

Chapter 1

Dataset

Before discussing the implementation of our solution to the VideoQA task, it is important to outline the data that is used to train and evaluate the model. Rather than using one of the datasets outlined in Chapter ??, we opted to create a new VideoQA dataset, which we name ‘OceanQA’. While it would have been preferable to use an existing dataset (and an existing implementation as a baseline) to allow a fair comparison, none of the existing datasets suited the project requirements¹, for the following two reasons:

1. Most of the existing VideoQA datasets use videos from real-world environments, where objects and events are usually more complex than computer-generated environments. Training models to work with real-world data therefore requires significant computational resources and can take days or even weeks. Given the time and resource limitations that exist for this project, it was sensible to avoid these datasets. More generally, creating a dataset gives greater flexibility over the size and complexity of the data; if faster training is required, we can simply create smaller images or use fewer objects in each video.
2. Hybrid models generally require some form of environment specification (or ontology) so that the model’s internal knowledge can be represented explicitly (see [10, 4, 11] for examples from VideoQA and VQA). Since most of the existing VideoQA datasets do not limit objects to be of specific types, or restrict object properties or video events to a given set, they cannot provide such a specification of the environment.

The dataset is generated programmatically and can therefore be made as large as required. However, in order to allow comparisons between models, we generate a fixed dataset of 1400 videos (each video contains 10 question-answer pairs) and use this data to train and evaluate the models outlined in subsequent Chapters. This dataset is divided into 1000 training, 200 validation and 200 testing videos.

¹The CLEVRER dataset [10] meets these requirements and would have been a good candidate for this project. Unfortunately, it was published in March 2020, six months after the project began.

The generated dataset can be used in either ‘full-data’ form or in ‘QA-data’ form. The full-data form contains videos, question-answer pairs and the ground truth of all the information in the videos; this means that every object property, relation and event is labelled by the dataset. This form gives the programmer complete access to the information in the video, which allows baseline models to be constructed, and may also allow internal parts of models to be evaluated. The QA-form, on the other hand, contains only videos and question-answer pairs. Since no additional data about the video is provided, this form reflects a ‘real’ VideoQA dataset, and should be used for evaluating the model as a whole.

Since the focus of this project is to investigate logical reasoning, we do not attempt to make the job of the neural network difficult by creating complex scenes; the dataset emulates a simple retro-game environment. Each image is also quite small at 256x256 pixels to allow faster network training. The remainder of this Chapter outlines the full details of the OceanQA dataset.

1.1 Videos

Each video in the OceanQA dataset is a sequence of 32 frames. Each frame contains a flat background and a maximum of 16 objects. Objects form the central component of each video, since all of the useful information in each video can be modelled by the following: properties of objects; binary relations between objects; and events, which relate to at least one object, occurring between two consecutive frames of the video. Each object is modelled using the following attributes: object type (or class), position, rotation and colour.

Figure 1.1 shows an example of each of the four possible classes of objects. Each class can be described as follows:

- **Octopus.** The ‘main’ character in the video - the octopus is the only non-static character and its properties change due to its actions. Each frame contains at most one octopus. The initial frame always contains a red octopus with a randomly assigned rotation.
- **Fish.** Fish are always silver, but can have any rotation. When the octopus comes close to the fish, the fish disappears (gets eaten).
- **Bag.** Similarly to fish, plastic bags are always white but can take any rotation. Bags are harmful to the octopus, so both objects disappear when close.
- **Rock.** Rocks can have four colours: brown, blue, purple and green, but always face upright. When an octopus comes near a rock the octopus’ colour will change (if necessary) to match that of the rock (it will be camouflaged).

Objects in each frame are always enclosed by a rectangular box. Object positions are given as $(x1, y1, x2, y2)$, where $(x1, y1)$ is the top left corner of the object, and $(x2, y2)$ is the bottom right².

²The y (vertical) direction is downward increasing.

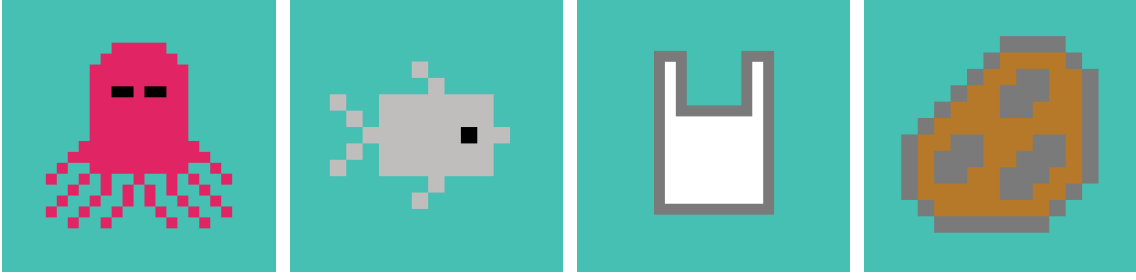


Figure 1.1: Examples of each object type in the videos (not to scale).

The dataset models a single binary relation between objects, ‘close to’. Internally this relation is defined as: object A is close to object B if, after expanding A by 5 pixels on each side, for each pair of parallel edges of A (edges in the horizontal direction and edges in the vertical direction), one of the edges either overlaps with B or is fully contained within B. The ‘close to’ relation is symmetric. The algorithm for determining closeness is outlined in pseudocode in Algorithm 1. The dataset does not consider any relations with an arity higher than two.

Algorithm 1 Determine whether one object is close to another

```

1: procedure CLOSETO(obj1, obj2)
2:   EXPANDEDGES(obj1) ▷ Add 5 pixels to each side
3:   overlapX  $\leftarrow$  obj2.x2  $\geq$  obj1.x2 and obj2.x1  $\leq$  obj1.x1
4:   overlapY  $\leftarrow$  obj2.y2  $\geq$  obj1.y2 and obj2.y1  $\leq$  obj1.y1
5:   for (x, y)  $\leftarrow$  obj2.corners do
6:     matchX  $\leftarrow$  obj1.x1  $\leq$  x  $\leq$  obj1.x2 or overlapX
7:     matchY  $\leftarrow$  obj1.y1  $\leq$  y  $\leq$  obj1.y2 or overlapY
8:     if matchX and matchY then
9:       return True
10:    end if
11:  end for
12:  return False
13: end procedure

```

Other than object properties and relations between objects, which encode visual and spatial information from a single frame, the other key data from the video are the events, which encode temporal information from the frames. Events occur between two consecutive frames, and each timestep can contain multiple (or no) events. Each video therefore contains 31 sets of events.

We choose to split events into two disjoint sets: actions and effects. Actions can be thought of as motions that an object can make to alter its position or rotation. We consider three such actions: move, rotate clockwise and rotate anticlockwise. Rotations have the intuitive effect on an object’s rotation. Move causes the object to move 15 pixels in the direction of its rotation. For example, an object at position (20, 100, 30, 110) with a ‘right-facing’ rotation will be at position (35, 100, 45, 110) after a move action.

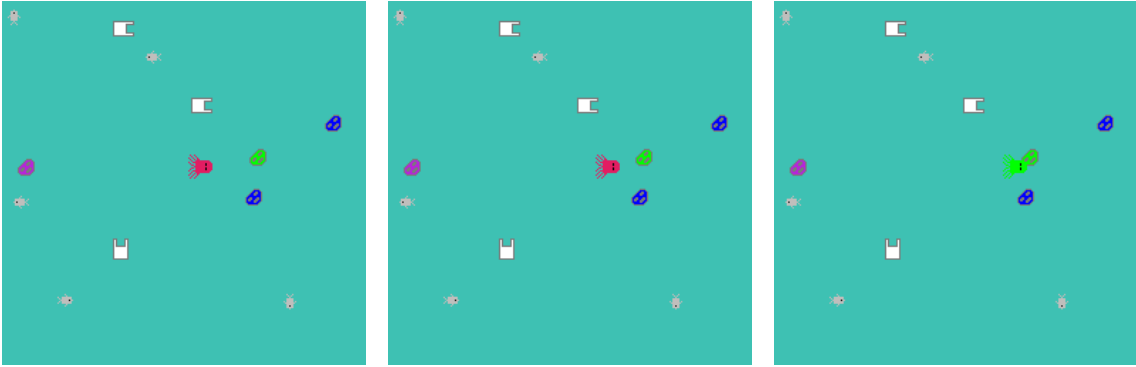


Figure 1.2: An example of an octopus moving close to a green rock and turning green.

Effects, on the other hand, are direct consequences of actions. Effects always alter an object’s state (its class, position, rotation and colour). The dataset contains three effects: change colour, eat a fish and eat a bag. The octopus is the only object which can change colour, and this only occurs when the octopus is close to a rock. As described above, the octopus takes the colour of the rock. The octopus then continues to keep this new colour, even after moving away from the rock. Figure 1.2 shows a snippet from a video where an octopus moves close to a rock and changes colour. A fish or bag is eaten when an octopus moves close, this causes the fish or bag to disappear. However, unlike a fish, when a bag is eaten the octopus also disappears. Examples of all events can be found in Appendix TODO.

To create a video, we first create an initial frame by randomly sampling the number of each object from a set of uniform distributions: $\mathcal{U}(5, 8)$ for fish, $\mathcal{U}(2, 3)$ for bags and $\mathcal{U}(3, 4)$ for rocks, where $\mathcal{U}(l, u)$ refers to a discrete uniform distribution with lower and upper bounds l and u , respectively. Where colours and rotations need to be chosen for objects, these are sampled uniformly from the set of colours and/or rotations that are possible for that type of object. To create the rest of the video, actions for the octopus are randomly sampled for each timestep. We then work out which effects have occurred in each frame and update the properties of each object accordingly. The octopus has a 0.1 probability of choosing to rotate at each timestep. However, if the octopus chooses to move but moving would cause the octopus to be outside the frame, the octopus will rotate instead. Clockwise and anticlockwise rotations are sampled uniformly.

1.2 Questions and Answers

As well as videos, VideoQA datasets must also contain question-answer pairs (QA pairs). Answers to questions ‘label’ part of the video. For example, if a question asks the model to find the event which occurs between two frames, the answer to that question labels the event. However, any event not associated with a QA pair is unlabelled. The shortage of labels is not specific to events; there can be as many as 512 object instances in a single video, but perhaps only a single QA

pair will label an object with a property value. Relations between objects face the same problem. For this reason, VideoQA datasets, unlike many other supervised learning datasets, can contain a lot of sparsely labelled data. On one hand, this may be beneficial to the model, since it only needs to understand the parts of the video which are mentioned in the questions, but, on the other, it can hinder the training of the model, since there is less data than would otherwise be available. Chapter ?? offers one solution to this problem.

The OceanQA dataset contains ten QA pairs per video, sampled randomly from seven question-types. The remainder of this section introduces each question type separately, along with a grammar for each question and answer presented in extended Backus-Naur form (EBNF). Since Natural Language Processing (NLP) is not the focus of this project, we choose to use a small, discrete set of structured question templates, which allow the model to easily extract relevant symbolic information from the questions. Allowing free-form natural language questions creates additional complexity and uncertainty for the model, which we felt would distract from the core focus on hybrid machine learning for VideoQA. Generating these free-form questions can also be very time and resource intensive. A further discussion on possible future work on hybrid models for VideoQA tasks which use free-form questions and answers is contained in Chapter ??.

Equation 1.1 formalises the grammar of the dataset’s objects, relations and events, which was implicitly described in Section 1.1. These EBNF grammar rules are used in the rules for the questions and answers.

$$\begin{aligned}
 \langle \text{object} \rangle &::= [\langle \text{property_value} \rangle] \langle \text{class} \rangle \\
 \langle \text{property_value} \rangle &::= \langle \text{rotation_value} \rangle \mid \langle \text{colour_value} \rangle \\
 \langle \text{rotation_value} \rangle &::= \text{upward-facing} \mid \text{right-facing} \mid \text{downward-facing} \mid \text{left-facing} \\
 \langle \text{colour_value} \rangle &::= \text{red} \mid \text{blue} \mid \text{purple} \mid \text{brown} \mid \text{green} \mid \text{silver} \mid \text{white} \\
 \langle \text{class} \rangle &::= \text{octopus} \mid \text{fish} \mid \text{bag} \mid \text{rock} \\
 \langle \text{property} \rangle &::= \text{rotation} \mid \text{colour} \\
 \langle \text{relation} \rangle &::= \text{close} \\
 \langle \text{event} \rangle &::= \langle \text{action} \rangle \mid \langle \text{effect} \rangle \\
 \langle \text{action} \rangle &::= \text{move} \mid \text{rotate clockwise} \mid \text{rotate anti-clockwise} \\
 \langle \text{effect} \rangle &::= \text{change colour} \mid \text{eat a fish} \mid \text{eat a bag} \\
 \langle \text{frame_idx} \rangle &::= 0 \mid 1 \mid \dots \mid 31
 \end{aligned}
 \tag{1.1}$$

As mentioned above there are seven question types. The first two of these are VQA questions, which only require the model to look at a single frame of the video. The remaining five question types require the model to reason across frames. We adapt the repetition count, repeating action and state transition questions from [3] (discussed in Chapter ??) to the OceanQA environment. We also add two further

video-specific questions: questions about actions between two frames and questions about changing property values. The full set of question and answer templates is as follows:

1. Property questions are designed to test a model’s understanding of object properties. In the training data 41% of the questions ask about colour, while 59% ask about rotation. All property values are represented in the training data, but not uniformly; some property values are very scarce. This reflects the underlying scarcity of objects with particular property values in the dataset; green octopuses, for example, are quite rare, while silver fish are very common. The EBNF grammar for property questions and answers is as follows:

$$\begin{aligned} \langle q_type_1 \rangle &::= \text{What } \langle property \rangle \text{ was the } \langle object \rangle \text{ in frame } \langle frame_idx \rangle? \\ \langle ans_type_1 \rangle &::= \langle property_value \rangle \end{aligned} \quad (1.2)$$

2. Relation questions test a model’s understanding of binary relations between objects. These questions only require a yes-or-no answer. The answer is ‘yes’ for roughly 16% of these questions in the training data. This imbalance reflects the lack of instances of binary relations between objects, since we only model the ‘close’ relation. Objects are selected randomly from a random frame of the video. The grammar for these questions and answers is as follows:

$$\begin{aligned} \langle q_type_2 \rangle &::= \text{Was the } \langle object \rangle \langle relation \rangle \text{ to the } \langle object \rangle \text{ in frame } \langle frame_idx \rangle? \\ \langle ans_type_2 \rangle &::= \text{yes} \mid \text{no} \end{aligned} \quad (1.3)$$

3. Action questions ask which action occurred between two frames of a video. ‘Move’ actions account for around 45% of answers to these questions, while ‘rotate clockwise’ and ‘rotate anticlockwise’ account for approximately 27.5% of questions each. The answer to these questions will never be ‘nothing’. The grammar for action questions and answers is as follows:

$$\begin{aligned} \langle q_type_3 \rangle &::= \text{Which action occurred immediately after frame } \langle frame_idx \rangle? \\ \langle ans_type_3 \rangle &::= \langle action \rangle \end{aligned} \quad (1.4)$$

4. Changed-property questions require the model to reason about how a property of the octopus changes from one frame to the next. The dataset guarantees that only a single (explicit) property changes immediately after $\langle frame_idx \rangle$. Approximately 78% of these questions ask about the colour of the octopus, while the remaining 22% ask about the rotation. The grammar is as follows:

$$\begin{aligned} \langle q_type_4 \rangle &::= \text{What happened to the octopus immediately after } \langle frame_idx \rangle? \\ \langle ans_type_4 \rangle &::= \text{Its } \langle property \rangle \text{ changed from } \langle property_value \rangle \text{ to } \langle property_value \rangle \end{aligned} \quad (1.5)$$

5. Repetition count questions ask the model to work out how many times an event occurs in a given video. This requires the model to be able to count the number of occurrences of an event. Events are sampled from a uniform distribution; if the event never occurs in the video, the answer is simply 0. Since each event can occur at most once per frame-interval, the answer is guaranteed to be between 0 and 30 (inclusive). The grammar for repetition count questions and answers is as follows:

$$\begin{aligned} \langle q_type_5 \rangle &::= \text{How many times does the octopus } \langle event \rangle? \\ \langle ans_type_5 \rangle &::= 0 \mid 1 \mid \dots \mid 30 \end{aligned} \quad (1.6)$$

6. Repeating action questions are similar to repetition count questions, but instead of asking the model for a number they ask the model to find the event which occurs a given number of times. Events cannot be sampled uniformly since, in a given video, multiple events may have the same count. Approximately 86% of questions are about actions, while the remaining 14% refer to effects. This imbalance is again due to the underlying scarcity of particular events in the dataset. There is a unique answer to every question. The grammar for repeating action questions and answers is as follows:

$$\begin{aligned} \langle q_type_6 \rangle &::= \text{What does the octopus do } \langle positive_int \rangle \text{ times?} \\ \langle ans_type_6 \rangle &::= \langle event \rangle \\ \langle positive_int \rangle &::= 0 \mid 1 \mid \dots \mid 30 \end{aligned} \quad (1.7)$$

7. State transition questions ask the model to find the action that occurs after a given event. 78% of the answers to these questions refer to the ‘move’ action, while ‘rotate clockwise’ and ‘rotate anticlockwise’ are the answers to 11% of the questions each. In addition to asking the model to reason temporally about actions and events, these questions also require the model to understand which instance of an event is the ‘nth’ occurrence. The ‘nth time’ section of the grammar is unused if there is only one occurrence of the event in the video. The grammar is as follows:

$$\begin{aligned} \langle q_type_7 \rangle &::= \text{What does the octopus do immediately after} \\ &\quad \langle action_noun \rangle \text{ [for the } \langle nth \rangle \text{ time]}? \\ \langle ans_type_7 \rangle &::= \langle action \rangle \\ \langle action_noun \rangle &::= \text{rotating clockwise} \mid \text{rotating anticlockwise} \mid \\ &\quad \text{eating a fish} \mid \text{eating a bag} \mid \text{changing colour} \\ \langle nth \rangle &::= \text{first} \mid \text{second} \mid \text{third} \mid \text{fourth} \mid \text{fifth} \end{aligned} \quad (1.8)$$

Property, relation and action questions are intended to be used to train and test the model’s understanding of object properties, binary relations between objects and actions between two consecutive frames, respectively. Since the property and relation questions contain the ‘ $\langle object \rangle$ ’ rule, knowledge of property values may be required to select the corresponding object from the frame. This means it would

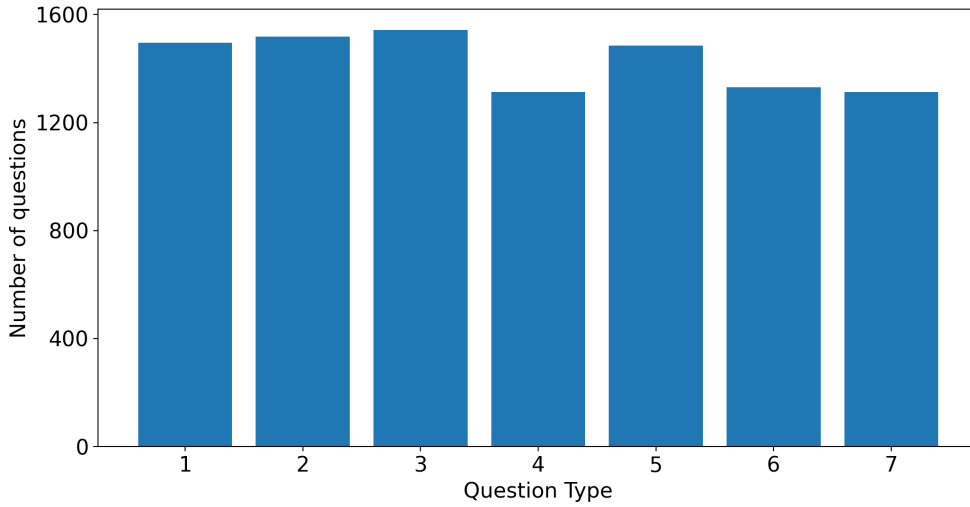


Figure 1.3: The number of QA pairs in the training data for each question type.

be helpful to train the model’s understanding of object properties before relations. Selecting the training data for object properties is, however, non-trivial. Chapter ?? discusses one possible solution to this problem. The final four question types are included to diversify the data the model can be trained on and to evaluate how well the model can combine its understanding of object properties, relations and events. None of the questions require any sort of external knowledge (other than the ability to count).

Questions are generated independently for each video, and questions are sampled uniformly from the seven question templates given above. This leads to an approximately uniform distribution of question types in the dataset, as shown in Figure 1.3. Since there are ten questions in each video the training, validation and testing datasets contain 10000, 2000 and 2000 questions, respectively. Appendix TODO contains a number of example images and QA pairs from the dataset.

1.3 Specification

Unlike end-to-end neural network models, hybrid models require knowledge that has been extracted from a video to be made explicit. For example, if a neural model is shown a video and asked for the colour of the octopus, it will first extract features from the video and then encode the video into a vector. This vector is combined with the encoding of the question, before a final section of the model produces either a short, free-form sentence or a probability distribution over possible answers. Relevant information from the video, including the colour of the octopus, is implicitly included in the vector encoding of the video. Hybrid models, however, usually require that this information is symbolic so that logical reasoning can be employed to find an answer to the question. For this reason it is necessary

to define a specification of the video environment, which outlines the set of concepts which the model must learn in order to accurately answer the questions. This specification can also provide a way to inject background knowledge about the environment into the model.

Each environment specification must contain the following information:

1. The number of frames in each video. This project assumes this to be fixed, but relaxing this assumption would not require any major changes to the construction of the dataset or the model outlined in the rest of this report.
2. A set of pairs, $\langle class, is_static \rangle$. Each element of the set contains the type of an object and a boolean corresponding to whether the object is capable of performing an action. All (relevant) object classes must be mentioned in this set.
3. A set of pairs, $\langle property, \{ property_value \} \rangle$, corresponding to a set of properties along with a set of their respective values. This representation for object properties assumes that each property has a discrete, finite set of values. Continuous properties are not considered in this project. Object class and position are assumed to be implicit properties since all objects must have a value for them. These implicit properties should not be listed here.
4. A set of binary relations.
5. A set of actions.
6. A set of effects of actions.

Although it has been described in detail already, Equation 1.9 provides the formal environment specification for the OceanQA dataset.

$$\begin{aligned}
 frames &::= 32 \\
 objects &::= \{ \langle octopus, false \rangle, \langle fish, true \rangle, \langle bag, true \rangle, \langle rock, true \rangle \} \\
 properties &::= \{ \langle colour, \{ red, blue, purple, brown, green, silver, white \} \rangle, \\
 &\quad \langle rotation, \{ upward-facing, right-facing, downward-facing, left-facing \} \rangle \} \\
 relations &::= \{ close \} \\
 actions &::= \{ move, rotate clockwise, rotate anticlockwise \} \\
 effects &::= \{ change colour, eat a fish, eat a bag \}
 \end{aligned}
 \tag{1.9}$$

As alluded to already, an environment specification, like the one shown in Equation 1.9, requires that objects, properties, relations, actions and effects are all discrete. Continuous variables could be modelled by discretising, but this may lead to lower accuracy and a large number of possible values. This may have to be treated as a trade-off for using hybrid models. On the other hand, this is one of the first pieces of work which explores the use of hybrid models in VideoQA tasks; future work may be able to overcome this problem.

Chapter 2

Hybrid Property, Event and Relation Learner

We outline a novel paradigm for solving the VideoQA problem. This approach merges deep learning and logic-based machine learning and inference methods in an attempt to learn the concepts required to answer the questions. We name this approach: Hybrid Property, Event and Relation Learner (H-PERL). Section 2.1 outlines the structure of an H-PERL model, while Section 2.2 discusses the implementation of a number of H-PERL components which are common to all models presented in the subsequent chapters.

2.1 H-PERL Model

H-PERL is a generic, pipelined architecture for VideoQA tasks. The pipeline is composed of a number of components which, when strung together, form a model. A model is, therefore, a set of component implementations, and can be thought of as an ‘instance’ of H-PERL. Chapters ?? and ?? each outline an H-PERL model for the OceanQA dataset.

2.1.1 Architecture

The H-PERL architecture assumes that all of the information in an environment which is required to answer the questions can be modeled using: objects; binary relations between objects; and events (which occur due to an object) between two consecutive frames in the video. For many environments, simple environments (like our OceanQA dataset) in particular, this assumption holds. However, for many VideoQA datasets, particularly those set in the real-world, this assumption may not be suitable. For example, extracting some objects from a video may be non-trivial (as in the case of abstract nouns), and yet information on these objects may still be required to answer the questions. Additionally, H-PERL is not capable of modelling relations between objects with an arity larger than 2.

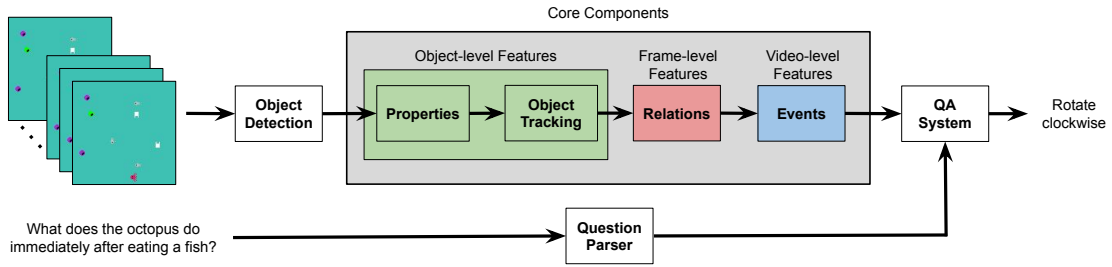


Figure 2.1: An overview of the H-PERL architecture for VideoQA. Green, red and blue shading indicates components which work to extract features at different levels of abstraction. Grey shading indicates the ‘core’ components of the architecture.

Figure 2.1 shows the components involved in a typical H-PERL model. During evaluation, information from the video and the question flow, from left to right, through the pipeline. Each H-PERL model assumes that the object detection, question parsing and QA system components are “pre-made” (either pre-trained or manually engineered). We refer to these as ‘non-core’ components. H-PERL allows the remaining ‘core’ components to be updated as the model is trained, although they don’t necessarily have to be. Each core component in the pipeline accumulates information. This means that each component guarantees that existing information (or features of the data) will not be overwritten (with the small exception of the event component when error correction is used, discussed further in Chapter ??). As shown in Figure 2.1 components in the pipeline work at different levels of abstraction; the object properties and tracking components work to extract object-level features, while the relations and events components work to extract frame and video-level features, respectively.

The following is a high-level description of the tasks each component is required to complete for the H-PERL architecture to work with high accuracy:

1. **Question Parser.** The QA parsing component is used to extract relevant pieces of information from the questions (and answers when training). For example, given the question: “What does the octopus do immediately after eating a fish?” and that the question is a state-transition question, the parsing component would extract, firstly, that the object in question was an octopus, and secondly, that the event was ‘eat a fish’. The parsing component, therefore, bridges the gap between the symbolic data, which the model works with, and the natural language questions and answers. When an H-PERL model is being evaluated, only the question needs to be parsed. However, when the model is training this component acts as a question-and-answer parser, since the model requires that the feedback that comes from the answer is also in symbolic form.
2. **Object Detector.** The detection component produces bounding boxes and classes for each object in each frame of the video. Any object not detected at this stage of the pipeline is assumed to be part of the background and is therefore ignored by the rest of the model.

3. **Property Extractor.** Given a set of images of objects from the detection component, the property extraction component assigns a value to every property listed in the environment specification for every object in the set.
4. **Object Tracker.** The object tracker is required to assign an identifier to each object in a given video. The object identifiers assigned in the initial frame of the video can be arbitrary, but then an algorithm is usually applied inductively to each remaining frame of the video in an attempt to assign each object the same identifier as it was given in the previous frame.
5. **Relation Classifier.** The job of the binary relation classifier is, given a symbolic representation of a video, to list all of the instances of binary relations between objects in the video. The set of possible binary relation is defined in the dataset's environment specification.
6. **Event Detector.** The event detection component produces a set of events for each pair of consecutive frames in the video. Each set of events can contain both actions and effects, and each event consists of an event name and an object identifier which signifies the non-static object that took part in that event.
7. **QA System.** The job of the QA system is to take all of the features which have been accumulated by previous components in the pipeline, along with a parsed question, and produce an answer to the question. The QA system is also given the question type, this allows it to make sense of the parsed question and to apply different reasoning for each question type.

While VideoQA tasks may require the pipeline to contain the above set of components, the H-PERL architecture can also be modified for solving related problems. For example, by removing the tracking and events components and keeping the rest of the pipeline the same, a modified version of the H-PERL architecture could be used to solve VQA (also known as FrameQA) tasks. In fact, the architecture could be modified to work with many different input types, including images, videos, text or speech, provided that the data contains something akin to an 'object', and that extracting these objects (and features from these objects) is feasible. We could also swap out the final component in the pipeline, the QA system, for another task-specific component. For example, we could create a video (or image) captioning architecture by using a captioning component rather than a QA component. Figure 2.2 shows a number of architectures similar to H-PERL that have been modified for other tasks.

At a high level, then, the H-PERL approach consists of three main stages:

1. Attention restriction
2. Feature extraction
3. Task-specific output generation

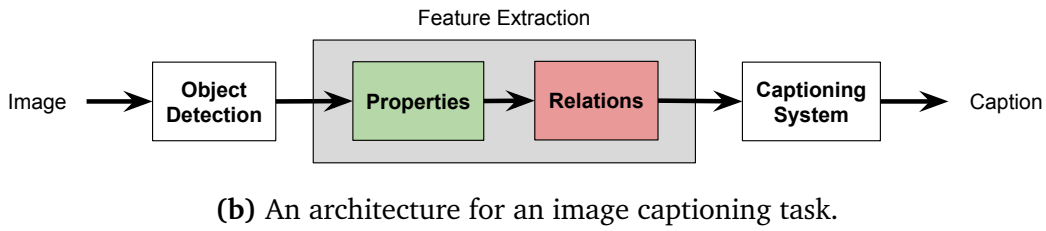
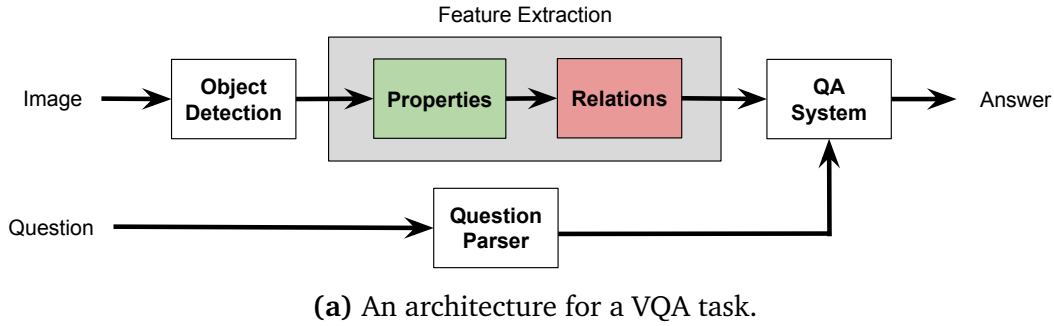


Figure 2.2: Potential modifications to the H-PERL architecture for solving other tasks. Grey shading indicates feature extraction components. Green and red shading indicate components which extract object-level and frame-level features, respectively.

These stages are similar to some approaches used in VQA tasks [8, 9], where the model’s attention is firstly restricted to specific parts of the input (in this case objects), before features are extracted from the input and an output is generated. However, unlike other approaches, H-PERL does not dynamically update the model’s attention based on the input question. Another key difference between H-PERL and other question-answering approaches is that extracted features are stored symbolically. This gives the observer a better understanding of what the model is learning and allows the injection of background or commonsense knowledge directly into the model.

2.1.2 Information Representation

The first component in the H-PERL pipeline, the object detector, takes a raw video as input and produces a set of bounding boxes (object positions) and object classes for each frame of the video. After object detection has been applied, data about the video is stored symbolically for each object in each frame of the video. Each of the subsequent components in the pipeline can access this symbolic data along with the raw image of each object (the raw video frames are discarded). Each of these components then accumulates symbolic information about the objects, frames or video. This is what is meant by extracting object, frame and video-level features.

Once the symbolic features have been extracted from the input, each of the components in the architecture need an agreed upon way of representing the information. This symbolic representation is as follows:

- For an object with identifier $\langle id \rangle$ in frame $\langle frame \rangle$, the object's properties, rotation and class (all referred to as $\langle property \rangle$), each with value $\langle value \rangle$, are represented as follows:

$$\text{obs}(\langle property \rangle(\langle value \rangle, \langle id \rangle), \langle frame \rangle) \quad (2.1)$$

As described in Chapter 1, the value of an object's rotation is given as $(x1, y1, x2, y2)$, where $(x1, y1)$ is the top left corner of the object, and $(x2, y2)$ is the bottom right.

- For two objects in frame $\langle frame \rangle$, with identifiers $\langle id1 \rangle$ and $\langle id2 \rangle$, a binary relation, $\langle relation \rangle$, between the objects is represented as follows:

$$\text{obs}(\langle relation \rangle(\langle id1 \rangle, \langle id2 \rangle), \langle frame \rangle) \quad (2.2)$$

- An event which occurs immediately after frame $\langle frame \rangle$, due to an object with identifier $\langle id \rangle$, is represented by one of the following rules, depending on whether the event was an action or effect:

$$\text{occurs_action}(\langle action \rangle(\langle id \rangle), \langle frame \rangle) \quad (2.3)$$

Using these representations an entire video can be encoded into a logic program. The final component in the pipeline, the QA system, takes as input a video encoding and a set of parsed questions, and constructs a set of ASP rules which are used, along with the video encoding, to find an answer to the question.

2.1.3 Requirements

To give a complete view of what is needed to construct an H-PERL model, we outline the minimum set of requirements and a number of assumptions. Each instance of an H-PERL model may require additional constraints to be applied on top of these. The set of requirements is as follows:

1. A set of pre-made, non-core components (question parser, object detector and QA system).
2. A VideoQA dataset, where each element is of the form:

$$\langle \text{video}, \{ \langle \text{question}, \text{answer}, \text{question type} \rangle \} \rangle$$

3. Environment specification for the given dataset.
4. (Optionally) Background knowledge of the environment, written in ASP.

H-PERL also requires that the following assumptions be made about the data:

1. All relevant information in the video can be modelled by object properties; binary relations between objects; and events occurring between consecutive frames of the video.
2. Each of the properties, relations and events components can be trained individually and directly using a specific type of question. For example, for the OceanQA dataset, we can use question types 1, 2 and 3 to train the properties, relations and events components, respectively.
3. As mentioned in Chapter 1, properties, relations and events must be discrete. There is also currently no way of modelling continuous variables in the data; for example, we cannot say that the octopus rotated clockwise by $\frac{1}{4}\pi$ radians.

These requirements and assumptions clarify some of the limitations of H-PERL models. Firstly, for some environments pre-trained object detectors (or data to train them with) may be difficult to find. Secondly, it may also be difficult, if not impossible, to construct or train a question parser for free-form, natural language question-answering datasets. Additionally, most QA datasets are not guaranteed to contain questions which can be used to directly train components of the model. Finally, many environments will simply be too complex to be accurately modelled by discrete properties, relations and events. These requirements and assumptions do, however, provide initial directions for future research in the area of Hybrid question-answering models. We discuss some potential extensions to H-PERL in Chapter ??.

As well as drawbacks, the H-PERL architecture does also provide a number of advantages. One of the most important advantages of hybrid models is the ability to encode commonsense knowledge or background knowledge of the environment directly into the model without having to learn it. Since the model accumulates video information symbolically, we can inject background knowledge at any stage of the pipeline. A number of further advantages of hybrid models are presented in the subsequent chapters, and these, along with the disadvantages, are summarised in Chapter ??.

2.2 Common Components

The two models, one hardcoded and one trained, which are outlined in subsequent chapters, contain a number of common components. Before describing the implementation of these models, we first outline the details of the components common to both. These components include the non-core components, which the H-PERL pipeline assumes are pre-made, as well as the object tracker, which uses the a similar algorithm for both the hardcoded and trained models.

2.2.1 Question Parser

As mentioned previously, the role of the question parsing component is to translate the natural language questions into a set of keywords (and their types). This extracted information is then used by the QA system to construct ASP rules with which the answers to the questions are found. During training, the question parser acts as a question-and-answer parser, as the model needs to collect training data from the answer.

Both the hardcoded and the trained H-PERL models use a very simple hand-engineered question parser. Since the question parser is told the type of the question it has been given, and since the OceanQA dataset uses templated questions, the question parser can simply extract the relevant parts of the question by looking at the corresponding template. This is implemented by splitting the question into a list of words and simply picking the correct word based on the template.

Take the question “What does the octopus do immediately after rotating clockwise for the second time?” as an example. Based on the question template shown in Chapter 1, the parsing component would extract the following from the question:

$$\begin{aligned} \text{object} &= \text{octopus} \\ \text{action} &= \text{rotate clockwise} \\ \text{occurrence} &= 2 \end{aligned} \tag{2.4}$$

Note that the action noun ‘rotating clockwise’ has been converted to the verb form, as given in the environment specification. In addition, the occurrence, ‘second’, is converted to a number, which is simpler for the QA system to work with. These conversions are all hand-engineered, as opposed to learnt.

The parsing of all other question types is done in the exactly same way. For the sake of brevity these are not shown. As a general rule, the information extracted from the questions is the parts of the question templates which are contained in $\langle \rangle$ brackets. For question types 4, 5, 6 and 7 the parser must also extract the *object* keyword, which is guaranteed to be octopus in all four of these question types.

2.2.2 Object Detector

Chapter ?? describes two object detection algorithms: Faster R-CNN and YOLO. Faster R-CNN has been shown to be the more accurate of the two [5]. YOLO’s key advantage, however, is its speed; it has been shown to be capable of processing five times as many frames per second as Faster R-CNN on the Pascal VOC dataset [6]. However, since OceanQA frames are fairly small and simple, and since we don’t need real-time performance for VideoQA tasks, we prefer accuracy over speed and therefore choose Faster R-CNN over YOLO for the object detection component.

The object detection network uses a pre-built Faster R-CNN model from the *Torchvision*¹ library. The network is trained on the full-data OceanQA dataset (where every object in every frame is fully labelled, as outlined in Chapter 1), using Faster R-CNN’s multi-task loss function [7]. The *PyCOCOTools*² library is used to implement the core training functionality.

As is standard in Faster R-CNN implementations, we provide the detection model with a feature extraction network. This network’s role is to extract features from the input frames and pass these to the pre-built detection network. We provide a three-layer convolutional network, described in Table 2.1, with randomly initialised weights. These weights, along with the weights of the detection network, are updated as the whole model is trained.

The full details of the object detector’s performance are given in Chapter ??

Layer	In Channels	Out Channels	Kernel Size	Stride	Padding
Conv1	3	32	3x3	1	1
Conv2	32	64	3x3	2	1
Conv3	64	128	3x3	2	1

Table 2.1: Specification of the feature extraction network. Batch normalisation is applied between layers Conv1 and Conv2, and between layers Conv2 and Conv3.

2.2.3 Object Tracker

The object tracker’s task is to assign every object in the video with an identifier. Ideally, an object is assigned the same identifier throughout the video, however, errors in the object detection can make this difficult. The tracker is given a video containing a list of frames, with each frame containing a set of detected objects.

We implement a simple object tracker which assigns integer identifiers, starting from zero, to each object in the initial frame of the video. The tracker then inductively applies an algorithm (discussed below) which attempts to match objects in the previous frame (we call these previous objects), which have been assigned an identifier, with objects in the next frame (next objects), which are unidentified. An object, *obj1*, matches with another object, *obj2*, if the following two conditions are met:

1. *obj1* and *obj2* have the same class.
2. The Euclidean distance between the top left corner of *obj1* and the top left corner of *obj2* is no more than 30 pixels.

Both the distance metric and the maximum distance value are considered hyperparameters of the tracking component. Many other values for these hyperparameters

¹Available at: <https://github.com/pytorch/vision>

²Available at: <https://github.com/cocodataset/cocoapi>

would work equally well with the OceanQA dataset, and if these values weren't known in advance, they could be found fairly easily through a combination of analysis of the dataset (the maximum distance must be at least 15 pixels) and trial-and-error.

After assigning identifiers to objects in the initial frame, the tracker iteratively applies the following algorithm to each subsequent frame:

1. Each object in the new frame votes for the object in the previous frame which matches with it best.
2. Each previous object collates all of the new objects which have voted for it and chooses the new object with which it best matches.
3. If a previous object does not get any votes, it is because there are no new objects which consider this object their best match. We therefore assume the previous object has disappeared. Each disappeared object is added to a *hidden objects* list, so that if the object reappeared in a subsequent frame the same identifier could be assigned.
4. Each remaining previous object is now matched with a single new object. The new object is assigned the identifier of the previous object it is matched with.
5. Some new objects - namely the which were rejected by a previous object in favour of another new object - remain unidentified. Since all previous objects have either passed on their identifier or disappeared, these remaining objects must be new. For each of these new objects, if it cannot be matched with an object in the *hidden objects* list, a previously unused identifier is created and assigned to the object, otherwise the object is assigned its best match in the *hidden objects* list.
6. Each object in the *hidden objects* list can only stay for a maximum of five frames from when they are added (although this duration is also considered a hyperparameter and could therefore be altered). The final step of the algorithm is to remove from this list any hidden object which has overstayed its welcome.

This simple, heuristic-based tracker works well for simple, well-defined environments such as OceanQA. In fact, Chapter ?? shows that the tracker can achieve perfect performance, assuming the object detection is completely accurate. However, this tracker may perform poorly in more complex environments, such as real-world VideoQA datasets. Other, more advanced object tracking algorithms may fare better, however. The Simple Online and Realtime Tracking (SORT) [1] algorithm is one of the best performers [2]. Since the H-PERL architecture is completely modular, our tracking algorithm could simply be swapped out for a more advanced alternative when the environment required it.

2.2.4 QA System

Bibliography

- [1] Alex Bewley et al. “Simple online and realtime tracking”. In: *2016 IEEE International Conference on Image Processing (ICIP)*. IEEE. 2016, pp. 3464–3468.
- [2] Gioele Ciaparrone et al. “Deep learning in video multi-object tracking: A survey”. In: *Neurocomputing* (2019).
- [3] 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.
- [4] Jiayuan Mao et al. “The neuro-symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision”. In: *arXiv preprint arXiv:1904.12584* (2019).
- [5] Joseph Redmon and Ali Farhadi. “Yolov3: An incremental improvement”. In: *arXiv preprint arXiv:1804.02767* (2018).
- [6] 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.
- [7] 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.
- [8] Kelvin Xu et al. “Show, attend and tell: Neural image caption generation with visual attention”. In: *International conference on machine learning*. 2015, pp. 2048–2057.
- [9] Zichao Yang et al. “Stacked attention networks for image question answering”. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2016, pp. 21–29.
- [10] Kexin Yi et al. “Clevrer: Collision events for video representation and reasoning”. In: *arXiv preprint arXiv:1910.01442* (2019).
- [11] Kexin Yi et al. “Neural-symbolic vqa: Disentangling reasoning from vision and language understanding”. In: *Advances in Neural Information Processing Systems*. 2018, pp. 1031–1042.