

Chapter 1

OceanQA 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 [3, 2, 4] 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 [3] 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.

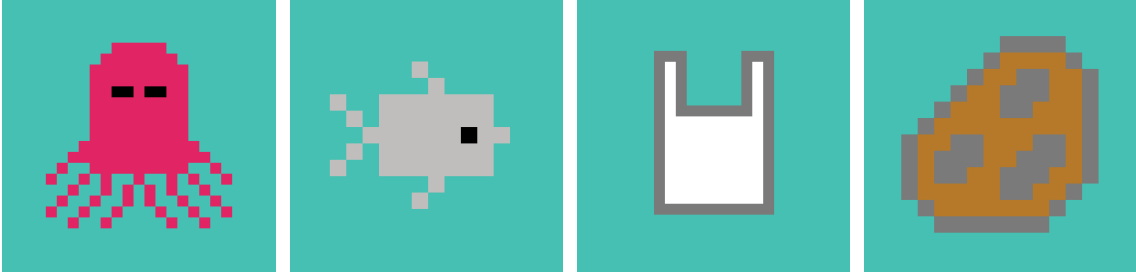


Figure 1.1: Examples of each object type in the videos (not to scale).

The dataset models a single binary relation between objects, *close*. 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’s rectangular box. The *close* 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.

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, also known as primitive events, can be thought of as motions that an object can make to alter its state (an object’s class, position, rotation and colour). 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 moving.

Algorithm 1 Determine whether one object is close to another

```

1: procedure CLOSE(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

```

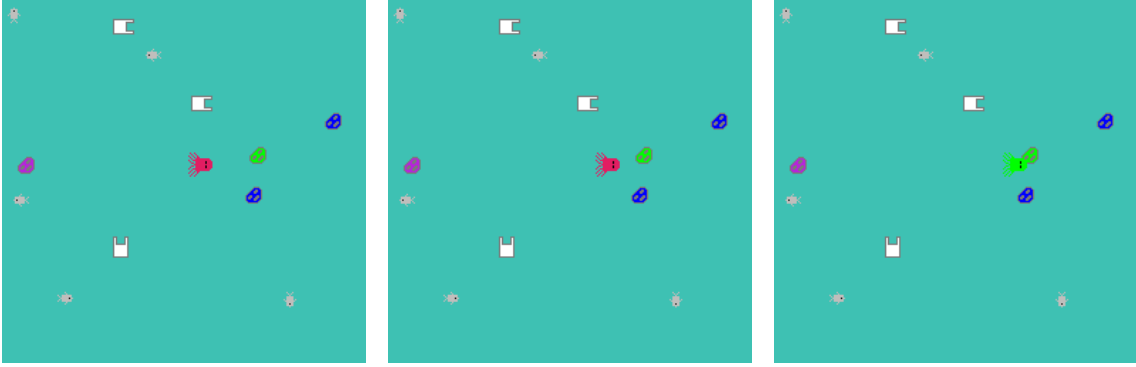


Figure 1.2: An example of an octopus moving close to a green rock and turning green.

Effects, on the other hand, happen as a consequence of some frame state being true; they are triggered by a particular frame state. Hence, effects are also known as triggered events. 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 ??.

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 2 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 [1] (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 62% of the questions ask about colour, while 38% 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 15% 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 44% of answers to these questions, while ‘rotate clockwise’ and ‘rotate anticlockwise’ account for approximately 28% 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 79% of these questions ask about the colour of the octopus, while the remaining 21% 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 84% of questions are about actions, while the remaining 16% 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. 75% of the answers to these questions refer to the *move* action, while *rotate clockwise* and *rotate anticlockwise* are the answers to 12.5% 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:

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

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 an $\langle object \rangle$, knowledge of property values may be required to select the corresponding object from the frame in order to compute the

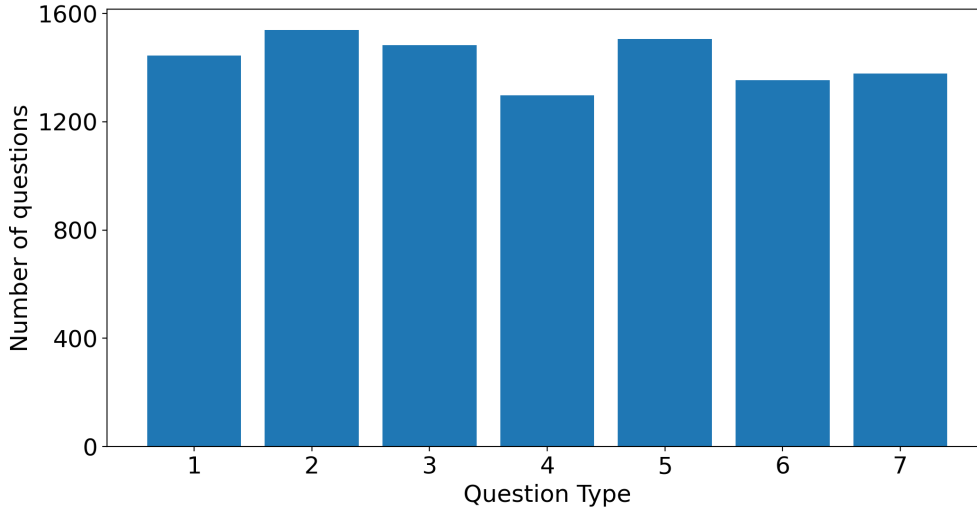


Figure 1.3: The number of QA pairs in the training data for each question type.

answer. This means it would 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 2 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 ?? 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

Trained Model

In contrast to the hardcoded model, the trained H-PERL model does not use manually engineered relations or events components, and must therefore rely on using components which can be trained. The trained model also uses the QA-data version of the OceanQA dataset, rather than the full-data version which the hardcoded model was able to use to train its properties component. This means that the trained model needs to rely on the data contained in QA pairs alone to train its components.

The model is trained using *curriculum learning* - the model learns concepts incrementally. In the case of H-PERL, this means that the model is, firstly, trained to understand object properties, then trained to understand relations, and, finally, is trained to understand events. After training a component, the component is used to extract its corresponding symbolic information from the training data, which can then be used by the next component to be trained.

2.1 Properties

As mentioned in Chapter 1, property questions in the OceanQA dataset ask the model to find a property value for a specific object. This object, however, can contain a reference to a property value. This means that, in some cases, knowledge of object properties is required in order to find the specified object in the frame. For example, if a question asked “What colour was the upward-facing fish in frame 12?”, and there three fish, each with unique rotations, in frame 12, one would need knowledge of object properties in order to select the correct image of the fish. This means the training data cannot be collated in the same way as the hardcoded model; the model needs a trained property extractor in order to find the images to train the property extractor with.

In this section we propose a solution to overcome this problem which utilises semi-supervised learning in order to label all of the objects in the dataset with property values. Once these labels have been found, the property component can be trained in the same way as the hardcoded model, outlined in Chapter ??.

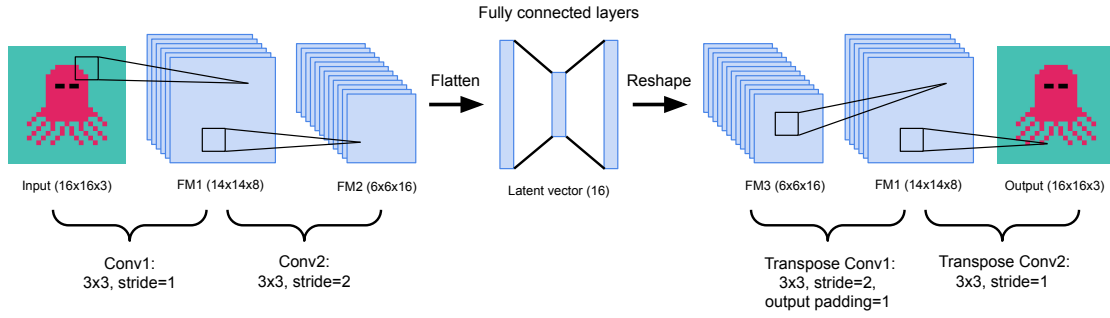


Figure 2.1: An illustration of the autoencoder architecture. Although not shown here, the network contains batch normalisation layers between each pair of convolution layers. FM stands for feature maps, and the dimensions of each feature maps are given as (*height, width, number of feature maps*).

The first step of the algorithm for finding these labels involves training an autoencoder neural network to extract a 16-dimensional latent vector from each object image. This network is trained in an unsupervised manner using a sample of 40,000 of the objects detected by the object detector in the training data, where each object type is equally represented in the sample. The architecture of the network is shown in Figure 2.1. Each object image is resized to 16x16 pixels and the network is trained with a learning rate of 0.001 for 5 epochs with a mean-absolute error (MAE) loss function.

The autoencoder allows the component to work with the latent vectors of objects rather than raw object images. At this stage, every object in the training data is encoded using the autoencoder and stored in vector form. After the autoencoder has been trained and all the objects have been encoded, the properties component splits objects in groups based on their class. The component then proceeds to individually apply the following three high-level steps, which are described in further detail in the sections below, to each object type, t_i :

1. The objects with type t_i are clustered using their latent encoding and each cluster is assigned an integer identifier, c_j . More detail on the clustering is outlined in Section 2.1.1.
2. The QA pairs will label a number of the objects in the training data with some or all of their property values. ASP can be used to find a mapping from each c_j to a set of property-value pairs. The ASP program used to conduct this mapping is described in Section 2.1.2.
3. Once this mapping is known, all of the objects of type t_i in the training data can be labelled. The algorithm for finding these labels is outlined in Section 2.1.3.

As mentioned above, when the property value labels have been found for all objects of all types, the property component can be trained using the same method as the hardcoded model, discussed in Chapter ??.

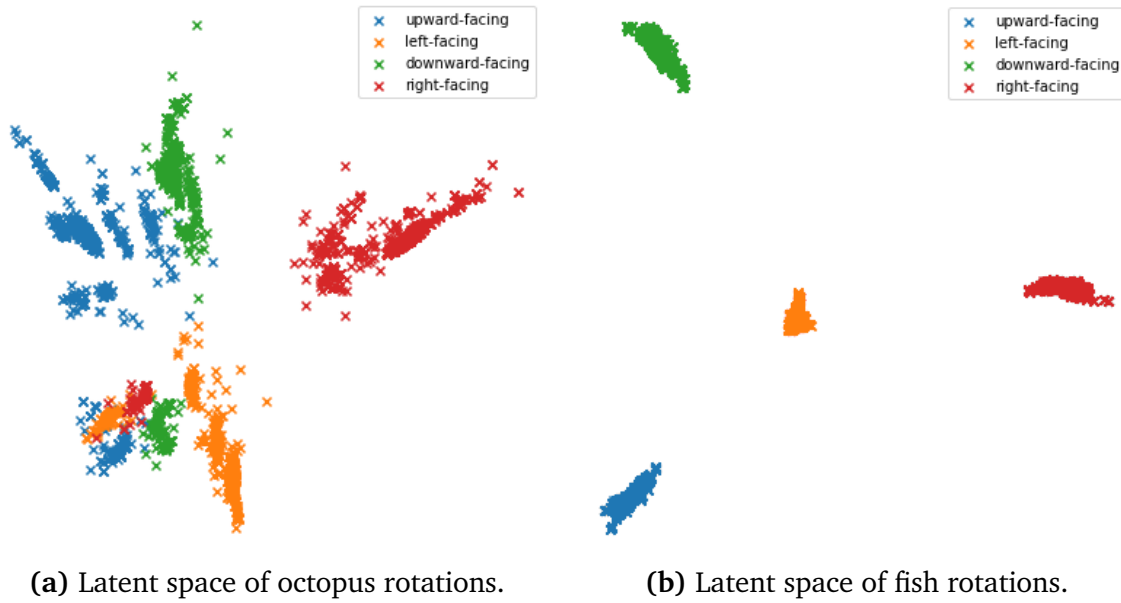


Figure 2.2: 2-dimensional representation of the latent space of octopus and fish images. The colours show different rotation values. Note, these property values come from the full-data version of the OceanQA dataset, not from clustering. Although not shown here, rocks and bags follow a very similar pattern to fish.

2.1.1 Clustering

An analysis of the full-data version of the OceanQA dataset (where objects come fully labelled with their properties) shows that the information encoded in the latent vector of each object separates objects into distinct groups based on their properties. This makes clustering a strong candidate for explicitly separating the objects into distinct groups. The Principle Component Analysis (PCA) projection of the latent vectors for octopus and fish into 2-dimensions is shown in Figure 2.2.

Clustering of the object latent vectors is done using the *Agglomerative Clustering* algorithm implemented by the *SciKitLearn*¹ library. While many other clustering algorithms are available, Agglomerative Clustering was found to work efficiently with a large number of samples. Unlike some other clustering algorithms, however, Agglomerative Clustering requires that the number of clusters to be produced is known beforehand.

Object Type	#Colours	#Rotations	#Clusters
Octopus	5	4	20
Fish	1	4	4
Bag	1	4	4
Rock	4	1	4

Table 2.1: Estimates of the number of clusters required for each class, calculated by multiplying the number of colour and rotation values found in the property questions.

¹Available at <https://scikit-learn.org/stable/modules/clustering.html>

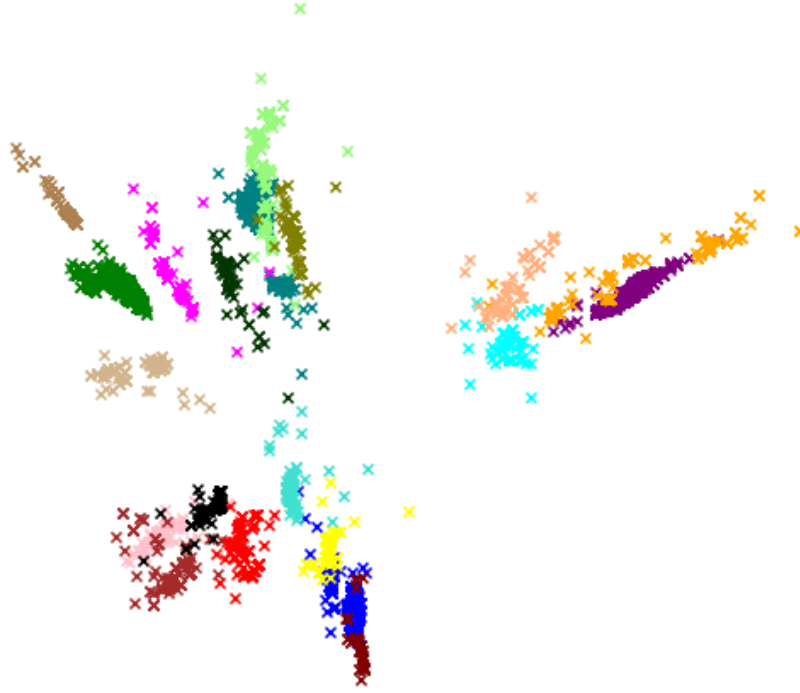


Figure 2.3: Output of the Agglomerative Clustering algorithm for the octopus images.

In order to calculate the number of clusters required for each class, we analyse the number of property values that are linked with that class in the QA pairs. For example, fish only take a single colour, silver, and so only this colour will be mentioned in the property questions that ask about fish. Additionally, fish can take four rotations, and so we expect that all four of these rotations will be mentioned in the questions. In total this gives us four possible property value combinations for fish. We repeat this process for all four object types and list the results in Table 2.1.

In theory, it is possible that this method of calculating the number of clusters required may underestimate, since not all property values may be mentioned in the QA pairs. In practice, however, this makes little difference because if a particular property value is not represented in the QA pairs it is likely to be very rare in the videos and in the evaluation questions.

Figure ?? shows the result of the clustering algorithm applied to the octopus images, projected into two dimensions. After clustering has been completed for a given class, we assign each cluster (each different colour shown in Figure 2.3) an integer identifier. The data collection problem outlined at the start of this chapter has now been reduced to finding the mapping from each cluster identifier to a set of property-value pairs, which optimises some objective function, for each object type t_i . For example, for a cluster c_j , we want to find the optimal values for colour and rotation which correspond to c_j . We denote the optimal mapping from the set of cluster identifiers C to the set of sets of property-value pairs P , for object type t_i as $f_{t_i}^* : C \rightarrow P$.

2.1.2 Property Value Optimisation

In order to find $f_{t_i}^*$, we must first define the objective function. Since we have a set of property questions and answers, we choose to maximise the number of questions answered correctly. It is now necessary to find the questions which correspond to object type t_i . We also use the question and answer parsing component (outlined in Chapter ??) to extract the relevant symbolic information from the question and from the answer. We denote the information extracted from QA pair k as (p_k, v_k, a_k) , where p_k is the property mentioned in the question (or *Null* if there isn't one), v_k is an optional property value mentioned in the question and a_k is the answer, which is a property value.

ASP is then used to conduct the optimisation. One ASP program is constructed for each object type. We outline the ASP program for the object type t_i in the following five stages:

1. Firstly, the search space for the optimisation is constructed using choice rules. For each cluster identifier c_j , the following two choice rules are added:

$$1\{colour_mapping(c_j, col_1); \dots; colour_mapping(c_j, col_n)\}1. \quad (2.1)$$

$$1\{rotation_mapping(c_j, rot_1); \dots; rotation_mapping(c_j, rot_m)\}1. \quad (2.2)$$

Where the set of possible colour values is: $\{col_1, \dots, col_n\}$ and the set of possible rotation values is: $\{rot_1, \dots, rot_m\}$.

These choice rules generate one answer set for each possible combination of property-value pairs. Notice that no attempt has been made to restrict the property-value pairs to those which were used to select the number of clusters, since the ASP optimisation will naturally choose the property-value pairs which answer the most questions correctly.

2. Secondly, for each QA pair, the object data from the frame given in the question is added to the program. Only the objects with type t_i are added. For example, if a question says “What colour was the leftward-facing fish in frame 23?” and three fish have been detected in frame 23, then only the information for those three fish are listed in the program (fish from other questions will of course be listed in the same program).

Each of the objects to be added will be given an integer identifier, this identifier is unrelated to the identifier assigned by the tracking component. Each object is also listed with the cluster identifier that it was assigned to. For each object listed in QA pair $\langle k \rangle$, which is part of cluster c_i and has identifier id , the following predicate is added to the ASP program:

$$obj(id, c_i, \langle k \rangle). \quad (2.3)$$

3. A number of helper ASP rules which convert object data in the form of the *obj* and *<prop>-mapping* predicates into the *holds* predicate, which is used

to answer the questions. These two helper rules are as follows:

$$\text{holds}(\text{colour}(\text{Val}, \text{Id}), Q) :- \text{obj}(\text{Id}, \text{Cluster}, Q), \quad (2.4)$$

$$\text{colour_mapping}(\text{Cluster}, \text{Val}).$$

$$\text{holds}(\text{rotation}(\text{Val}, \text{Id}), Q) :- \text{obj}(\text{Id}, \text{Cluster}, Q), \quad (2.5)$$

$$\text{rotation_mapping}(\text{Cluster}, \text{Val}).$$

Rules 2.4 and 2.5 collate the data from the object's cluster and from the property value that that cluster has been assigned and convert this data into *holds* form. The *holds* predicate can then be used to answer the questions.

4. Next, the ASP rules corresponding to the questions and the ASP facts corresponding to the answers to those questions are added to the program. We firstly find the tuples (p_k, v_k, a_k) which correspond to the object type t_i are extracted by the question and answer parsing component from the $\langle k \rangle^{th}$ QA pair. The following rule and fact are then added if v_k is not *Null*:

$$\text{answer}(\langle k \rangle, p_k, V) :- \text{holds}(p_{v_k}(v_k, \text{Id}), \langle k \rangle), \quad (2.6)$$

$$\text{holds}(p_k(V, \text{Id}), \langle k \rangle), \text{obj}(\text{Id}, -, \langle k \rangle).$$

$$\text{answer}(\langle k \rangle, p_{v_k}, V) :- \text{holds}(p_{v_k}(V, \text{Id}), \langle k \rangle), \quad (2.7)$$

$$\text{holds}(p_k(v_k, \text{Id}), \langle k \rangle), \text{obj}(\text{Id}, -, \langle k \rangle).$$

$$\text{expected}(\langle k \rangle, p_k, a_k). \quad (2.8)$$

$$\text{expected}(\langle k \rangle, p_{v_k}, v_k). \quad (2.9)$$

Where p_{v_k} is the property which corresponds to property value v_k .

For example, the question “What colour is the downward-facing bag in frame 7?” and answer “white”, would be converted into the following:

$$\text{answer}(\langle k \rangle, \text{colour}, V) :- \text{holds}(\text{rotation}(\text{down}, \text{Id}), \langle k \rangle), \quad (2.10)$$

$$\text{holds}(\text{colour}(V, \text{Id}), \langle k \rangle), \text{obj}(\text{Id}, -, \langle k \rangle).$$

$$\text{answer}(\langle k \rangle, \text{rotation}, V) :- \text{holds}(\text{rotation}(V, \text{Id}), \langle k \rangle), \quad (2.11)$$

$$\text{holds}(\text{colour}(\text{white}, \text{Id}), \langle k \rangle), \text{obj}(\text{Id}, -, \langle k \rangle).$$

$$\text{expected}(\langle k \rangle, \text{colour}, \text{white}). \quad (2.12)$$

$$\text{expected}(\langle k \rangle, \text{rotation}, \text{down}). \quad (2.13)$$

Two rules are created for questions where v_k is not *Null* because the question gives away two pieces of information: the colour and the rotation. However, when v_k is *Null*, the following ASP rule and fact are used instead:

$$\text{answer}(\langle k \rangle, p_k, V) :- \text{holds}(p_k(V, \text{Id}), \langle k \rangle), \text{obj}(\text{Id}, -, \langle k \rangle). \quad (2.14)$$

$$\text{expected}(\langle k \rangle, p_k, a_k). \quad (2.15)$$

5. The final part of the ASP program is the set of weak constraints used to find the optimal mapping. We firstly define the helper rule, *mapping*, which is used to collate both the colour and rotation property values for each cluster:

$$\begin{aligned} \text{mapping}(C, Col, Rot) :- & \text{colour_mapping}(C, Col), \\ & \text{rotation_mapping}(C, Rot). \end{aligned} \quad (2.16)$$

The weak constraints are then defined as follows:

$$:\sim \text{answer}(Q, Prop, Val), \text{expected}(Q, Prop, Val).[-1@2, Q, Prop, Val] \quad (2.17)$$

$$\begin{aligned} :\sim & \text{mapping}(C1, Col, Rot), \text{mapping}(C2, Col, Rot), \\ & C1 \neq C2.[1@1, C1, C2, Col, Rot] \end{aligned} \quad (2.18)$$

The body of Rule 2.17 is satisfied when question Q is answered correctly. Since we are looking to maximise the number of questions answered correctly and ASP always minimises weak constraints, we give this rule a negative weight. This rule is given the higher priority of the two.

Rule 2.18, on the other hand, says that we prefer answer sets where mappings are unique. This rule could also be used as a hard constraint to rule out any answer sets where the mappings are not unique. However, we choose not to enforce this constraint so that the ASP optimiser could choose an answer set where more questions are answered correctly, at the expense of non-unique mappings. The reason for this decision is the following:

If a group of objects with the same property values is split between two or more clusters, and these objects are commonly asked about in the QA pairs, then the ASP optimiser has the ability to assign multiple object clusters to the same property values, if it leads to a larger number of questions being answered correctly.

We therefore always try to ensure whatever rules the model learns it learns them in order to maximise the number of questions answered correctly. Hence, Rule 2.17 is given the highest priority.

The ASP program is run for each object type t_i , and $\langle prop \rangle_m \text{mapping}$ predicates contained in the optimal answer set in each case are then used to find the optimal mapping $f_{t_i}^*$, which is then stored and used to label all objects in the training data.

2.1.3 Data Labelling

Once $f_{t_i}^*$ has been found for every object type t_i , all of the objects in the training data can be labelled with property values. Labelling an object's property values requires that the centre point of each cluster be computed. For the sake of efficiency, the cluster centres are precomputed for each object type. The centre for a cluster c_j is computed as the average of the object latent vectors assigned to c_j .

After clusters centres have been computed, we assign property values to an object obj_k with type t_i as follows:

1. The object’s image is encoded into a latent vector, v_k , using the autoencoder.
2. For each cluster c_j which corresponds to object type t_i , we compute the cosine distance between v and the centre of c_j . The object is assigned to the cluster with the smallest distance, which we denote c^* .
3. Using the mapping found by the ASP optimisation for object type t_i , we can simply look up the optimal set of property values that correspond to c^* .

Using this method, property values can be assigned to all detected objects in the training videos using only each object’s image. After property values have been assigned, we train the property component in exactly the same way as Chapter ??.

Although the method for training the property component outlined in this section works well when objects are simple and uniform, it is unlikely to scale to more complex, *real-world* datasets. This is because, when training an autoencoder in an unsupervised way, it can learn to extract a lot of noise from the images. When the latent vectors of object images are noisy, clustering will not be as successful in splitting the objects into groups based on their property values.

The speed of the optimisation may also be a concern if the number of property values is large, since the size of the search space scales exponentially with the average number of property values. However, for the purposes of the OceanQA dataset neither of these potential drawbacks causes any problems. The details of the evaluation of the trained property component, along with the entire trained model, is available in Chapter ??.

2.2 Relations

Unlike the relations component in the hardcoded model, which used a manually engineered algorithm, the relations component in the trained model must learn definitions of binary relations between objects from data. As with the trained properties component outlined above, the relations component uses the QA-data version of the OceanQA dataset. Before the relations component is trained, the trained properties component is used to label all of the objects in the training data with property values.

Instead of using manually engineered functions, we opt to use a neural network consisting of fully-connected layers as the core implementation of the relations component. For a generic environment, one neural network would be created for each binary relation, however, since OceanQA only has one relation, only one network is required for this component. Each of these networks is assigned a binary relation, r , to learn. Each network therefore takes a pair of objects as input, and learns to classify the objects as either being related by r or not.

1	0	0	0	0	1	0	0	0	0	0	0	0	1	0	0.21	0.73	0.27	0.79
Type = octopus				Colour = blue						Rotation = downward				Position				

Figure 2.4: The vector encoding for a blue, downward-facing octopus, roughly in position (54, 187, 70, 203). In the component itself the position is given to a much higher degree of accuracy than shown here.

Since fully-connected layers work with numbers, rather than symbolic data, we need to find a way to encode a pair of objects into a vector. To do this we encode each object in the pair separately (an example encoding of an object is shown in Figure 2.4) and then concatenate the two encodings together. The encoding of an object is the result of concatenating each of the following:

1. A one-hot encoding of the object's type.
2. A one-hot encoding of object's colour.
3. A one-hot encoding of object's rotation.
4. The position tuple for the object, where each coordinate has been divided by 256 to produce a number of between 0 and 1.

The neural network takes the vector encoding of the pair of objects as input and passes this vector through a series of fully-connected layers. The details of the layers are shown in Table 2.2. As shown in the table, the output of the network is a single neuron. The sigmoid function is applied to the output of the network to ensure the value is between 0 and 1. If the value of the output of a network learning relation r is denoted o_r , then the relation classifier component's output for a pair of objects, denoted $r(obj1, obj2)$, is as follows:

$$r(obj1, obj2) = \begin{cases} true & \text{if } o_r \geq 0.5 \\ false & \text{otherwise} \end{cases} \quad (2.19)$$

Since the OceanQA dataset has only a single binary relation, *close*, only one neural network is required to be trained for the relation component. Training this network is fairly straightforward; the process is comprised of the following steps:

Layer	Input Size	Output Size
Fully-connected 1	38	1024
Fully-connected 2	1024	256
Fully-connected 3	256	64
Fully-connected 4	64	1

Table 2.2: Details of the layers of each relation classification network. The input is the size of two encoded objects, and the output is used for binary classification. Dropout of 0.2 is also applied between each pair of layers.

1. All of the relation questions in the QA-data version of the training dataset are collated.
2. From each QA pair, a frame number, two objects (including types and, optionally, property values) and the answer are extracted using the question-and-answer parsing component.
3. The frame number and object information is used to look up any missing property values for each object. The property component has already been applied to the dataset, so the property values of all detected objects are guaranteed to be labelled. The full set of property values, as well as the object type and position from the detector, is used to construct the object encoding. The two encodings are concatenated to produce the object-pair encoding.
4. Finally, the set of object-pair encodings, and their associated answers, is used to train the network corresponding to the *close* relation.

After the network has been trained, the component can be used to classifier binary relations between objects. When being applied to a given video, the relations component classifies every pair of objects in every frame of the video as being either close or not using Equation 2.22. If two objects in frame $\langle frame \rangle$ with identifiers $\langle id1 \rangle$ and $\langle id2 \rangle$, respectively, are deemed to be close by the component, then the $obs(close(\langle id1 \rangle, \langle id2 \rangle), \langle frame \rangle)$ predicate is stored. For any pair of objects deemed to not be close, nothing is stored; the closed-world assumption (CWA) - that is, that any predicate which cannot be proved to be true is false - is made for relations between objects.

Applying the network to every pair of objects in every frame of a video could become a burden on the speed of the relations component, especially when there are many objects in each frame. Batching pairs of objects together when applying them to the network helps to improve efficiency, but it has its limits. Chapter ?? shows that the relations component is the second slowest, after the detector, of all the components in the trained model. However, for the purposes of this project, the relations component is fast enough. Chapter ?? also shows the full details of this components performance with perfect inputs.

2.3 Events

The final of the three core components in the trained model is the event component. As with the properties and relations component, the event component uses the QA-data version of the dataset for training. However, unlike the previous components, both the properties and relations components are applied to the training data in order to extract relevant symbolic information, prior to training the event component. During training, the event component can therefore access data that has been extracted by all preceeding components in the architecture: the object detector, properties, tracking and relations components.

Since a lot of symbolic information will already have been extracted from the video during evaluation by preceeding components, the event component does not need to look at raw pixels in the video at all. Instead, it is much more efficient to use the data that has already been gleaned from the video to work out which events occur. Unlike the hardcoded events component, the trained component cannot assume to have access to an \mathcal{AL} model of the environment. Instead, it must learn its own symbolic rules using the training data. One option is to learn an \mathcal{AL} model and use this to find which events occur in the same way as the hardcoded component. However, these \mathcal{AL} rules can be quite complex, so we do not consider this option for this project. Instead, we opt to learn ASP rules which, given a symbolic encoding for a pair of frames, return the action which occurs between the two frames. Once the actions for the entire video have been found, additional ASP rules can be used to find the effects.

2.3.1 Terminology

The ASP rules for finding the actions which occur in a video are found using a greedy search algorithm. The algorithm is run three times; once for each each action. Before outlining this algorithm, we define the terminology used as follows:

- **Feature** :- High-level names, whose values categorise the state of an object. We make use of the following features: x -position, y -position, colour, rotation and disappear (defined below). Each object has a specific value for each feature, which come (directly or indirectly) from the information extracted from the video.
- **Feature Sort** :- Each feature can be categorised into the following three *sorts*, based on the set of possible values that an object can hold for that feature:
 1. **Binary**. Disappear is the only binary feature.
 2. **Discrete** (excluding binary). Colour and rotation are discrete features.
 3. **Continuous**. x -position and y -position are both continuous features.
- **Feature Operation** :- An operation that is applied to the value of a feature. Each feature operation is usually intended to be applied to a single feature *sort*. For example, a discrete feature operation cannot be applied to continuous feature values.
- **Operation Result** :- The result of applying a feature operation to a feature value. Every operation result must be a binary value, since we want to use these values in the body of ASP rules.

The goal of the training algorithm is, given a dataset $\{(F_i, F_{i+1}, a_i)\}$, where F_i and F_{i+1} is the set of feature values for frame i and $i + 1$, respectively, and a_i is the action which occurs between frame i and $i + 1$, find a set of rules with feature

operations in the body which best explain the data. In this case, explaining the data means maximising the number of questions answered correctly, where each QA pair contains a frame i (using which we can look up F_i and F_{i+1} in the extracted data) and an answer (which provides a_i).

As mentioned above, each feature operation is associated with a feature sort. For each sort, we now define the set of possible operations that can be applied to feature values of that sort. For every object in the video, a feature operation σ is applied to the pair $(f_{v_i}, f_{v_{i+1}})$, where f_{v_i} and $f_{v_{i+1}}$ are the object's values for feature f in frame i and $i + 1$, respectively. Every set of feature operations also includes the \top operation, which always evaluates to true.

Firstly, for binary features we only consider values in the initial frame, however, it would be very simple to expand this to include both frames. In the following, we show the the set of possible feature operations for a binary feature, Σ_{binary} . Each operation takes the pair $(f_{v_i}, f_{v_{i+1}})$ as input:

$$\Sigma_{binary} = \left\{ \begin{array}{l} \sigma(f_{v_i}, f_{v_{i+1}}) := f_{v_i} \\ \sigma(f_{v_i}, f_{v_{i+1}}) := \neg f_{v_i} \\ \sigma(f_{v_i}, f_{v_{i+1}}) := \top \end{array} \right\} \quad (2.20)$$

Secondly, a pair of discrete feature values $(f_{v_i}, f_{v_{i+1}})$ can have the following feature operations applied:

$$\Sigma_{discrete} = \left\{ \begin{array}{l} \sigma(f_{v_i}, f_{v_{i+1}}) := f_{v_i} = v_n, f_{v_{i+1}} = v_m \text{ for each } (v_n, v_m) \\ \sigma(f_{v_i}, f_{v_{i+1}}) := f_{v_i} = f_{v_{i+1}} \\ \sigma(f_{v_i}, f_{v_{i+1}}) := f_{v_i} \neq f_{v_{i+1}} \\ \sigma(f_{v_i}, f_{v_{i+1}}) := \top \end{array} \right\} \quad (2.21)$$

Where (v_n, v_m) loops over all possible combinations of feature values.

Finally, a pair of continuous feature values $(f_{v_i}, f_{v_{i+1}})$ can take the following feature operations:

$$\Sigma_{continuous} = \left\{ \begin{array}{l} \sigma(f_{v_i}, f_{v_{i+1}}) := f_{v_i} > f_{v_{i+1}} \\ \sigma(f_{v_i}, f_{v_{i+1}}) := f_{v_i} < f_{v_{i+1}} \\ \sigma(f_{v_i}, f_{v_{i+1}}) := f_{v_i} = f_{v_{i+1}} \\ \sigma(f_{v_i}, f_{v_{i+1}}) := f_{v_i} \neq f_{v_{i+1}} \\ \sigma(f_{v_i}, f_{v_{i+1}}) := \top \end{array} \right\} \quad (2.22)$$

The \top feature operations always gives an operation result of *true*. Intuitively, a \top feature operation has no effect on the outcome of the rule it is in. Since the algorithm has the option to choose the \top feature operation for any feature, we can enforce that all features appear in every rule body at least once. Any feature that the algorithm doesn't care about can be assigned a \top feature operation. For simplicity, we also enforce that each feature can appear at most once in a rule body. The search algorithm therefore tries to learn a set of rules, where each rule takes the following form:

$$a_i \text{ if } \sigma_{x_position}(f_{v_i}, f_{v_{i+1}}), \sigma_{y_position}(f_{v_i}, f_{v_{i+1}}), \sigma_{colour}(f_{v_i}, f_{v_{i+1}}), \quad (2.23) \\ \sigma_{rotation}(f_{v_i}, f_{v_{i+1}}), \sigma_{disappear}(f_{v_i}, f_{v_{i+1}})$$

Where each $\sigma_{\langle f \rangle}$ is a member of $\langle f \rangle$'s set of possible feature operations.

Each feature operation is also assigned a *weight*. The weight of a rule is then the sum of the weights of its operations. Secondary goals for the optimisation algorithm are, firstly, to minimise the number of rules and, secondly, to minimise the sum of weights of the rules.

Before outlining an Inductive Logic Programming (ILP) search algorithm for finding a set of rules, we first give an example which explains the terminology.

Example. Consider an element of the dataset which includes an action, *move*, and an octopus which is upward-facing, blue and has position (100, 40, 116, 56) in frame 6, and is upward-facing, purple and has position (100, 25, 116, 41) in frame 7. For simplicity we only consider the *x*-position and *y*-position features. Each position feature only uses the top left corner of the object's bounding box, since both corners move by the same amount. The position feature operations are given as follows:

<i>x</i> -position operation \rightarrow result	<i>y</i> -position operation \rightarrow result
$100 > 100 \rightarrow false$	$40 > 25 \rightarrow true$
$100 < 100 \rightarrow false$	$40 < 25 \rightarrow false$
$100 = 100 \rightarrow true$	$40 = 25 \rightarrow false$
$100 \neq 100 \rightarrow false$	$40 \neq 25 \rightarrow true$
$\top \rightarrow true$	$\top \rightarrow true$

Clearly, given only this single example, there are many possible rules which the algorithm could learn, many of which will be undesirable. However, when there are a more significant number of examples, the algorithm should be able to learn rules which correctly model the environment.

2.3.2 Search Strategy

Equation 2.23 shows that, for each rule, the training algorithm needs to search over the set of all possible combinations of feature operations, where each feature is assigned exactly one operation in the rule body. Even for a relatively simple environment, such as OceanQA, this is a very large search space, and when the possibility of searching over an arbitrary number of rules is taken into account the space becomes even larger. The number of feature operations for each feature and the total search space size is shown in Table 2.3.

To ensure the training algorithm remains tractable, we implement a greedy search over features, rather than a complete search over the entire space. Our search algorithm takes each feature individually and finds the operation for the given feature which maximises the number of action questions answered correctly. After each optimisation, we store the set of operations which have been selected as the *accumulated operations*. When the next feature is optimised, each rule in the new hypothesis can choose to use either the full set of stored operations or none - using none means setting all previously optimised feature operations to \top .

Feature	Search Size
x -position	5
y -position	5
colour	52
rotation	19
disappear	3
Total	74100

Table 2.3: Number of feature operations for each feature. The total search size is then calculated by multiplying all these sizes together.

Algorithm 2 gives a high-level overview of the steps required to find the optimal hypothesis - the set of rules which correctly answers the most questions - for every action. Algorithm 2 makes use of the `OPTIMISEFEATURE` function, which takes as input: an action, a set of feature operations, the accumulated feature operations and the training data. This function searches over all of the feature operations to find the set of rules which maximise the number of questions answered correctly. It also searches over the number of rules to be used, as well as whether each of those rules should use an accumulated operation set or not. We implement this function using an ASP optimisation program, the encoding for which is outlined in Section 2.3.3.

In the remainder of this section we provide an inductive description for the task which the `OPTIMISEFEATURE` function must solve.

Consider an optimisation which uses features f_1, \dots, f_{feat} , and that the algorithm is currently optimising feature f_i . The algorithm will therefore be given: an action a ; the set of feature operations Σ_s , where s is the sort of f_i ; the *accumulated operations* acc_1, \dots, acc_m from the previous feature optimisation, where each acc_j is the set of feature operations used by rule r_j in the previous optimisation; and the training data D . In the case where $i = 1$, the *accumulated operations* contains a

Algorithm 2 Finding the optimal hypothesis for each action

```

1: procedure OPTIMISERULES(actions, features, data)
2:   hyps  $\leftarrow \{\}$ 
3:   for  $a \leftarrow \text{actions}$  do
4:     acc  $\leftarrow \{\}$ 
5:     for  $f \leftarrow \text{features}$  do
6:        $s \leftarrow \text{SORTOF}(f)$ 
7:       acc  $\leftarrow \text{OPTIMISEFEATURE}(\text{action}, \Sigma_s, \text{acc}, \text{data})$ 
8:     end for
9:     hyp  $\leftarrow \text{GENHYP}(a, \text{acc})$ 
10:    hyps  $\leftarrow \text{hyps} + (a, \text{hyp})$ 
11:  end for
12:  return hyps
13: end procedure

```

single rule whose feature operations are all set to \top . The search space for feature optimisation f_i can then be constructed as follows:

1. At the highest level of the search, the algorithm needs to choose how many rules to learn. Therefore we start by creating one node in the search tree for each possible number of rules in $\{1, 2, \dots, R_{max}\}$, where R_{max} is a hyperparameter of the search.
2. Then, for each possible number of rules, each rule can choose to use at most one of $\{acc_1, \dots, acc_m\}$. Therefore, for each previous node n , which allows R rules, each rule in $\{r_1, \dots, r_R\}$ chooses between m accumulators or none, meaning $(m + 1)^R$ child nodes must be added to n , since $(m + 1)^R$ is the number of possible combinations of accumulator choices.
3. Finally, after choosing which accumulator to use (or none), each rule chooses a feature operation for the feature f_i . This means, for each leaf node n constructed so far, we add a child node to n for each possible feature operation $\sigma_i \in \Sigma_s$.

Once the search space has been constructed, the optimiser evaluates each of the leaf nodes in the tree outlined above, and chooses the node which maximises the number of action questions answered correctly. The `OPTIMISEFEATURE` function then returns accumulated set of feature operations, $\{acc_1, \dots, acc_i\}$, so that they can be used for the next feature optimisation.

2.3.3 ASP Encoding

This section outlines the details the ASP program which is used in place of the `OPTIMISEFEATURE` function described above. As expected, the program creates the search space using choice rules, and evaluates each possible answer set by counting the number of action questions it answers correctly.

Bibliography

- [1] 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.
- [2] Jiayuan Mao et al. “The neuro-symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision”. In: *arXiv preprint arXiv:1904.12584* (2019).
- [3] Kexin Yi et al. “Clevrer: Collision events for video representation and reasoning”. In: *arXiv preprint arXiv:1910.01442* (2019).
- [4] 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.