

Learning to Play General Video-Games via an Object Embedding Network

William Woof & Ke Chen

School of Computer Science, The University of Manchester

Email: {william.woof, ke.chen}@manchester.ac.uk

Abstract—Deep reinforcement learning (DRL) has proven to be an effective tool for creating general video-game AI. However most current DRL video-game agents learn end-to-end from the video-output of the game, which is superfluous for many applications and creates a number of additional problems. More importantly, directly working on pixel-based raw video data is substantially distinct from what a human player does. In this paper, we present a novel method which enables DRL agents to learn directly from object information. This is obtained via use of an object embedding network (OEN) that compresses a set of object feature vectors of different lengths into a single fixed-length unified feature vector representing the current game-state and fulfills the DRL simultaneously. We evaluate our OEN-based DRL agent by comparing to several state-of-the-art approaches on a selection of games from the GVG-AI Competition. Experimental results suggest that our object-based DRL agent yields performance comparable to that of those approaches used in our comparative study.

I. INTRODUCTION

General video-game AI (GVG-AI) is an area regarding the development of general algorithms that enable AI agents to play a wide range of different video-games with minimal tailoring to specific games. While developing techniques for General AI is a key focus of research into GVG-AI, general video-game playing agents also have a number of applications within the games industry. Besides from the obvious applications, such as a replacement to hand-coded in-game AI, GVG-AI can also either be used as a development tool or as a proxy for human play-testers. Such agents could be employed effectively in a wide array of applications from testing game balance [1] to evaluating procedurally generated content [2]. As well as applications within games and games design, GVG-AI also has wider implications for the field of AI, as techniques which work well on video-games can often also be applied to real-world problems.

One promising area in the search for general video-game players is Deep Reinforcement Learning (DRL). DRL agents have been successfully applied to a wide range of video-games ranging from 2D arcade games [3] to challenging 3D shooters [4]. These agents learn through interacting with the game autonomously, using a deep neural network to select actions based on the current state of the game. During this process the agent receives rewards (usually dictated by the in-game score) which indicate how well it is performing. By using an appropriate reinforcement learning algorithm the agent is able to modify its neural network in order to maximise this

reward signal. However, such agents are far from perfect, and can sometimes be difficult to apply in practice – a problem compounded by the fact that they typically take a long time to train. A significant consideration into the design of these agents is how information from the game is presented to their neural networks (*i.e.* the representation given to the agent), as well as the design of the networks themselves.

Current state-of-the-art DRL approaches to video-games learn directly from raw video data, using deep *convolutional neural networks* (CNNs) [3]. While widely applicable, this approach is subject to limitations for certain applications. For example, many games (e.g., *Starcraft*) feature controllable cameras, meaning much of the game-state is obscured from the agent at any given point. Working around this, e.g., by putting the camera under the agent’s control, adds additional complexity to the agent. Additionally, in many cases it may be undesirable, or even impossible to produce a video-output for the agent to consume. Rendering videos for multiple agents may be prohibitively expensive, and in some cases there may be no obvious way to produce a good visual representation for NPCs (non-player characters). More importantly, interpreting raw video data at a pixel level is substantially different from how human players appear to play, as studied in [5].

Unlike many other reinforcement learning tasks, the ground-truth information about the current state of the environment is often available in video-games, although such information needs to be organised and presented to an agent in some way. Hence, the use of this direct information about the current game-state could be alternative to working directly on raw video data. For instance, Samothrakis et al [6] employ a fixed set of general features, *i.e.*, distance to the nearest enemy, number of tokens collected and so on, to encapsulate the current game-state. While this approach often works well, those general features have to be handcrafted, which is laborious and requires human expertise. Moreover, this approach is relatively game-specific and hence generally inappropriate to GVG-AI. To overcome this limitation, various game-independent object representations, have been employed [7], [8]. In an object representation, each game-state observation is given as a list of objects, and their classes and attributes. For example, the state of one round of the game *Pong* might be represented by two objects of the class `bat`, with attributes of `x-coord`, `y-coord`, and `player`, and an object of the class `ball` with attributes `x-coord`, `y-coord`, `x-velocity`, and `y-velocity`. Many video-games rely on objects for their internal representation of the game-state. For instance, the

popular *Unity 3D* game engine relies heavily on game objects, and even the early *Atari* console used a primitive sprite-based system. In general, the use of object representations not only leads to an effective approach to representing game-states across a wide variety of video games but also has a number of practical benefits. For example, it allows the use of different subsets of objects for different agents. Also, objects provide useful anchor points for applying various advanced reinforcement learning techniques such as hierarchical reinforcement learning [9], intrinsic motivation [10], and planning [11]. Given there are a different number of objects in different states of game, however, how to structure this information in a way that can be input into a conventional deep neural network is a key issue of using object representations with DRL. Previous solutions to this problem mimic an image representation by using an ‘‘object perception grid’’, where objects are overlaid onto a grid and mapped to the nearest cell determined by their x and y co-ordinates within the game. The number of objects mapped to each cell is then used as an input for a conventional neural network (either fully-connected or convolutional). Unfortunately, such a solution requires selecting an appropriate grid size manually, and entails a large input space, increasing the required neural network complexity. In general, it also still suffers from many of the same problems as encountered by using raw image representations, such as a restricted field of view.

In this paper, we present a novel approach to address the object representation issues in GVG-AI. To overcome all the aforementioned limitations, we adopt a specific type of neural network architecture, set networks [12]–[14], to develop our *object embedding network* (OEN). This network can not only take a list of object-feature vectors of arbitrary lengths as input to produce just a single, yet unified, fixed-length representation of all the objects within the current game-state, but also be trained on a given task simultaneously. Hence, our OEN-based approach provides an alternative way to apply DRL algorithms within video-games, based on object information. Our approach is generally motivated by recent advances within approaches to relational reasoning [13] and dynamics prediction [14], which suggest that working with objects, rather than raw data, can help scale up deep learning to more complicated tasks in a similar fashion to human information processing.

Our main contributions in this paper are summarised as follows:

- 1) We propose an OEN model, based on set networks, for learning directly from sets of object feature vectors.
- 2) We develop an OEN-based GVG-AI agent for playing general video games.
- 3) We evaluate our approach on selected games from the GVG-AI competition and demonstrate that it performs comparably to a variety of other popular approaches for representing game states.

The rest of this paper is organised as follows. Section II reviews related work. Section III presents our OEN model. Section IV describes our object-based approach to GVG-AI. Section V describes our experimental settings and reports experimental results. Finally, Section VI discusses issues and

implications arising from this study.

II. RELATED WORK

A. Deep Reinforcement Learning for video-games

Reinforcement learning is a sub-field of machine learning where agents autonomously interact with an environment and seek to maximise some reward signal. *Deep reinforcement learning* (DRL) is an extension of classical ‘tabular’ approaches to reinforcement learning which enable these techniques to be applied to more complicated problems. By using the in-game score counter as a reward signal, DRL can be applied to develop agents for playing video-games [3], making video-games a popular test-bed for new DRL algorithms. Those algorithms work by using a deep neural network to ‘score’ possible actions given a particular game-state, which is then trained according to a loss function. Such a loss function for DRL may be formulated based on the agent’s past experience as well as the reward obtained.

The *deep-Q network* (DQN) algorithm [3] is a pioneering work in applying DRL to general video games playing, where the DQN was trained to play a variety of games for the *Atari 2600* games console. In the original DQN algorithm, the game-state is presented to the agent as a series of four images from four consecutive frames of video output, which is interpreted by a deep *convolutional neural network* (CNN) with four input channels. The success of this DQN algorithm has led to a number of alternative DRL algorithms for video-game playing [15], [16] for the *Atari* system, and these algorithms have also been applied to a variety of other video games [4], [17]. While the DQN algorithm and its variants can be easily adapted to a variety of input types and network architectures, these mainly use the same input format; i.e. a sequence of frames from a game video stream. A notable exception is the use of object perception grids, which is described below in Section II-B.

B. Object-oriented reinforcement learning

Reinforcement learning from objects has previously been studied within object-oriented reinforcement learning [18], which allows for exploiting structural information at an object level. In object-oriented reinforcement learning, the state-space is expressed in terms of a set of objects. These objects all belong to some class from a fixed set of classes. Each object is an ordered tuple of object attributes, where the domain of these attributes is determined by the object class. For example, for a simple empty 5×5 grid-world, we might formulate the state-space with a single class, *Agent*, of attributes x and y , $Dom(x) = Dom(y) = \{1, 2, \dots, 5\}$. Then the initial state s_0 would be a single object agent from the *Agent* class with attribute values $x = y = 0$.

Conventional approaches to solving these problems usually involve planning algorithms, and first order logic [19], or otherwise rely on the discrete nature of the environment, which is generally incompatible with DRL approaches. In particular, structuring this information in a way that can be used by a neural network is a challenging problem. Very recently, this problem has been addressed via an object perception grid

representation, e.g., [7], [8]. In this representation, objects are mapped onto grid squares, based on their x and y attributes. Each of these grid squares is then treated as an input neuron, which is set to 1 if an object of a given class is present at that cell, and 0 otherwise. The full representation is then given by multiple input grids, one for each possible object class. This effectively produces an image-like representation, which can be fed into a CNN. Apart from removing some of the complexity of the input, this approach suffers from most of the same problems as a visual representation. Moreover, this representation also requires the designer to select a grid size and coarsity. While many games may feature a natural choice for these, for many games it may require careful selection and additional tuning to select these parameters appropriately. Additionally this process can only be applied where objects are related by a clear 2D structure.

III. OBJECT EMBEDDING NETWORK

In this section, we present an *object embedding network* (OEN) to learn a unified object representation from arbitrary sets of objects characterised by a variety of object features without being limited to 2D spatial structures. Our OEN model is based on an emerging class of deep neural-network architecture, set networks [12], which were recently developed to tackle the input data in a set form. In general, objects in a game-state naturally stand in a set form. By using the same principles behind set networks, our OEN transforms an arbitrary number of object feature vectors corresponding to a game state into a single fixed-length “unified” object representation. This unified representation can be used for a variety of different purposes. Our OEN can be trained to learn a unified representation and fulfil a specific learning objective simultaneously.

At game-state s_t , assume that there is a set of objects, $O^{(t)} = \{o_k^{(t)}\}_{k \in 1, \dots, K_t}$, where each object, $o_k^{(t)}$, can be characterised by a feature vector (or a number of attributes), $\mathbf{x}_k^{(t)}$. Hence, the feature vectors of all K_t objects collectively form a set, $X^{(t)} = \{\mathbf{x}_k^{(t)}\}_{k \in 1, \dots, K_t}$, for game-state s_t . Our problem is how to learn a fixed-length unified feature vector that retain as much representative information conveyed by K_t objects as possible for arbitrary K_t .

A common way to get representative information of a set of vectors is to compute some statistic about the set. In practice, this can be achieved using simple arithmetic pooling functions, e.g., max or sum pooling, applied element-wise, which condense an input set of vectors into a single fixed-length vector of the same dimension. However, simply applying simple pooling functions over a set of object feature vectors is likely to incur a loss of important information. For example, if the object features consist of x and y co-ordinates then taking the mean of all object feature vectors simply ends up with the average position of all objects. While this is useful information, it does not convey the important information, e.g., “is object $o_i^{(t)}$ next to object $o_j^{(t)}$?”. Hereinafter, we drop out the explicit game-state index, t , to facilitate our presentation.

Motivated by set networks [12], we deal with this problem by embedding raw feature vectors of objects into a higher

dimensional space, which allows for retaining non-trivial information after pooling. This can be achieved by applying a proper “embedding” function E to the feature vector of each object:

$$E(X) := \{E(\mathbf{x}_1), \dots, E(\mathbf{x}_K)\}.$$

Let Π to denote a pooling function, a unified representation of the object feature set, X , is then achieved by

$$\mathbf{r}(X) = \Pi_{k \in 1, \dots, K} E(\mathbf{x}_k).$$

However, finding a proper embedding function explicitly is extremely difficult in general. Also, the optimal choice of such an embedding function is task-dependent. Instead of using an explicit embedding function, we can employ a neural network to learn an optimal embedding function. Furthermore, for a specific task based on the unified representation, we can incorporate another neural network for fulfilling the given learning objective and learning the optimal embedding function simultaneously.

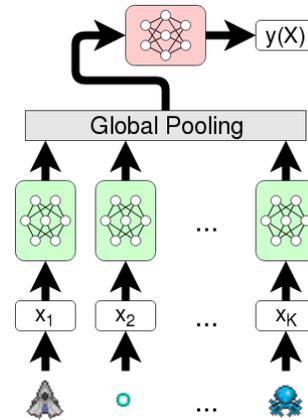


Fig. 1: Object embedding network architecture.

As illustrated in Fig. 1, a generic object embedding network (OEN) consists of *embedding network* (shown in green), *global pooling function* and *task network* (shown in pink). For a set of K objects in a game-state, K identical embedding networks $E(X; \theta_E)$ are employed for embedding different objects, respectively, where θ_E is a collective notation of parameters shared by all K embedding networks. The global pooling function condenses the embedding object representations produced by K embedding networks to yield a fixed-length unified object representation: $\mathbf{r}(X) = \Pi_{k \in 1, \dots, K} E(\mathbf{x}_k, \theta_E)$. Then, the unified representation is fed to the task network denoted by $P(\mathbf{r}(X); \theta_P)$ where θ_P is a collective notation of parameters in the task network.

Parameter estimation in the OEN is done by optimising a loss function defined on training data, \mathcal{D} , given for a specific task, $L(\mathcal{D}; \theta_E, \theta_P)$:

$$\{\theta_E^*, \theta_P^*\} = \operatorname{argmin}_{\theta_E, \theta_P} L(\mathcal{D}; \theta_E, \theta_P).$$

For some loss function L , e.g., in a supervised learning task, a prediction-error based loss function can be used. In

Section IV, we detail a loss function defined on transitions (s_t, a_t, r_t, s_{t+1}) , drawn from experience-replay, for our reinforcement learning tasks.

Contextual information regarding the relationship between an object and other co-occurring ones at the same game-state can play an important role. We can exploit this by replacing our embedding function $E(\mathbf{x})$ with a “contextual” embedding function $E(\mathbf{x}, X)$ which takes into account information from the wider set when embedding each object. To this end, different techniques have been proposed in set networks, e.g., [12], [14], [20], to explore this contextual information. Motivated by the work of [12], [20], we adopt a simple global-context based method to explore the contextual information in our work. In this method, some statistic Π of feature vectors of all the objects in the set X is first estimated by $\hat{\mathbf{x}} = \Pi_{\mathbf{x} \in X} \mathbf{x}$. Then, the feature vector of each object, \mathbf{x} , is concatenated with this statistic vector, $\hat{\mathbf{x}}$, to form a “contextualised” feature vector of the object: $(\mathbf{x}, \hat{\mathbf{x}})$. Instead of the feature vector of each object, \mathbf{x} , its contextualised feature vector, $(\mathbf{x}, \hat{\mathbf{x}})$, is fed to the embedding network in the OEN. In particular we adopt the same “equivariant” transformation proposed in [12] which is given by:

$$f_{equiv}(\mathbf{x}, X) = \mathbf{x} - \text{maxpool}(X).$$

By using this method, our OEN model can be extended to explore contextual information without altering its general architecture and learning algorithms.

IV. MODEL DESCRIPTION

In this section, we present our method for establishing an OEN-based DRL agent for playing general video games. We first describe object and feature extraction required by the OEN and then propose our OEN implementation and its deep Q-learning algorithm.

A. Object and feature extraction

For our agent, the process of identification and extraction of objects is handled by the environment. That is, we assume that object extraction can be done directly via access to the ground truth of the environment. Hence, the wide-spread use of object-oriented programming languages should help with this process, as many objects are likely to be treated as such in code. Thus, the description of a game state is given in an object-oriented format; i.e., observation O from a game-state is given as a list of objects, $O = \{o_1, \dots, o_K\}$.

Given a list of objects in this format, we still need to characterise those objects via a number of attributes to meet our requirement of our OEN-based DRL. In essence, this is a feature extraction process to obtain object feature vectors from the raw objects based on their attributes, which leads to a pre-processing function `process_observation` required by our OEN. A natural solution to feature extraction is concatenating all real-valued attributes of an object into a single feature vector. However, this solution results in a problem; while our OEN can handle only fixed-length object feature vectors, the number of attributes used to characterise an object

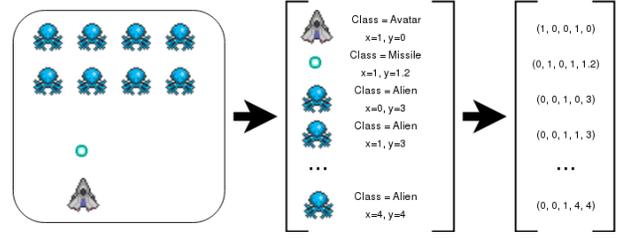


Fig. 2: Exemplar object feature extraction process.

is not fixed and different objects could have a different number and types of attributes. Hence, the user must select a set of attributes applicable to all the objects for a fixed-length feature vector. In this manner, a set of fixed-length feature vectors, $X = \{\mathbf{x}_1, \dots, \mathbf{x}_K\}$, can be extracted for a list of objects, O , in a game-state. This set of feature vectors are fed into our OEN, which acts as a value-network for the agent.

Fig. 2 depicts an exemplar object feature extraction process. As seen in Fig. 2, a game-state is broken down into a (finite) set of objects, along with a number of attributes for each object, as chosen by the user, e.g., position, class of object, and so on. This list of objects is then converted into a list of feature vectors by mapping each object to a fixed-length vector based on its attributes. In this example, a one-hot vector of the object class concatenated with the object’s co-ordinates yields a 5-dimensional object feature vector.

Algorithm 1 Q-learning algorithm for OEN

```

// Initialise agent
initialise empty replay memory  $M$ 
initialise  $\theta$  for OEN Q-function  $Q(s, a; \theta)$ 
 $\theta_{target} \leftarrow \theta$ 
// Initialise environment
env.reset()
step  $\leftarrow 0$ 
while step  $\leq$  max_step do
  step  $\leftarrow$  step + 1
  // Get action from Q-function
   $s_t \leftarrow$  env.get_state()
   $O_t \leftarrow$  process_observation( $s_t$ )
   $a_t \leftarrow$  epsilon_greedy( $Q(O_t, a; \theta)$ )
  // Step environment and observe result
   $r_t, s_{t'} \leftarrow$  env.apply_action( $a_t$ )
   $O_{t'} \leftarrow$  process_observation( $s_{t'}$ )
   $T_{t'} \leftarrow$  env.has_ended()
  add tuple ( $O_t, a_t, r_t, O_{t'}, T_{t'}$ ) to  $M$ 
  // Train network
  sample tuple ( $O'_{t'}, a'_{t'}, r'_{t'}, O'_{t'}, T'_{t'}$ ) from  $M$ 
   $Q_{targ} \leftarrow r'_{t'} + T'_{t'} \cdot \gamma \cdot \max_a Q(O'_{t'}, a; \theta_{target})$ 
  update  $\theta$  via gradient descent on  $(Q(O'_{t'}, a'_{t'}; \theta) - Q_{targ})^2$ 

if step % 1000 == 0 then
   $\theta_{target} \leftarrow \theta$ 
end if
end while

```

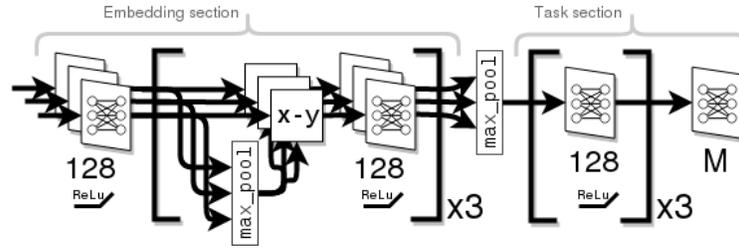


Fig. 3: The object embedding network used to implement our DRL agent.

B. OEN-based Q-learning

Motivated by the deep set network of [12], we develop an object embedding network to implement our DRL agent, as illustrated in Fig. 3. In this OEN, the embedding network consists of four layer of 128 ReLU units, with an “equivariant” transformation $f_{equiv}(\mathbf{x}, X) = \mathbf{x} - \text{maxpool}(X)$ between each layer. To generate a unified object representation, all embedding representations are pooled by element-wise max pooling across the whole set. The task network consists of three fully connected hidden layers of 128 ReLU units, and a final output layer of M linear units corresponding to value functions of M possible actions used in playing the given video games. It is worth mentioning that our OEN-based DRL implementation is largely identical to the DQN-based DRL [3] apart from two aspects: a) we use the OEN shown in Fig. 3, while the DQN uses a deep CNN as a learning model, and b) our OEN works on object feature vectors, while the DQN works on raw video data. Thus, the same deep Q-learning algorithm can be adapted to train our OEN-based DRL agent.

In a Q-learning based DRL algorithms, the Q-values are estimated using a value function given by a neural network $Q(s, a; \theta)$, where s is a game-state, a is a selected action and θ is a collective notation of all the parameters in this neural network. The policy for the agent is then given by selecting the action which maximises $Q(s, a; \theta)$ for the agent’s current state. In the DRL learning algorithm proposed by Mnih et al [3], the DQN actually outputs a value vector, $\mathbf{Q}(s; \theta)$, for all the actions simultaneously, where $Q(s, a_m; \theta)$ is the m th element of this output vector, reflecting the value of the m th action. Given a sequence (s_t, a_t, r_t, s_{t+1}) , we can obtain an estimate for $Q(s_t, a_t)$ using the Bellman equation:

$$\hat{Q}(s_t, a_t) = r_t + \gamma \max_a Q(s_{t+1}, a).$$

This estimate is used as a target Q_{target} for $Q(s_t, a_t)$. Thus, we define the Q-learning loss function at game-state t as follows:

$$L(\mathcal{D}_t; \theta) = (Q(s_t, a_t; \theta) - Q_{target})^2,$$

where $\mathcal{D}_t = (s_t, a_t, r_t, s_{t+1})$ is training data retrieved from an experience-replay memory [3]. Since the new estimate Q_{target} depends heavily on the previous values $Q(s_{t+1}, a; \theta)$, a separate network parameterised with the known θ_{target} is used to obtain Q_{target} estimates. θ_{target} is then updated once the new parameters θ in the OEN are achieved during the learning. To optimise the loss function for the Q-learning, we employ the Adam optimizer [21], a gradient-based optimiser. Network

parameter update is also done in mini-batches of multiple state-actions and targets simultaneously. Since multiple sequences of different lengths are not readily expressible as fixed-size tensors (which is required by most deep-learning libraries), for each batch we pad each sequence with zero vectors until they have the equal length. A mask of these zero elements is produced and used to nullify their contribution to the output of the OEN. For clarity, we describe the detailed Q-learning algorithm used for training our OEN in Algorithm 1.

V. EXPERIMENT

In this section, we evaluate the performance of our OEN-based DRL agent on five selected games used in the GVG-AI competition [22] by comparing with two baseline agents that use different representations of the game-state, and measure average performance during training on a variety of games. To ensure a fair comparison between the different forms of representations, for each agent we change only the agent’s neural network and the format of the observation presented to the agent, keeping the rest of the agent design the same.

A. Test environment

In order to adequately test our approach, we require a corpus of distinct video-games, preferably unified under a single framework. A common choice for this is the *arcade learning environment* (ALE) [23], however, the ALE does not provide access to ground truth information about object attributes, hence we instead look to the GVG-AI Competition.

The GVG-AI competition [22] is a regular competition challenging researchers to build AI agents capable of playing a variety of different video-games. These games are specified via *video-game description language* (VGDL) [24] where games are defined by a block-based sprite system, and are often based on well-known titles. Importantly, games from previous rounds of the competition are released to the public, proving a large collection of games, all running within the same framework. See Figure 4 for some examples of games from the GVG-AI Competition.

Importantly, information about the current game-state is presented to the agent in the form of a state-observation which includes a list of information about the various sprites (i.e., game objects) within the game. Additionally, while not explicitly provided to the agent, a video output is also produced for human consumption.

As well as giving direct access to in-game objects, VGDL provides a number of other benefits as a reinforcement learning test-bed, including:



Fig. 4: Five games from Test Set 1 of the GVG-AI Competition used in our experiments: (a) Aliens. (b) Boulderdash. (c) Missile Command. (d) Survive Zombies. (e) Zelda. Missile Command, Boulderdash and Zelda are based on classic arcade games of the same names, while Aliens is loosely based on the game *Space Invaders*.

- Existing tasks can easily be modified to test how this affects the agent.
- New games/tasks can be quickly synthesised for particular purposes.
- A large available pool of pre-existing tasks in the form of GVG-AI competition games.
- Direct access to the underlying mechanics and ontology, which may be useful as a ground-truth for investigating things such as model-based agents and transfer learning.

For our environment we adapted the original `py-vgdl` code¹, adding in some extra functionality from the competition version such as support for sprite images. Another option would have been to use the code provided for the GVG-AI Competition, although when we started this work there was no python client available.

B. Experimental Settings

To evaluate the performance of each agent we train the agent for 2,000,000 steps, testing the agent for 50 episodes (without training) every 50,000 steps, and record the reward obtained over each of these episodes. While there are a number of different indicators of agent performance, e.g., percentage of games won, we select average episode rewards as this most closely reflects the reinforcement learning objective of the agent, hence is less sensitive to factors such as a mis-specified reward. Additionally, different users may have different criteria for how they want agents to perform during training, i.e., some may be interested in short-term performance after a certain number of training steps, while others may be interested only in the ‘asymptotic’ final performance, hence we record agent performance throughout training.

We modified our environment code to include support for three different forms of observation types:

- 1) Image representation: A sequence of raw pixel images of the game screen, appropriately sized to be close to the 84×84 post-processed resolution of the original DQN algorithm.
- 2) Object representation: A list of game objects, each given in the form of vector with: a one-hot vector of object class, object co-ordinates, object orientation, and values of given object resources.
- 3) Feature representation: A list of the shortest distances from each object class to the player avatar, plus a list of any additional avatar resources. This is the same as the features described in [6].

In order to simplify our input space, and since certain objects in the game are irrelevant (or invisible) to the agent, for each game we defined a list of the important object classes, and ignore any objects from classes not on that list for both the objects and features representation. We also remove the frame-skip functionality from our agent as GVG-AI Competition games have a slower update rate, so being able to select actions at each step is important.

For each of these three observation types, we modify our baseline agent as follows, giving us three different agents²:

- 1) For the image representation we use the same CNN as used in [16] with a final linear layer of M outputs. Similar to the original DQN algorithm we also compile states from two consecutive frames (we do not use four frames as sprites are visible every frame in VGDL which is not the case for certain ALE games³, and this reduces computational burden).
- 2) For the object representation we use the OEN described in Figure 3.
- 3) For the feature representation, we use a simple fully-connected neural network with layers of 64, 64, 128, and M ReLU units, respectively.

Where M is the number of possible actions in the given environment. All the agent hyper-parameters use in our experiments are as follows: $\gamma=0.99$, $\epsilon\text{-start}=0.5$, $\epsilon\text{-final}=0.1$, $\epsilon\text{-anneal-step}=500,000$, $\text{replay-memory-size}=50,000$ $\text{learning-rate}=0.00025$, $\text{mini-batch-size}=32$, the agent is trained every four steps, and the parameters are initialised randomly with a Gaussian distribution: $\mathcal{N}(0, 0.1)$.

We select five games from Test Set 1⁴ of the GVG-AI competition as our test set: Aliens, Boulderdash, Missile Command, Survive Zombies, and Zelda. For each of these we use the first level (i.e., `level_0`) as our game environment. We train our agents on each game four times, using a different seed for agent initialisation each time. Initial environment seeds are reset to the same value for each agent, ensuring that there are no differences in agent performance due to different environment initialisations.

C. Results

Full results of our experiments are shown in Figure 5. We also report the best mean test score on each game for each

²Our agent & experimental code is available at: <https://github.com/EndingCredits/Object-Based-RL>

³Due to the limited sprite buffers of the *Atari* console, a common optimisation is to draw certain sprites only every other frame.

⁴Found here: http://www.gvgai.net/training_set.php?tg=1

¹Available here: https://github.com/EndingCredits/gym_vgdl

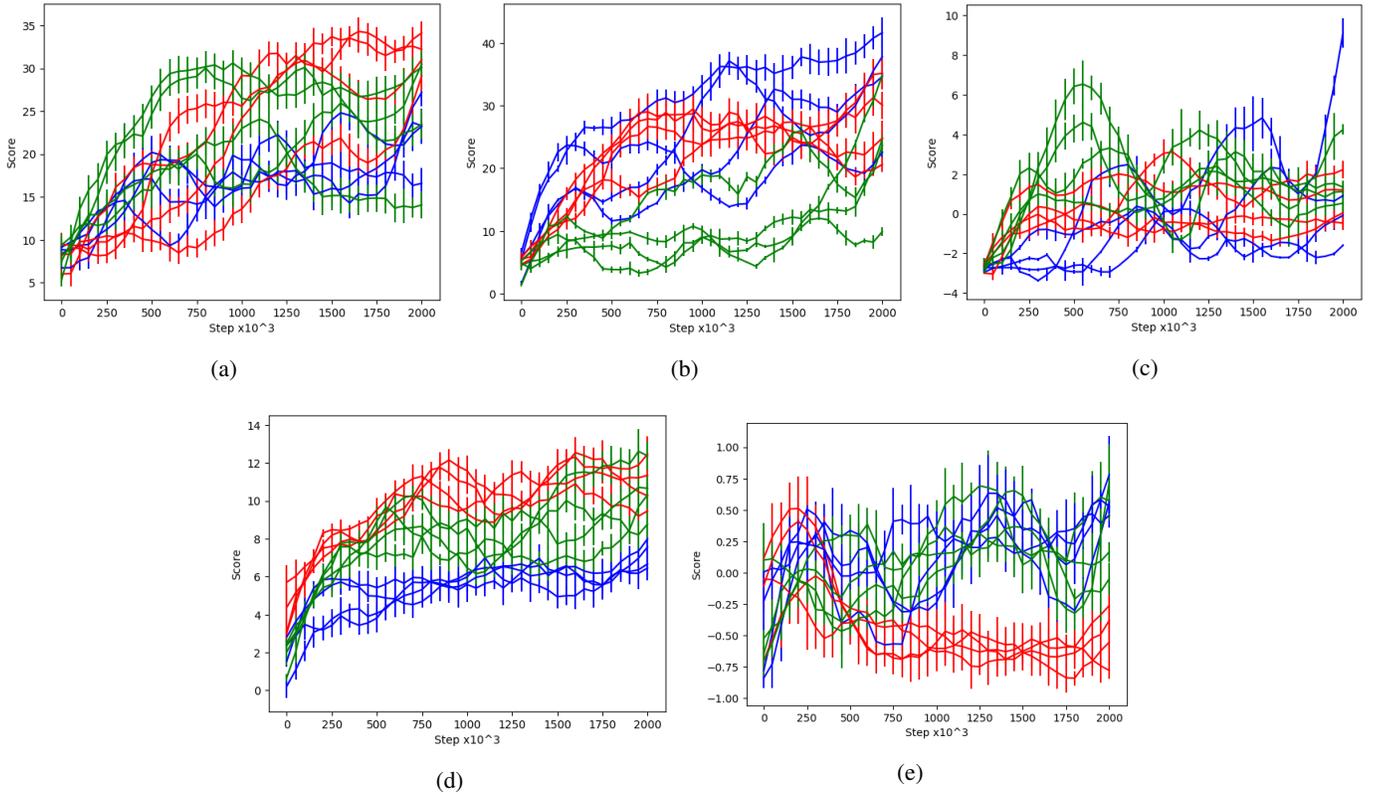


Fig. 5: Average episode rewards over 50 episodes at various points during training for five games: (a) Aliens. (b) Boulderdash. (c) Missile Command. (d) Survive Zombies. (e) Zelda. Results are smoothed using a fourth-order Savitzky-Golay filter with a window size of 21 to improve readability. Lines in Blue are agents with feature representation & fully connected network, Red are image representation & CNN, and Green are object representation & OEN (ours). Best viewed in colour.

TABLE I: Best mean score for each agent over 50 episodes.

	Aliens	Boulderdash	Missile Command	Survive Zombies	Zelda
Image	35.98	41.90	5.72	16.40	1.16
Features	29.82	43.34	9.12	10.42	1.58
Objects (ours)	37.30	35.90	10.44	14.88	1.68

agent in Table I, as these give an idea of the theoretical max performance of each agent type accounting for variability in agent parameters (although clearly these results are subject to sample bias, and are likely to be overestimates).

Due to the unpredictable nature of deep reinforcement learning we observed a large variance in agent performance between episodes but also between the average of different tests, making it difficult to compare individual agent results. Additionally, due to time and computational constraints we were only able to train for two million steps (comparable to eight million frames with frame-skip). This is significantly fewer than the tens and hundreds of millions of frames which many agents from the literature are trained for, meaning the results reported here may only be representative of the early stages of training. Nevertheless, we do observe certain patterns. In particular, there are notable difference in performance between representations. This is not surprising, as it is well known that choice of representation and network

architecture has a big impact on performance across other areas of deep learning. However, this difference in performance is not consistent across all games; different games seem to favour different representations. Surprisingly, the “features” baseline outperforms each other agents on certain games, despite often obscuring information about the game state (for example, the agent is given distances to certain objects, but not the direction to them, or the number of them).

It is observed from all the experimental results reported above that our object-based agent is capable of learning in all five games we tested it on. Additionally, across all the games, our agent performs comparably with the other two approaches. To this end, our experiments demonstrate that our OEN-based DRL agent can be an effective alternative to the existing agents for playing general video games.

VI. DISCUSSION

While we believe our method is generally widely applicable, a fundamental assumption made for our approach is that the game-state can be expressed in terms of objects. In some games object information may be unavailable, thus our approach cannot be applied in those games. Nevertheless, our object-based approach could work on those games, e.g., Starcraft, where image-based approaches, e.g., DQN [3], are difficult to apply due to unavailability of a complete visual-output

snapshot of the full game-state (e.g. due to presence of a controllable camera, or mouse-over options). In general, object extraction from a game-state can be done either directly via access to the ground-truth of the environment or from video or some other sources. When the ground-truth of the environment is available, object extraction is usually straightforward by utilising the object-oriented nature of most games (although this may also rely on certain domain knowledge to identify which objects are relevant to the agent). Otherwise, objects could be extracted directly from video data based on the state-of-the-art semantic image segmentation techniques. As a game environment is usually much simpler than natural images, existing semantic image segmentation techniques should be sufficient for this task. Another limitation of our approach is that it requires the user to find a fixed number of attributes to form feature vectors for all objects applicable to the OEN. In future this requirement for a fixed number of attributes across all object classes could be removed by pre-embedding all objects into a fixed-size space, or using class-specific embedding functions.

The use of relational information between objects is important when expressing the game-state. Indeed, Liang et.al. [25] showed that simple relationships between objects form a good feature set for reinforcement learning in video-games. In our work, we use only a simple method to exploit contextual information which has limited ability to capture these relations. In set networks, there are more sophisticated methods to exploit the contextual information, e.g., those used in [13], [14]. To achieve more effective unified object representations, our OEN model described in Section IV-B can be improved by adopting those techniques developed for set networks.

To conclude, we have presented a novel approach to learning directly from semi-structured object information via an OEN for playing general video games. A comparative study based on five GVG-AI competition games suggests that our approach yields performance comparable to two state-of-the-art approaches in general. In our ongoing research, we aim to address the issues and the limitations discussed above for improvement.

REFERENCES

- [1] A. Champeandard, "Making designers obsolete? evolution in game design," <http://aigamedev.com/open/interview/evolution-in-cityconquest/>, 2012, [Online; accessed 2016-02-18].
- [2] T. S. Nielsen, G. A. Barros, J. Togelius, and M. J. Nelson, "General video game evaluation using relative algorithm performance profiles," in *Applications of Evolutionary Computation*. Springer, 2015, pp. 369–380.
- [3] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski et al., "Human-level control through deep reinforcement learning," *Nature*, vol. 518, no. 7540, pp. 529–533, 2015.
- [4] G. Lample and D. S. Chaplot, "Playing fps games with deep reinforcement learning," 2017.
- [5] R. Dubey, P. Agrawal, D. Pathak, T. L. Griffiths, and A. A. Efros, "Investigating human priors for playing video games," *arXiv preprint arXiv:1802.10217*, 2018.
- [6] S. Samothrakis, D. Perez-Liebana, S. M. Lucas, and M. Fasli, "Neuroevolution for general video game playing," in *Computational Intelligence and Games (CIG), 2015 IEEE Conference on*. IEEE, 2015, pp. 200–207.
- [7] K. Kuanusont, S. M. Lucas, and D. Pérez-Liebana, "General video game ai: Learning from screen capture," in *Evolutionary Computation (CEC), 2017 IEEE Congress on*. IEEE, 2017, pp. 2078–2085.
- [8] K. Narasimhan, R. Barzilay, and T. Jaakkola, "Deep transfer in reinforcement learning by language grounding," *arXiv preprint arXiv:1708.00133*, 2017.
- [9] N. Topin, N. Haltmeyer, S. Squire, J. Winder, M. desJardins, and J. MacGlashan, "Portable option discovery for automated learning transfer in object-oriented markov decision processes." in *IJCAI*, 2015, pp. 3856–3864.
- [10] T. D. Kulkarni, K. Narasimhan, A. Saeedi, and J. Tenenbaum, "Hierarchical deep reinforcement learning: Integrating temporal abstraction and intrinsic motivation," in *Advances in neural information processing systems*, 2016, pp. 3675–3683.
- [11] M. Garnelo, K. Arulkumaran, and M. Shanahan, "Towards deep symbolic reinforcement learning," *arXiv preprint arXiv:1609.05518*, 2016.
- [12] M. Zaheer, S. Kottur, S. Ravanbakhsh, B. Poczos, R. R. Salakhutdinov, and A. J. Smola, "Deep sets," in *Advances in Neural Information Processing Systems*, 2017, pp. 3394–3404.
- [13] A. Santoro, D. Raposo, D. G. Barrett, M. Malinowski, R. Pascanu, P. Battaglia, and T. Lillicrap, "A simple neural network module for relational reasoning," in *Advances in neural information processing systems*, 2017, pp. 4974–4983.
- [14] N. Watters, A. Tacchetti, T. Weber, R. Pascanu, P. Battaglia, and D. Zoran, "Visual interaction networks," *arXiv preprint arXiv:1706.01433*, 2017.
- [15] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *International Conference on Machine Learning*, 2016, pp. 1928–1937.
- [16] A. Pritzel, B. Uria, S. Srinivasan, A. Puigdomènech, O. Vinyals, D. Hassabis, D. Wierstra, and C. Blundell, "Neural episodic control," *arXiv preprint arXiv:1703.01988*, 2017.
- [17] C. Tessler, S. Givony, T. Zahavy, D. J. Mankowitz, and S. Mannor, "A deep hierarchical approach to lifelong learning in minecraft." in *AAAI*, vol. 3, 2017, p. 6.
- [18] C. Diuk, A. Cohen, and M. L. Littman, "An object-oriented representation for efficient reinforcement learning," in *Proceedings of the 25th international conference on Machine learning*. ACM, 2008, pp. 240–247.
- [19] D. E. Hershkowitz, J. MacGlashan, and S. Tellex, "Learning propositional functions for planning and reinforcement learning," in *2015 AAAI Fall Symposium Series*, 2015.
- [20] C. R. Qi, H. Su, K. Mo, and L. J. Guibas, "Pointnet: Deep learning on point sets for 3d classification and segmentation," *arXiv preprint arXiv:1612.00593*, 2016.
- [21] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.
- [22] D. Perez-Liebana, J. Liu, A. Khalifa, R. D. Gaina, J. Togelius, and S. M. Lucas, "General video game ai: a multi-track framework for evaluating agents, games and content generation algorithms," *arXiv preprint arXiv:1802.10363*, 2018.
- [23] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling, "The arcade learning environment: An evaluation platform for general agents," *Journal of Artificial Intelligence Research*, 2012.
- [24] T. Schaul, "A video game description language for model-based or interactive learning," in *Computational Intelligence in Games (CIG), 2013 IEEE Conference on*. IEEE, 2013, pp. 1–8.
- [25] Y. Liang, M. C. Machado, E. Talvitie, and M. Bowling, "State of the art control of atari games using shallow reinforcement learning," *arXiv preprint arXiv:1512.01563*, 2015.