motion-appearance co-memory networks for video question answering”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2018, pp. 6576–6585.
- [14] Michael Gelfond and Yulia Kahl. *Knowledge representation, reasoning, and the design of intelligent agents: The answer-set programming approach*. Cambridge University Press, 2014.
- [15] Ross Girshick. “Fast r-cnn”. In: *Proceedings of the IEEE international conference on computer vision*. 2015, pp. 1440–1448.
- [16] Ross Girshick et al. “Rich feature hierarchies for accurate object detection and semantic segmentation”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2014, pp. 580–587.
- [17] Sergio Guadarrama et al. “Youtube2text: Recognizing and describing arbitrary activities using semantic hierarchies and zero-shot recognition”. In: *Proceedings of the IEEE international conference on computer vision*. 2013, pp. 2712–2719.
- [18] Kaiming He et al. “Deep residual learning for image recognition”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 770–778.
- [19] Yunseok Jang et al. “Tgif-qa: Toward spatio-temporal reasoning in visual question answering”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2017, pp. 2758–2766.
- [20] Rudolph Emil Kalman et al. “A new approach to linear filtering and prediction problems [J]”. In: *Journal of basic Engineering* 82.1 (1960), pp. 35–45.
- [21] Kyung-Min Kim et al. “DeepStory: Video Story QA by Deep Embedded Memory Networks”. In: *IJCAI17*. AAAI Press, 2017, 20162022.
- [22] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. “ImageNet Classification with Deep Convolutional Neural Networks”. In: *Advances in Neural Information Processing Systems 25*. Ed. by F. Pereira et al. Curran Associates, Inc., 2012, pp. 1097–1105.
- [23] Harold W Kuhn. “The Hungarian method for the assignment problem”. In: *Naval research logistics quarterly* 2.1-2 (1955), pp. 83–97.
- [24] Mark Law, Alessandra Russo, and Krysia Broda. “Inductive Learning of Answer Set Programs”. In: *Logics in Artificial Intelligence*. Springer International Publishing, 2014, pp. 311–325.

- [25] Mark Law, Alessandra Russo, and Krysia Broda. “Inductive learning of answer set programs from noisy examples”. In: *arXiv preprint arXiv:1808.08441* (2018).
- [26] Mark Law, Alessandra Russo, and Krysia Broda. “Iterative learning of answer set programs from context dependent examples”. In: *Theory and Practice of Logic Programming* 16.5-6 (2016), pp. 834–848.
- [27] Mark Law, Alessandra Russo, and Krysia Broda. “Learning weak constraints in answer set programming”. In: *Theory and Practice of Logic Programming* 15.4-5 (2015), pp. 511–525.
- [28] Mark Law, Alessandra Russo, and Krysia Broda. *The ILASP system for learning Answer Set Programs*. www.ilasp.com. 2015.
- [29] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. “Deep learning”. In: *Nature* 521 (2015), pp. 436–444.
- [30] Jie Lei et al. “Tvqa: Localized, compositional video question answering”. In: *arXiv preprint arXiv:1809.01696* (2018).
- [31] Yuncheng Li et al. “TGIF: A new dataset and benchmark on animated GIF description”. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 2016, pp. 4641–4650.
- [32] Pengfei Liu, Xipeng Qiu, and Xuanjing Huang. “Recurrent neural network for text classification with multi-task learning”. In: *arXiv preprint arXiv:1605.05101* (2016).
- [33] Tegan Maharaj et al. “A Dataset and Exploration of Models for Understanding Video Data through Fill-in-the-Blank Question-Answering”. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 7359–7368.
- [34] Kenneth Marino et al. “OK-VQA: A Visual Question Answering Benchmark Requiring External Knowledge”. In: *CoRR* abs/1906.00067 (2019).
- [35] Volodymyr Mnih, Nicolas Heess, Alex Graves, et al. “Recurrent models of visual attention”. In: *Advances in neural information processing systems*. 2014, pp. 2204–2212.
- [36] Stephen Muggleton. “Inductive logic programming”. In: *New generation computing* 8.4 (1991), pp. 295–318.
- [37] Ren Quiniou et al. “Intelligent Adaptive Monitoring for Cardiac Surveillance”. In: *Computational Intelligence in Healthcare 4: Advanced Methodologies*. Ed. by Isabelle Bichindaritz et al. Springer Berlin Heidelberg, 2010, pp. 329–346.
- [38] Joseph Redmon and Ali Farhadi. “Yolov3: An incremental improvement”. In: *arXiv preprint arXiv:1804.02767* (2018).

- [39] Joseph Redmon et al. “You only look once: Unified, real-time object detection”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 779–788.
- [40] Shaoqing Ren et al. “Faster r-cnn: Towards real-time object detection with region proposal networks”. In: *Advances in neural information processing systems*. 2015, pp. 91–99.
- [41] Adrian Rosebrock. *Intersection over Union (IoU) for object detection*. Nov. 7, 2016. URL: <https://www.pyimagesearch.com/2016/11/07/intersection-over-union-iou-for-object-detection> (visited on 01/15/2020).
- [42] Nicholas Poul Schultz-Møller, Matteo Migliavacca, and Peter Pietzuch. “Distributed Complex Event Processing with Query Rewriting”. In: *Proceedings of the Third ACM International Conference on Distributed Event-Based Systems*. Association for Computing Machinery, 2009.
- [43] Karen Simonyan and Andrew Zisserman. “Very deep convolutional networks for large-scale image recognition”. In: *arXiv preprint arXiv:1409.1556* (2014).
- [44] Gursimran Singh. “Spatio-temporal relational reasoning for video question answering”. PhD thesis. University of British Columbia, 2019.
- [45] Christian Szegedy et al. “Going deeper with convolutions”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2015, pp. 1–9.
- [46] Makarand Tapaswi et al. “MovieQA: Understanding Stories in Movies through Question-Answering”. In: *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2015), pp. 4631–4640.
- [47] Du Tran et al. “Learning spatiotemporal features with 3d convolutional networks”. In: *Proceedings of the IEEE international conference on computer vision*. 2015, pp. 4489–4497.
- [48] Dejing Xu et al. “Video Question Answering via Gradually Refined Attention over Appearance and Motion”. In: *MM ’17*. 2017.
- [49] Jun Xu et al. “MSR-VTT: A Large Video Description Dataset for Bridging Video and Language”. In: *IEEE International Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [50] Yunan Ye et al. “Video question answering via attribute-augmented attention network learning”. In: *Proceedings of the 40th International ACM SIGIR Conference on Research and Development in Information Retrieval*. ACM. 2017, pp. 829–832.

- [51] Youngjae Yu et al. “End-to-End Concept Word Detection for Video Captioning, Retrieval, and Question Answering”. In: *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)* (2016), pp. 3261–3269.
- [52] Kuo-Hao Zeng et al. “Leveraging Video Descriptions to Learn Video Question Answering”. In: *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence*. AAAI Press, 2017, 4334–4340.
- [53] Linchao Zhu et al. “Uncovering the temporal context for video question answering”. In: *International Journal of Computer Vision* 124.3 (2017), pp. 409–421.

Appendix A ASP encoding of \mathcal{AL}

In this appendix we present a method for encoding an \mathcal{AL} system description in ASP, as described in [14]. We need to encode three different parts of the system description: the signature, the \mathcal{AL} statements and some domain-independent axioms.

We encode the signature of the system description, $sig(SD)$, as follows:

- For each constant symbol c of sort *sort_name* other than *fluent*, *static* or *action*, $sig(SD)$ contains

$$sort_name(c). \quad (7)$$

- For every static g of SD, $sig(SD)$ contains

$$static(g). \quad (8)$$

- For every inertial fluent f of SD, $sig(SD)$ contains

$$fluent(inertial, f). \quad (9)$$

- For every defined fluent f of SD, $sig(SD)$ contains

$$fluent(defined, f). \quad (10)$$

- For every action a of SD, $sig(SD)$ contains

$$action(a). \quad (11)$$

In the following we refer to the ASP encoding of the \mathcal{AL} system description as $\Pi(SD)$, where $\Pi(SD)$ includes $sig(SD)$. We introduce a relation $holds(f, i)$ which says that fluent f is true at timepoint i . We also introduce the notation $h(l, i)$ where l is a domain literal and i is a step, which will not be used in the ASP program, but will instead be replaced by either $holds(f, i)$ if $l = f$, or by $\neg holds(f, i)$ if $l = \neg f$.

We encode the \mathcal{AL} statements as follows:

- If the maximum number of steps is $< max >$, then $\Pi(SD)$ includes

$$\#const\ n = < max > . \quad (12)$$

$$step(0..n). \quad (13)$$

- For every causal law, a **causes** l **if** p_0, \dots, p_m $\Pi(SD)$ contains

$$h(l, I + 1) :- h(p_0, I), \dots, h(p_m, I), occurs(a, I), I < n. \quad (14)$$

- For every state constraint, l **if** p_0, \dots, p_m $\Pi(SD)$ contains

$$h(l, I) :- h(p_0, I), \dots, h(p_m, I). \quad (15)$$

- For every executability condition, **impossible** a_0, \dots, a_k **if** p_0, \dots, p_m $\Pi(SD)$ contains

$$\neg occurs(a_0, I); \dots; \neg occurs(a_k, I) :- h(p_0, I), \dots, h(p_m, I). \quad (16)$$

The $;$ in the head of rule 16 stands for logical disjunction and can be read as meaning at least one of a_0, \dots, a_k must not occur at timepoint I .

In order to complete our encoding of an \mathcal{AL} system description we need to add a number of domain-independent axioms. These axioms are not specific to any system or task, but rather convey commonsense knowledge that should apply to many systems. It is worth noting, however, that in certain situations some or all of these axioms may not make sense and should not be used.

We encode the domain-independent knowledge as follows:

- The inertia axiom states that inertial fluents will keep their state unless explicitly changed:

$$\begin{aligned} holds(F, I+1) &:- fluent(inertial, F), \\ &holds(F, I), \\ ¬ \neg holds(F, I+1), \\ &I < n. \end{aligned} \quad (17)$$

$$\begin{aligned} \neg holds(F, I+1) &:- fluent(inertial, F), \\ &\neg holds(F, I), \\ ¬ holds(F, I+1), \\ &I < n. \end{aligned} \quad (18)$$

- The closed-world assumption (CWA) for defined fluents states that defined fluents which are not known to be true are assumed to be false.

$$\neg holds(F, I) :- fluent(defined, F), step(I), not holds(F, I). \quad (19)$$

- The CWA for actions states that actions that are not known to occur are assumed to not occur.

$$\neg occurs(A, I) :- action(A), step(I), not occurs(A, I). \quad (20)$$

Finally, in order to make use of system description encoding we would include information on which events occurred, using the $occurs(A, I)$ predicate, and information on the different states of the system using $holds(F, I)$ and $\neg holds(F, I)$.