

# RECURRENT ENVIRONMENT SIMULATORS

**Silvia Chiappa, Sébastien Racaniere, Daan Wierstra & Shakir Mohamed**

DeepMind, London, UK

{csilvia, sracaniere, wierstra, shakir}@google.com

## ABSTRACT

Models that can simulate how environments change in response to actions can be used by agents to plan and act efficiently. We improve on previous environment simulators from high-dimensional pixel observations by introducing recurrent neural networks that are able to make temporally and spatially coherent predictions for hundreds of time-steps into the future. We present an in-depth analysis of the factors affecting performance, providing the most extensive attempt to advance the understanding of the properties of these models. We address the issue of computationally inefficiency with a model that does not need to generate a high-dimensional image at each time-step. We show that our approach can be used to improve exploration and is adaptable to many diverse environments, namely 10 Atari games, a 3D car racing environment, and complex 3D mazes.

## 1 INTRODUCTION

In order to plan and act effectively, agent-based systems require an ability to anticipate the consequences of their actions within an environment, often for an extended period into the future. Agents can be equipped with this ability by having access to models that can simulate how the environments changes in response to their actions. The need for environment simulation is widespread: in psychology, model-based predictive abilities form sensorimotor contingencies that are seen as essential for perception (O'Regan & Noë, 2001); in neuroscience, environment simulation forms part of deliberative planning systems used by the brain (Niv, 2009); and in reinforcement learning, the ability to imagine the future evolution of an environment is needed to form predictive state representations (Littman et al., 2002) and for Monte Carlo planning (Sutton & Barto, 1998).

Simulating an environment requires models of temporal sequences that must possess a number of properties to be useful: the models should make predictions that are accurate, temporally and spatially coherent over long time periods; and allow for flexibility in the policies and action sequences that are used. In addition, these models should be general-purpose and scalable, and able to learn from high-dimensional perceptual inputs and from diverse and realistic environments. A model that achieves these desiderata can empower agent-based systems with a vast array of abilities, including counterfactual reasoning (Pearl, 2009), intuitive physical reasoning (McCloskey, 1983), model-based exploration, episodic control (Lengyel & Dayan, 2008), intrinsic motivation (Oudeyer et al., 2007), and hierarchical control.

Deep neural networks have recently enabled significant advances in simulating complex environments, allowing for models that consider high-dimensional visual inputs across a wide variety of domains (Wahlström et al., 2015; Watter et al., 2015; Sun et al., 2015; Patraucean et al., 2015). The model of Oh et al. (2015) represents the state-of-the-art in this area, demonstrating high long-term accuracy in deterministic and discrete-action environments.

Despite these advances, there are still several challenges and open questions. Firstly, the properties of these simulators in terms of generalisation and sensitivity to the choices of model structure and training are poorly understood. Secondly, accurate prediction for long time periods into the future remains difficult to achieve. Finally, these models are computationally inefficient, since they require the prediction of a high-dimensional image each time an action is executed, which is unnecessary in situations where the agent is interested only in the final prediction after taking several actions.

In this paper we advance the state-of-the-art in environment modelling. We build on the work of Oh et al. (2015), and develop alternative architectures and training schemes that significantly improve performance, and provide in-depth analysis to advance our understanding of the properties of these

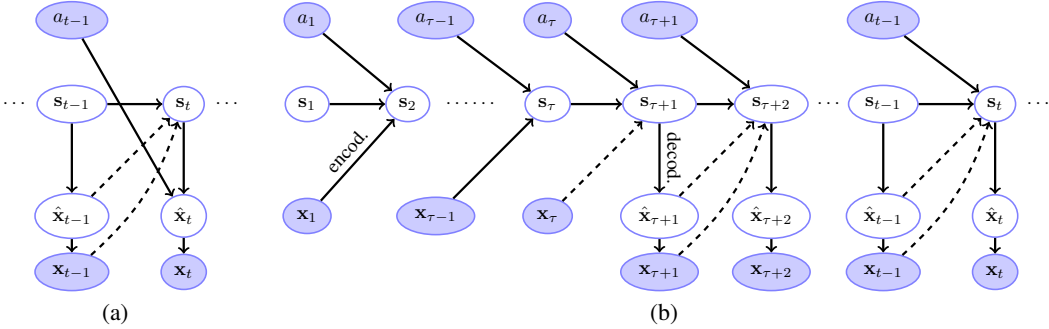


Figure 1: Graphical model representing (a) the recurrent structure used in Oh et al. (2015) and (b) our recurrent structure. Filled and empty nodes indicate observed and hidden variables respectively.

models. We also introduce a simulator that does not need to predict visual inputs after every action, reducing the computational burden in the use of the model. We test our simulators on three diverse and challenging families of environments, namely Atari 2600 games, a first-person game where an agent moves in randomly generated 3D mazes, and a 3D car racing environment; and show that they can be used for model-based exploration.

## 2 RECURRENT ENVIRONMENT SIMULATORS

An environment simulator is a model that, given a sequence of actions  $a_1, \dots, a_{\tau-1} \equiv a_{1:\tau-1}$  and corresponding observations  $\mathbf{x}_{1:\tau}$  of the environment, is able to predict the effect of subsequent actions  $a_{\tau:\tau+\tau'-1}$ , such as forming predictions  $\hat{\mathbf{x}}_{\tau+1:\tau+\tau'}$  or state representations  $\mathbf{s}_{\tau+1:\tau+\tau'}$  of the environment.

Our starting point is the recurrent simulator of Oh et al. (2015), which is the state-of-the-art in simulating deterministic environments with visual observations (frames) and discrete actions. This simulator is a recurrent neural network with the following backbone structure:

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}, \mathcal{C}(\mathbb{I}(\hat{\mathbf{x}}_{t-1}, \mathbf{x}_{t-1}))), \quad \hat{\mathbf{x}}_t = \mathcal{D}(\mathbf{s}_t, a_{t-1}).$$

In this equation,  $\mathbf{s}_t$  is a hidden state representation of the environment, and  $f$  a non-linear deterministic state transition function. The symbol  $\mathbb{I}$  indicates the selection of the predicted frame  $\hat{\mathbf{x}}_{t-1}$  or real frame  $\mathbf{x}_{t-1}$ , producing two types of state transition called **prediction-dependent transition** and **observation-dependent transition** respectively.  $\mathcal{C}$  is an encoding function consisting of a series of convolutions, and  $\mathcal{D}$  is a decoding function that combines the state  $\mathbf{s}_t$  with the action  $a_{t-1}$  through a multiplicative interaction, and then transforms it using a series of full convolutions to form the predicted frame  $\hat{\mathbf{x}}_t$ .

The model is trained to minimise the mean squared error between the observed time-series  $\mathbf{x}_{\tau+1:\tau+\tau'}$ , corresponding to the evolution of the environment, and its prediction. In a probabilistic framework, this corresponds to maximising the log-likelihood in the graphical model depicted in Fig. 1(a). In this graph, the link from  $\hat{\mathbf{x}}_t$  to  $\mathbf{x}_t$  represents stochastic dependence, as  $\mathbf{x}_t$  is formed by adding to  $\hat{\mathbf{x}}_t$  a Gaussian noise term with zero mean and unit variance, whilst all remaining links represent deterministic dependences. The dashed lines indicate that only one of the two links is active, depending on whether the state transition is prediction-dependent or observation-dependent.

The model is trained using stochastic gradient descent, in which each mini-batch consists of a set of segments of length  $\tau + T$  randomly sub-sampled from  $\mathbf{x}_{1:\tau+\tau'}$ . For each segment in the mini-batch, the model uses the first  $\tau$  observations to evolve the state and forms predictions of the last  $T$  observations only. *Training* comprises three phases differing in the use of prediction-dependent or observation-dependent transitions (after the first  $\tau$  transitions) and in the value of the *prediction length*  $T$ . In the first phase, the model uses observation-dependent transitions and predicts for  $T = 10$  time-steps. In the second and third phases, the model uses prediction-dependent transitions and predicts for  $T = 3$  and  $T = 5$  time-steps respectively. During *evaluation* or *usage*, the model can only use prediction-dependent transitions.

## ACTION-DEPENDENT STATE TRANSITION

A strong feature of the model of Oh et al. (2015) described above is that the actions influence the state transitions only *indirectly* through the predictions or the observations. Allowing the actions to condition the state transitions directly could potentially enable the model to incorporate action information more effectively. We therefore propose the following backbone structure:

$$\mathbf{s}_t = f(\mathbf{s}_{t-1}, a_{t-1}, \mathcal{C}(\mathbb{I}(\hat{\mathbf{x}}_{t-1}, \mathbf{x}_{t-1}))), \quad \hat{\mathbf{x}}_t = \mathcal{D}(\mathbf{s}_t).$$

In the graphical model representation, this corresponds to replacing the link from  $a_{t-1}$  to  $\hat{\mathbf{x}}_t$  with a link from  $a_{t-1}$  to  $\mathbf{s}_t$  as in Fig. 1(b).

## SHORT-TERM VERSUS LONG-TERM ACCURACY

The last two phases in the training scheme of Oh et al. (2015) described above are used to address the issue of poor accuracy that recurrent neural networks trained using only observation-dependent transitions display when asked to predict several time-steps ahead. However, the paper does not analyse nor discuss alternative training schemes.

In principle, the highest accuracy should be obtained by training the model as closely as possible to the way it will be used, and therefore by using a number of prediction-dependent transitions which is as close as possible to the number of time-steps the model will be asked to predict for. However, prediction-dependent transitions increase the complexity of the objective function such that alternative schemes are most often used (Talvitie, 2014; Bengio et al., 2015; Oh et al., 2015). Current training approaches are guided by the belief that using the observation  $\mathbf{x}_{t-1}$ , rather than the prediction  $\hat{\mathbf{x}}_{t-1}$ , to form the state  $\mathbf{s}_t$  has the effect of reducing the propagation of the errors made in the predictions, which are higher at earlier stages of the training, enabling the model to *correct* itself from the mistakes made up to time-step  $t - 1$ . For example, Bengio et al. (2015) introduce a scheduled sampling approach where at each time-step the type of state transition is sampled from a Bernoulli distribution, with parameter annealed from an initial value corresponding to using only observation-dependent transitions to a final value corresponding to using only prediction-dependent transitions, according to a schedule selected by validation.

Our analysis of different training schemes on Atari, which considered the interplay among *warm-up* length  $\tau$ , prediction length  $T$ , and number of prediction-dependent transitions, suggests that, rather than as having a corrective effect, observation-dependent transitions should be seen as restricting the time interval in which the model considers its predictive abilities, and therefore focuses resources. Indeed we found that, the higher the number of *consecutive* prediction-dependent transitions, the more the model is encouraged to focus on learning the global dynamics of the environment, which results in higher long-term accuracy. The highest long-term accuracy is always obtained by a training scheme that uses only prediction-dependent transitions even at the early stages of the training. Focussing on learning the global dynamics comes at the price of shifting model resources away from learning the precise details of the frames, leading to a decrease in short-term accuracy. Therefore, for complex games for which reasonable long-term accuracy cannot be obtained, training schemes that mix prediction-dependent and observation-dependent transitions are preferable. It follows from this analysis that percentage of consecutive prediction-dependent transitions, rather than just percentage of such transitions, should be considered when designing training schemes.

From this viewpoint, the poor results obtained in Bengio et al. (2015) when using only prediction-dependent transitions can be explained by the difference in the type of the tasks considered. Indeed, unlike our case in which the model is tolerant to some degree of error such as blurriness in earlier predictions, the discrete problems considered in Bengio et al. (2015) are such that one prediction error at earlier time-steps can severely affect predictions at later time-steps, so that the model needs to be highly accurate short-term in order to perform reasonably longer-term. Also, Bengio et al. (2015) treated the prediction used to form  $\mathbf{s}_t$  as a fixed quantity, rather than as a function of  $\mathbf{s}_{t-1}$ , and therefore did not perform exact maximum likelihood.

## PREDICTION-INDEPENDENT STATE TRANSITION

In addition to potentially enabling the model to incorporate action information more effectively, allowing the actions to directly influence the state dynamics has another crucial advantage: it allows to consider the case of a state transition that does not depend on the frame, *i.e.* of the form  $\mathbf{s}_t = f(\mathbf{s}_{t-1}, a_{t-1})$ , corresponding to removing the dashed links from  $\hat{\mathbf{x}}_{t-1}$  and from  $\mathbf{x}_{t-1}$  to  $\mathbf{s}_t$  in

Fig. 1(b). We shall call such a model *prediction-independent simulator*, referring to its ability to evolve the state without using the prediction during usage. Prediction-independent state transitions for high-dimensional observation problems have also been considered in Srivastava et al. (2015).

A prediction-independent simulator can dramatically increase computational efficiency in situations in which the agent is interested in the effect of a sequence of actions rather than of a single action. Indeed, such a model does not need to project from the lower dimensional state space into the higher dimensional observation space through the set of convolutions, and vice versa, at each time-step.

### 3 PREDICTION-DEPENDENT SIMULATORS

We analyse simulators with state transition of the form  $\mathbf{s}_t = f(\mathbf{s}_{t-1}, a_{t-1}, \mathcal{C}(\mathbb{I}(\hat{\mathbf{x}}_{t-1}, \mathbf{x}_{t-1})))$  on three families of environments with different characteristics and challenges, namely Atari 2600 games from the arcade learning environment (Bellemare et al., 2013), a first-person game where an agent moves in randomly generated 3D mazes (Beattie et al., 2016), and a 3D car racing environment called TORCS (Wymann et al., 2013). We use two evaluation protocols. In the first one, the model is asked to predict for 100 or 200 time-steps into the future using actions from the test data. In the second one, a human uses the model as an interactive simulator. The first protocol enables us to determine how the model performs within the action policy of the training data, whilst the second protocol enables us to explore how the model generalises to other action policies.

As state transition, we used the following *action-conditioned* long short-term memory (LSTM) (Hochreiter & Schmidhuber, 1997):

$$\text{Encoding: } \mathbf{z}_{t-1} = \mathcal{C}(\mathbb{I}(\hat{\mathbf{x}}_{t-1}, \mathbf{x}_{t-1})), \quad (1)$$

$$\text{Action fusion: } \mathbf{v}_t = \mathbf{W}^h \mathbf{h}_{t-1} \otimes \mathbf{W}^a \mathbf{a}_{t-1}, \quad (2)$$

$$\begin{aligned} \text{Gate update: } \mathbf{i}_t &= \sigma(\mathbf{W}^{iv} \mathbf{v}_t + \mathbf{W}^{iz} \mathbf{z}_{t-1}), \quad \mathbf{f}_t = \sigma(\mathbf{W}^{fv} \mathbf{v}_t + \mathbf{W}^{fz} \mathbf{z}_{t-1}), \\ \mathbf{o}_t &= \sigma(\mathbf{W}^{ov} \mathbf{v}_t + \mathbf{W}^{oz} \mathbf{z}_{t-1}), \end{aligned} \quad (3)$$

$$\text{Cell update: } \mathbf{c}_t = \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv} \mathbf{v}_t + \mathbf{W}^{cz} \mathbf{z}_{t-1}), \quad (4)$$

$$\text{State update: } \mathbf{h}_t = \mathbf{o}_t \otimes \tanh(\mathbf{c}_t), \quad (5)$$

where  $\otimes$  denotes the Hadamard product,  $\sigma$  the logistic sigmoid function,  $\mathbf{a}_{t-1}$  is a one-hot vector representation of  $a_{t-1}$ , and  $\mathbf{W}$  are parameter matrices. In Eqs. (2)–(5),  $\mathbf{h}_t$  and  $\mathbf{c}_t$  are the LSTM state and cell forming the model state  $\mathbf{s}_t = (\mathbf{h}_t, \mathbf{c}_t)$ ; and  $\mathbf{i}_t$ ,  $\mathbf{f}_t$ , and  $\mathbf{o}_t$  are the input, forget, and output gates respectively (for simplicity, we omit the biases in their updates). The vectors  $\mathbf{h}_t$  and  $\mathbf{v}_t$  had dimension 1024 and 2048 respectively. Details about the encoding and decoding functions  $\mathcal{C}$  and  $\mathcal{D}$  for the three families of environments can be found in Appendix B.1, B.2 and B.3. We used a warm-up phase of length  $\tau = 10$  and we did not backpropagate the gradient to this phase.

#### 3.1 ATARI

We considered the 10 Atari games Freeway, Ms Pacman, Qbert, Seaquest, Space Invaders, Bowling, Breakout, Fishing Derby, Pong, and Riverraid. Of these, the first five were analysed in Oh et al. (2015) and are used for comparison. The remaining five were chosen to better test the ability of the model in environments with other challenging characteristics, such as scrolling backgrounds (Riverraid), small/thin objects that are key aspects of the game (lines in Fishing Derby, ball in Pong and Breakout), and sparse-reward games that require very long-term predictions (Bowling). We used training and test datasets consisting of five and one million  $210 \times 160$  RGB images respectively, with actions chosen from a trained DQN agent (Mnih et al., 2015) according to an  $\epsilon = 0.2$ -greedy policy. Such a large number of training frames ensured that our simulators did not strongly overfit to the training data (see training and test lines in Figs. 2 and 3, and the discussion in Appendix B.1).

#### SHORT-TERM VERSUS LONG-TERM ACCURACY

Below we summarise our results on the interplay among warm-up length  $\tau$ , prediction length  $T$ , and number of prediction-dependent transitions – the full analysis is given in Appendix B.1.1.

The warm-up and prediction lengths  $\tau$  and  $T$  regulate degree of accuracy in two different ways. 1) The value of  $\tau + T$  determines how far into the past the model can access information – this is the case irrespectively of the type of transition used, although when using prediction-dependent transitions

information about the last  $T$  time-steps of the environment would need to be inferred. Accessing information far back into the past can be necessary even when the model is used to perform one-step ahead prediction only. 2) The higher the value of  $T$  and the number of prediction-dependent transitions, the more the corresponding objective function encourages long-term accuracy. This is achieved by guiding the one-step ahead prediction error in such a way that further predictions will not be strongly affected, and by teaching the model to make use of information from the far past. The more precise the model is in performing one-step ahead prediction, the less noise guidance should be required. Therefore, models with very accurate convolutional and transition structures should need less encouragement.

**Increasing the percentage of consecutive prediction-dependent transitions increases long-term accuracy, often at the expense of short-term accuracy.** We found that using only observation-dependent transitions leads to poor performance in most games. Increasing the number of consecutive prediction-dependent transitions produces an increase in long-term accuracy, but also a decrease in short-term accuracy usually corresponding to reduction in sharpness. For games that are too complex, although the lowest long-term prediction error is still achieved with using only prediction-dependent transitions, reasonable long-term accuracy cannot be obtained, and training schemes that mix prediction-dependent and observation-dependent transitions are therefore preferable.

To illustrate these results, we compare the following training schemes for prediction length  $T = 15$ :

- **0% PDT:** Only observation-dependent transitions.
- **33% PDT:** Observation and prediction-dependent transitions for the first 10 and last 5 time-steps respectively.
- **0%-20%-33% PDT:** Only observation-dependent transitions in the first 10,000 parameter updates; observation-dependent transitions for the first 12 time-steps and prediction-dependent transitions for the last 3 time-steps for the subsequent 100,000 parameters updates; observation-dependent transitions for the first 10 time-steps and prediction-dependent transitions for the last 5 time-steps for the remaining parameter updates (adaptation of the training scheme of Oh et al. (2015) to  $T = 15$ ).
- **46% PDT Alt.:** Alternate between observation-dependent and prediction-dependent transitions from a time-step to the next.
- **46% PDT:** Observation and prediction-dependent transitions for the first 8 and last 7 time-steps respectively.
- **67% PDT:** Observation and prediction-dependent transitions for the first 5 and last 10 time-steps respectively.
- **0%-100% PDT:** Only observation-dependent transitions in the first 1000 parameter updates; only prediction-dependent transitions in the subsequent parameter updates.
- **100% PDT:** Only prediction-dependent transitions.

For completeness, we also consider a training scheme as in Oh et al. (2015), which consists of three phases with  $T = 10, T = 3, T = 5$ , and 500,000, 250,000, 750,000 parameter updates respectively. In the first phase  $\mathbf{s}_t$  is formed by using the observed frame  $\mathbf{x}_{t-1}$ , whilst in the two subsequent phases  $\mathbf{s}_t$  is formed by using the predicted frame  $\hat{\mathbf{x}}_{t-1}$ .

In Figs. 2 and 3 we show the prediction error averaged over 10,000 sequences<sup>1</sup> for the games of Bowling<sup>2</sup>, Fishing Derby, Pong and Seaquest. More specifically, Fig. 2(a) shows the error for predicting up to 100 time-steps ahead after the model has seen 200 million frames (corresponding to half million parameter updates using mini-batches of 16 sequences), using actions and warm-up frames from the test data, whilst Figs. 2(b)-(c) and 3 show the error at time-steps 5, 10 and 100 versus number of frames seen by the model.

These figures clearly show that long-term accuracy generally improves with increasing number of consecutive prediction-dependent transitions. When using alternating (46% PDT Alt.), rather than consecutive (46% PDT), prediction-dependent transitions long-term accuracy is worse, as we are effectively asking the model to predict at most two time-steps ahead. We can also see that

<sup>1</sup>We define the prediction error as  $\frac{1}{3 \times 10,000} \sum_{n=1}^{10,000} \|\mathbf{x}_t^n - \hat{\mathbf{x}}_t^n\|^2$ .

<sup>2</sup>In this game, the player is given two chances to roll a ball down an alley in an attempt to knock down as many of the ten pins as possible, after which the score is updated and the knocked pins are relocated. Knocking down every pin on the first shot is a strike, while knocking every pin down in both shots is a spare. The player's score is determined by the number of pins knocked down, as well as the number of strikes and spares acquired.

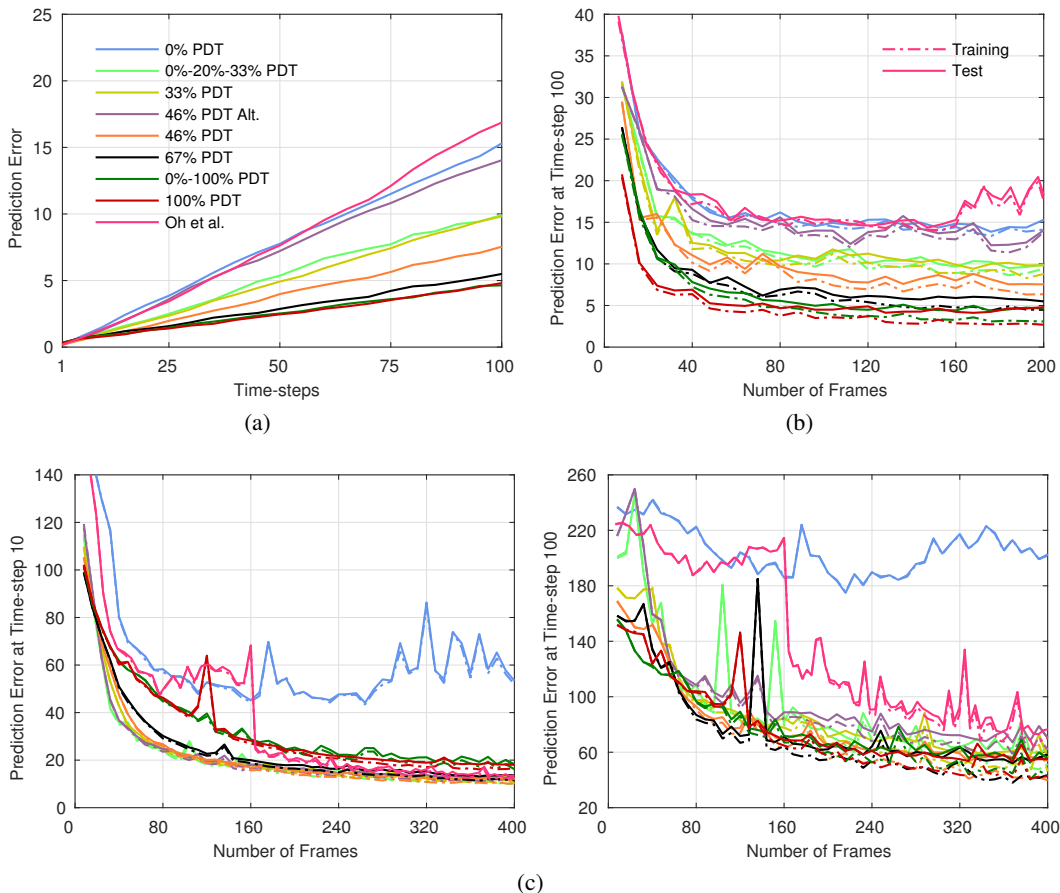


Figure 2: Prediction error averaged over 10,000 sequences on (a)-(b) Bowling and (c) Fishing Derby for different training schemes. The same color and line code is used in all figures. (a): Prediction error vs time-steps after the model has seen 200 million frames. (b)-(c): Prediction error vs number of frames seen by the model at time-steps 10 and 100.

using more prediction-dependent transitions produces lower short-term accuracy and/or slower short-term convergence. Finally, the figures show that using a training phase with only observation-dependent transitions that is too long, as in Oh et al. (2015), can be detrimental: the models reaches at best a performance similar to the 46% PDT Alt. training scheme (the sudden drop in prediction error corresponds to transitioning to the second training phase), but is most often worse.

By looking at the predicted frames we could notice that, in games containing balls and paddles, using only observation-dependent transitions gives rise to errors in reproducing the dynamics of these objects. Such errors decrease with increasing prediction-dependent transitions. In other games, using only observation-dependent transitions causes the model to fail in representing moving objects, except for the agent in most cases. Training schemes containing more prediction-dependent transitions encourage the model to focus more on learning the dynamics of the moving objects and less on details that would only increase short-term accuracy, giving rise to more globally accurate but less sharp predictions. Finally, in games that are too complex, the strong emphasis on long-term accuracy produces predictions that are overall not sufficiently good.

More specifically, from the videos available at<sup>3</sup> [PDTvsODT](#), we can see that using only observation-dependent transitions has a detrimental effect on long-term accuracy for Fishing Derby, Ms Pacman,

<sup>3</sup>Highlighted names like these are direct links to folders containing videos. Each video consists of 5 randomly selected 200 time-steps ahead predictions separated by black frames (the generated frames are shown on the left, whilst the real frames are shown on the right – the same convention will be used throughout the paper). Shown

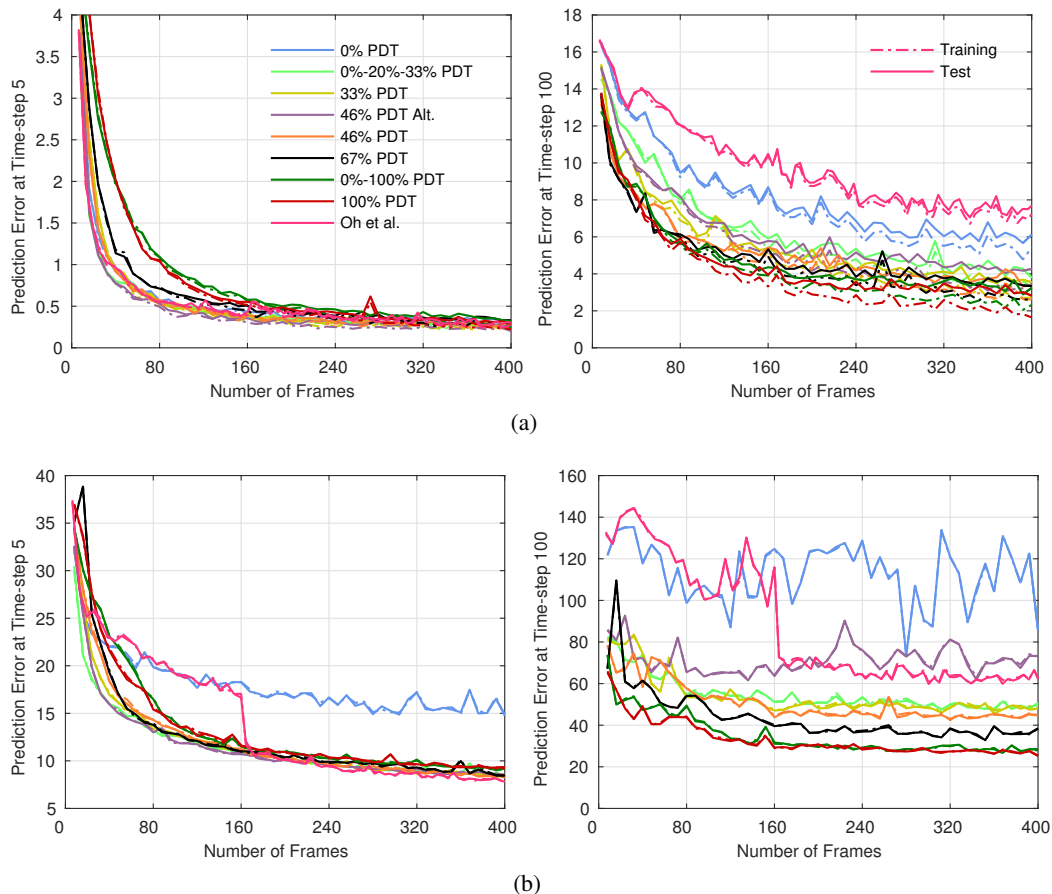


Figure 3: Prediction error on (a) Pong and (b) Seaquest for different training schemes.

Qbert, Riverraid, Seaquest and Space Invaders. The most salient features of the videos are: consistent inaccuracy in predicting the paddle and ball in Breakout; reset to a new life after a few time-steps in Ms Pacman; prediction of background only after a few time-steps in Qbert; no generation of new objects or background in Riverraid; quick disappearance of existing fish and no appearance of new fish from the sides of the frame in Seaquest. For Bowling, Freeway, and Pong, long-term accuracy is generally good, but the movement of the ball is not always correctly predicted in Bowling and Pong and the chicken sometimes disappears in Freeway. On the other hand, using only prediction-dependent transitions results in good long-term accuracy for Bowling, Fishing Derby, Freeway, Pong, Riverraid, and Seaquest: the model accurately represents the paddle and ball dynamics in Bowling and Pong; the chicken hardly disappears in Freeway, and new objects and background are created and most often correctly positioned in Riverraid and Seaquest.

The trading-off of long for short-term accuracy when using more prediction-dependent transitions is particularly evident in the videos of Seaquest: the higher the number of such transitions, the better the model learns the dynamics of the game, with new fish appearing in the right location more often. However, this comes at the price of reduced sharpness, mostly in representing the fish.

This trade-off causes problems in Breakout, Ms Pacman, Qbert, and Space Invaders, so that schemes that also use observation-dependent transitions are preferable for these games. For example, in Breakout, the model fails at representing the ball, making the predictions not sufficiently good. Notice that the prediction error (see Fig. 15) is misleading in terms of desired performance, as the 100%PDT training scheme performs as well as other mixing schemes for long-term accuracy – this highlights the difficulties in evaluating the performance of these models.

are 15 frames per seconds. Videos associated with the material discussed in this and following sections can also be found at <https://sites.google.com/site/resvideos1729>.

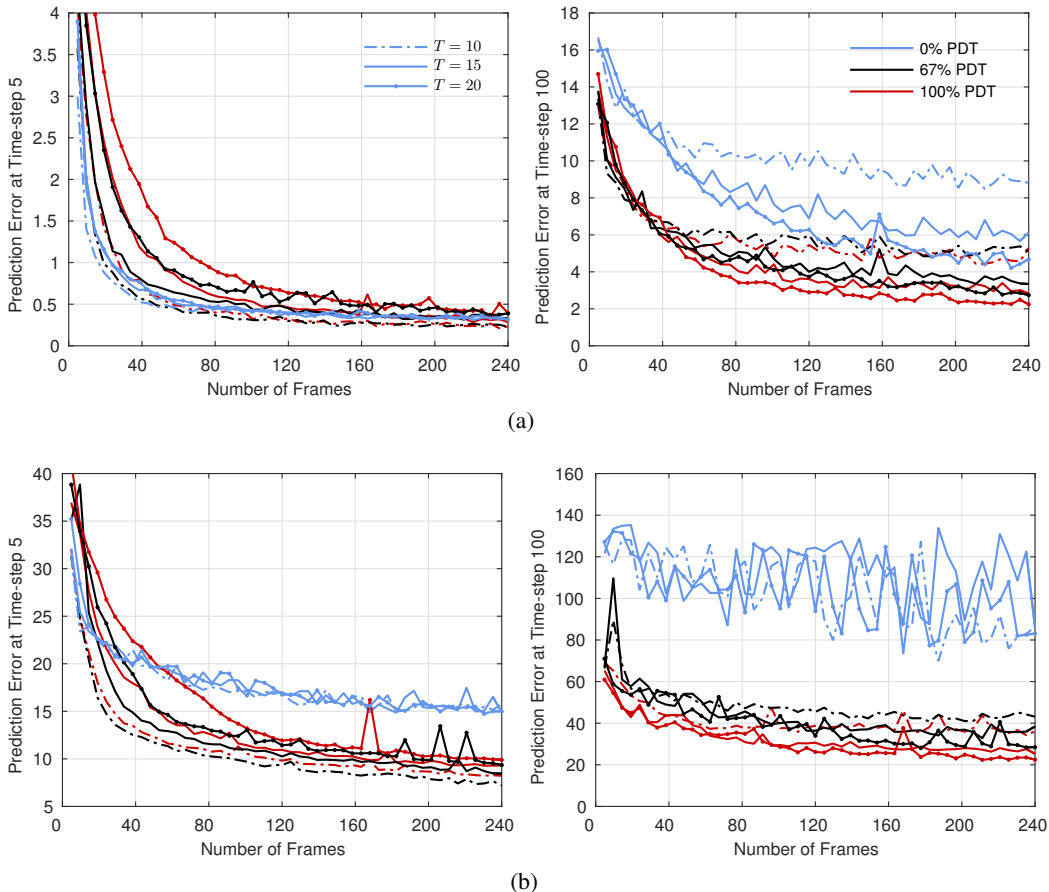


Figure 4: Prediction error vs number of frames seen by the model (excluding warm-up frames) for (a) Pong and (b) Seaquest, using prediction lengths  $T = 10, 15$ , and  $20$ , and training schemes 0%PDT, 67%PDT, and 100%PDT.

**Increasing the prediction length  $T$  increases long-term accuracy when using prediction-dependent transitions.** In Fig. 4, we show the effect of using different prediction lengths  $T \leq 20$  on the training schemes 0%PDT, 67%PDT, and 100%PDT for Pong and Seaquest. In Pong, with the 0%PDT training scheme, using higher  $T$  improves long-term accuracy: this is a game for which this scheme gives reasonable accuracy and the model is able to benefit from longer history. This is however not the case for Seaquest (or other games as shown in Appendix B.1.1). On the other hand, with the 100%PDT training scheme, using higher  $T$  improves long-term accuracy in most games (the difference is more pronounced between  $T = 10$  and  $T = 15$  than between  $T = 15$  and  $T = 20$ ), but decreases short-term accuracy. Similarly to above, reduced short-term accuracy corresponds to reduced sharpness: from the videos available at  $T < 20$  we can see, for example, that the moving caught fish in Fishing Derby, the fish in Seaquest, and the ball in Pong are less sharp for higher  $T$ .

**Truncated backpropagation still enables increase in long-term accuracy.** Due to memory constraints, we could only backpropagate gradients over sequences of length up to 20. To use  $T > 20$ , we split the prediction sequence into subsequences and performed parameter updates separately for each subsequence. For example, to use  $T = 30$  we split the prediction sequence into two successive subsequences of length 15, performed parameter updates over the first subsequence, initialised the state of the second subsequence with the final state from the first subsequence, and then performed parameter updates over the second subsequence. This approach corresponds to a form of truncated backpropagation through time (Williams & Zipser, 1995) – the extreme of this strategy (with  $T$  equal to the length of the whole training sequence) was used by Zaremba et al. (2014).



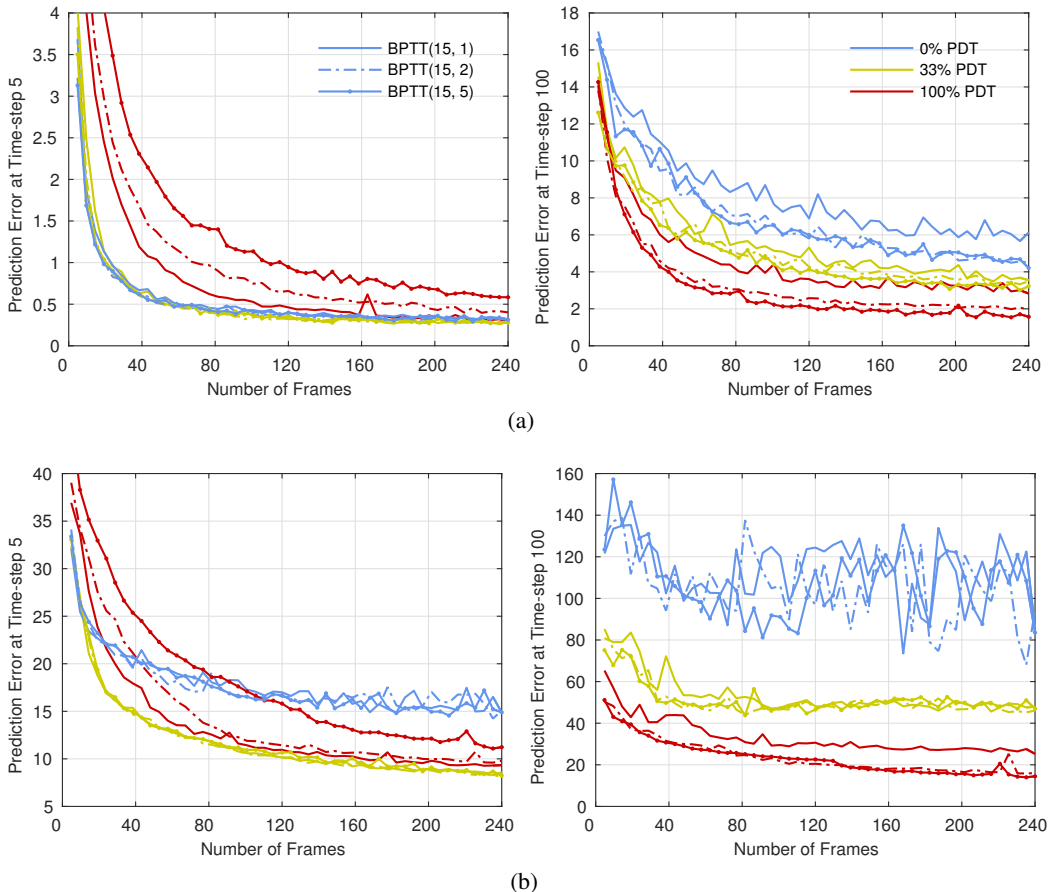


Figure 5: Prediction error vs number of frames seen by the model (excluding warm-up frames) for (a) Pong and (b) Seaquest, using BPTT(15, 1), BPTT(15, 2), and BTT(15, 5), and training schemes 0%PDT, 33%PDT, and 100%PDT.

In Fig. 5, we show the effect of using 2 and 5 subsequences of length 15 (indicated by BPTT(15, 2) and BTT(15, 5)) on the training schemes 0%PDT, 33%PDT, and 100%PDT for Pong and Seaquest. We can see that the 0%PDT and 33%PDT training schemes display no difference in accuracy for different values of  $T$ . On the other hand, with the 100%PDT training scheme, using more than one subsequence improves long-term accuracy (the difference is more pronounced between  $T = 15$  and  $T = 30$  than between  $T = 30$  and  $T = 75$ ), but decreases short-term accuracy (the difference is small at convergence between  $T = 15$  and  $T = 30$ , but big between  $T = 30$  and  $T = 75$ ). The decrease in accuracy with 5 subsequences is drastic in some games.

For Riverraid, using more than one subsequence with the 33%PDT and 100%PDT training schemes improves long-term accuracy dramatically, as shown in Fig. 6, as it enables correct prediction after a jet loss. Interestingly, for the 100%PDT training scheme, using  $\tau = 25$  with prediction length  $T = 15$  (black line) does not give the same amount of gain as when using BPTT(15, 2), even if history length  $\tau + T$  is the same. This would seem to suggest that some improvement in BPTT(15, 2) is due to encouraging longer-term accuracy, indicating that this can be achieved even when not fully backpropagating the gradient.

From the videos available at  $T > 20$ , we can see that with  $T = 75$  the predictions in some of the Fishing Derby videos are faded, whilst in Pong the model can suddenly switch from one dynamics to another for the ball and the opponent’s paddle.

In conclusion, using higher  $T$  through truncated backpropagation can improve performance. However, in schemes that use many prediction-dependent transitions, a high value of  $T$  can lead to poor predictions.

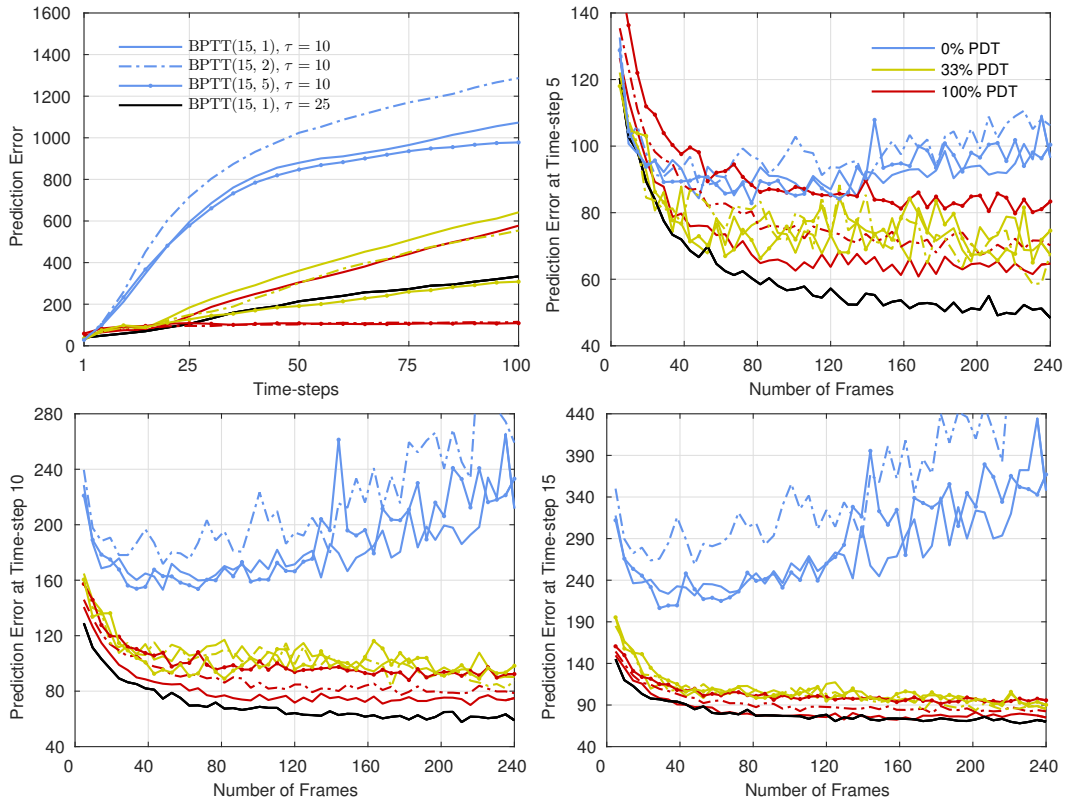


Figure 6: Prediction error vs number of frames seen by the model for Riverraid, using BPTT(15, 1), BPTT(15, 2), and BTT(15, 5), and training schemes 0%PDT, 33%PDT, and 100%PDT. The black line is obtained with the 100%PDT training scheme.

### EVALUATION THROUGH HUMAN PLAY

Whilst we cannot expect our simulators to generalise to structured sequences of actions never chosen by the DQN and that are not present in the training data, such as moving the agent up and down the alley in Bowling, it is reasonable to expect some degree of generalisation in the action-wise simple environments of Breakout, Freeway and Pong.

We tested these three games by having humans using the models as interactive simulators. We generally found that models trained using only prediction-dependent transitions were more fragile to states of the environment not experienced during training, such that the humans were able to play these games for longer with simulators trained with mixing training schemes. This seems to indicate that models with higher long-term test accuracy are at higher risk of overfitting to the training policy.

In Fig. 7(a), we show some salient frames from a game of Pong played by a human for 500 time-steps (the corresponding video is available at [Pong-HPlay](#)). The game starts with score (2,0), after which the opponent scores five times, whilst the human player scores twice. As we can see, the scoring is updated correctly and the game dynamics is accurate. In Fig. 7(b), we show some salient frames from a game of Breakout played by a human for 350 time-steps (the corresponding video is available at [Breakout-HPlay](#)). As for Pong, the scoring is updated correctly and the game dynamics is accurate. These images demonstrate some degree of generalisation of the model to a human style of play.

### EVALUATION OF STATE TRANSITIONS STRUCTURES

In Appendix B.1.2 and B.1.3 we present an extensive evaluation of different action-dependent state transitions, including convolutional transformations for the action fusion, and gate and cell updates, and different ways of incorporating action information. We also present a comparison between action-dependent and action-independent state transitions.

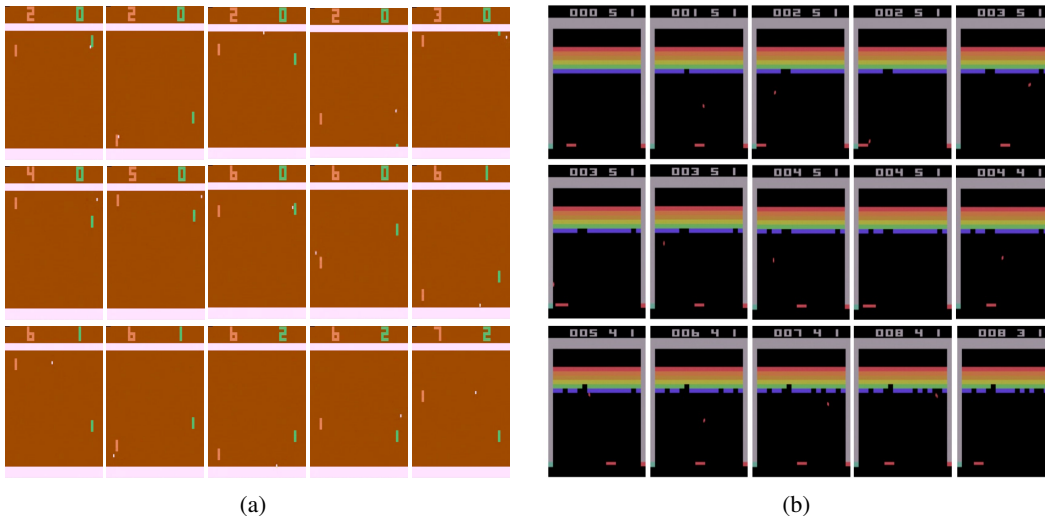


Figure 7: Salient frames extracted from (a) 500 frames of Pong and (b) 350 frames of Breakout generated using our simulator with actions taken by a human player (larger versions can be found in Figs. 47 and 48).

Some action-dependent state transitions give better performance than the baseline (Eqs. (1)–(5)) in some games. For example, we found that increasing the state dimension from 1024 to the dimension of the convolved frame, namely 2816, might be preferable. Interestingly, this is not due to an increase in the number of parameters, as the same gain is obtained using convolutions for the gate and cell updates. These results seem to suggest that high-dimensional sparse state transition structures could be a promising direction for further improvement. Regarding different ways of incorporation action information, we found that using local incorporation such as augmenting the frame with action information and indirect action influence gives worse performance than direct and global action influence, but that there are several ways of incorporating action information directly and globally that give similar performance.

### 3.2 3D ENVIRONMENTS

Both TORCS and the 3D maze environments highlight the need to learn dynamics that are temporally and spatially coherent: TORCS exposes the need to learn fast moving dynamics and consistency under motion, whilst 3D mazes are partially-observed and therefore require the simulator to build an internal representation of its surrounding using memory, as well learn basic physics, such as rotation, momentum, and the solid properties of walls.

**TORCS.** The data was generated using an artificial agent controlling a fast car without opponents (more details are given in Appendix B.2).

When using actions from the test set (see Fig. 49 and the corresponding video at [TORCS](#)), the simulator was able to produce accurate predictions for up to several hundreds time-steps. As the car moved around the racing track, the simulator was able to predict the appearance of new features in the background (towers, sitting areas, lamp posts, etc.), as well as model the jerky motion of the car caused by our choices of random actions. Finally, the instruments (speedometer and rpm) were correctly displayed.

The simulator was good enough to be used interactively for several hundred frames, using actions provided by a human. This showed that the model had learnt well how to deal with the car hitting the wall on the right side of the track. Some salient frames from the game are shown in Fig. 8 (the corresponding video can be seen at [TORCS-HPlay](#)).

**3D Mazes.** We used an environment that consists of randomly generated 3D mazes, containing textured surfaces with occasional paintings on the walls: the mazes were all of the same size, but

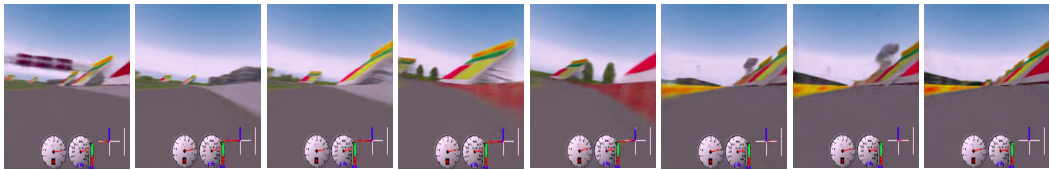


Figure 8: Salient frames highlighting coherence extracted from 700 frames of TORCS generated using our simulator with actions taken by a human player.



Figure 9: Predicted (left) and real (right) frames at time-steps 1, 25, 66, 158 and 200 using actions from the test data.

differed in the layout of rooms and corridors, and in the locations of paintings (see Fig. 11(b) for an example of layout). More details are given in Appendix B.3.

When using actions from the test set, the simulator was able to very reasonably predict frames even after 200 steps. In Fig. 9 we compare predicted frames to the real frames at several time-steps (the corresponding video can be seen at [3DMazes](#)). We can see that the wall layout is better predicted when walls are closer to the agent, and that corridors and far away-walls are not as long as they should be. The lighting on the ceiling is correct on all the frames shown.

When using the simulator interactively with actions provided by a human, we could test that the simulator had learnt consistent aspects of the maze: when walking into walls, the model maintained their position and layout (in one case we were able to walk through a painting on the wall – paintings are rare in the dataset and hence it is not unreasonable that they would not be maintained when stress testing the model in this way). When taking 360° spins, the wall configurations were the same as previously generated and not regenerated afresh, and shown in Fig. 10 (see also [3DMazes-HPLay](#)). The coherence of the maze was good for nearby walls, but not at the end of long-corridors.

### 3.3 MODEL-BASED EXPLORATION

The search for exploration strategies better than  $\epsilon$ -greedy is an active area of research. Various solutions have been proposed, such as density based or optimistic exploration (Auer et al., 2002). Oh et al. (2015) considered a memory-based approach that steers the agent towards previously unobserved frames. In this section, we test our simulators using a similar approach, but select a group of actions rather than a single action at a time. Furthermore, rather than a fixed 2D environment, we consider the more challenging 3D mazes environment. This also enables us this present a qualitative analysis, as we can exactly measure and plot the proportion of the maze visited over time. Our aim is to be quantitatively and qualitatively better than random exploration (using dithering of 0.7, as this lead to the best possible random agent).

We used a 3D maze simulator to predict the outcome of sequences of actions, chosen with a hard-coded policy. Our algorithm (see below) did  $N$  Monte-Carlo simulations with randomly selected sequences of actions of fixed length  $d$ . At each time-step  $t$ , we stored the last 10 observed frames in an episodic memory buffer and compared predicted frames to those in memory.

```

for  $t = 1, \text{episodeLength}, d$  do
  for  $n = 1, N$  do
    Choose random actions  $A^n = a_{t:t+d-1}$ ;
    Predict  $\hat{\mathbf{x}}_{t+1:t+d}^n$ ;
  end
  Follow actions in  $A^{n_0}$  where
   $n_0 = \text{argmax}_n \min_{j=0,10} \|\hat{\mathbf{x}}_{t+d}^n - \mathbf{x}_{t-j}\|$ 
end

```

Our method (see Fig. 11(a)) covered 50% more of the maze area after 900 time-steps than random exploration. These results were obtained with 100 Monte-Carlo simulations and sequences of 6 actions (more details are given in Appendix B.4). Comparing typical paths chosen by the random explorer and by our explorer (see Fig. 11(b)), we see the our explorer has much smoother trajectories.

This is a good local exploration strategy that leads to faster movement through corridors. To

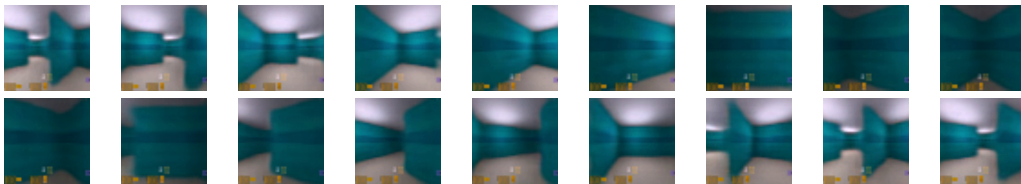


Figure 10: Salient frames highlighting wall-layout memory after 360° spin generated using our simulator with actions taken by a human player.

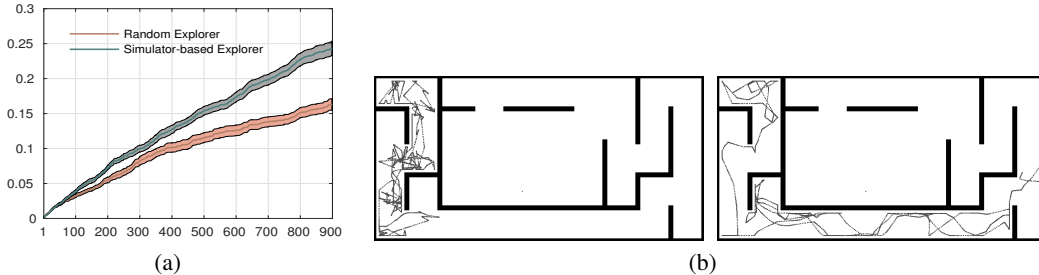


Figure 11: (a) Average ratio over 10 mazes (shaded is the 68% confidence interval) of area visited by the random agent and an agent using our model. (b) Typical example of paths followed by (left) the random agent and (right) our agent (see the Appendix for more examples).

transform this into a good global exploration strategy, our explorer would have to be augmented with a better memory in order to avoid going down the same corridor twice. These sorts of smooth local exploration strategies could also be useful in navigation problems.

#### 4 PREDICTION-INDEPENDENT SIMULATORS

A prediction-independent simulator has state transitions of the form  $\mathbf{s}_t = f(\mathbf{s}_{t-1}, a_{t-1})$ , which therefore do not require the high-dimensional predictions. In the Atari environment, for example, this avoids having to project from the state space of dimension 1024 into the observation space of dimension 100,800 ( $210 \times 160 \times 3$ ) through the decoding function  $\mathcal{D}$ , and vice versa through the encoding function  $\mathcal{C}$  – in the used structure this enables saving around 200 million flops at each time-step.

For the state transition, we found that a working structure was to use Eqs. (1)–(5) with  $\mathbf{z}_t = \mathbf{h}_t$  and with different parameters for the warm-up and prediction phases. As for the prediction-dependent simulator, we used a warm-up phase of length  $\tau = 10$ , but we did backpropagate the gradient back to time-step five in order to learn the encoding function  $\mathcal{C}$ .

Our analysis on Atari (see Appendix C) suggests that the prediction-independent simulator is much more sensitive to changes in the state transition structure and in the training scheme than the prediction-dependent simulator. We found that using prediction length  $T = 15$  gave much worse long-term accuracy than with the prediction-dependent simulator. This problem could be alleviated with the use of prediction length  $T = 30$  through truncated backpropagation.

Fig. 12 shows a comparison of the prediction-dependent and prediction-independent simulators using  $T = 30$  through two subsequences of length 15 (we indicate this as BPTT(15, 2), even though in the prediction-independent simulator we did backpropagate the gradient to the warm-up phase).

When looking at the videos available at [PI-Simulators](#), we can notice that the prediction-independent simulator tends to give worse type of long-term prediction. In Fishing Derby for example, in the long-term the model tends to create fish of smaller dimension in addition to the fish present in the real frames. Nevertheless, for some difficult games the prediction-independent simulator achieves better performance than the prediction-dependent simulator. More investigation about alternative

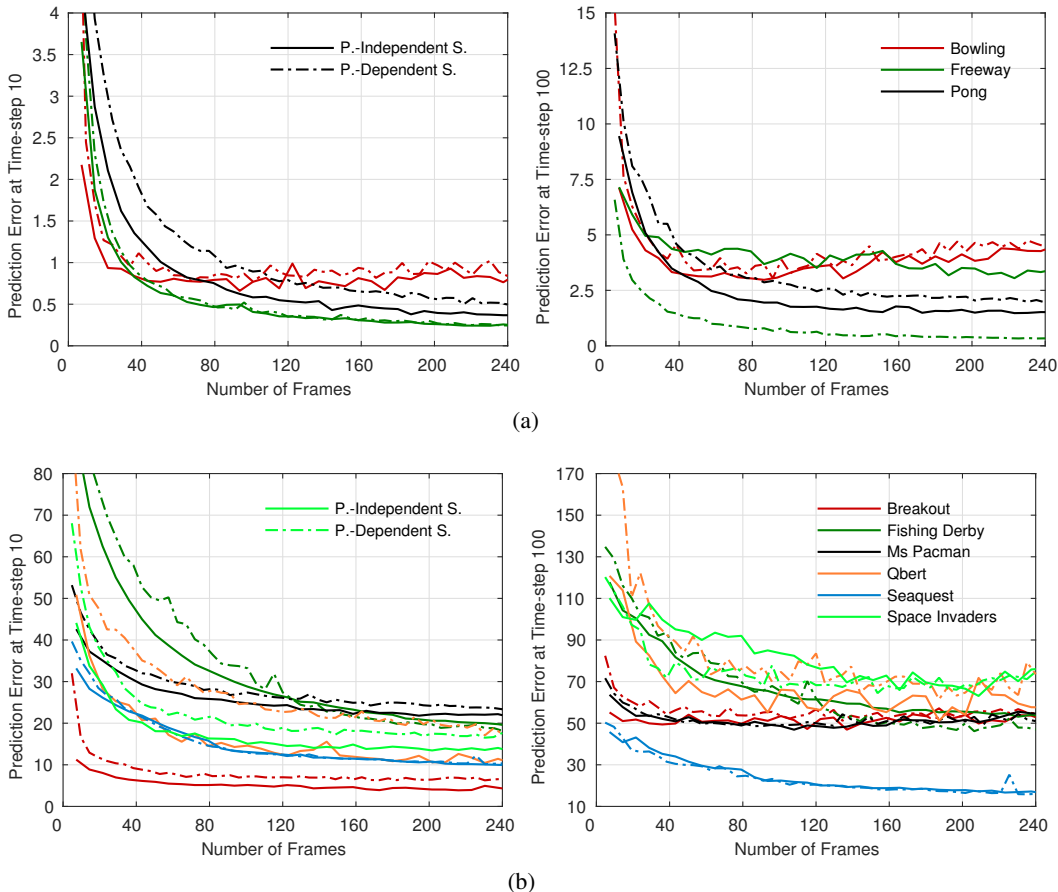


Figure 12: Prediction error vs number of frames seen by the model (excluding warm-up frames) for the prediction-dependent and prediction-independent simulators using BPTT(15, 2) for (a) Bowling, Freeway, Pong and (b) Breakout, Fishing Derby, Ms Pacman, Qbert, Seaquest, Space Invaders (the prediction-dependent simulator is trained with the 0%-100%PDT training scheme).

state transitions and training schemes would need to be performed to obtain the same overall level of accuracy as with the prediction-dependent simulator.

## 5 DISCUSSION

In this paper we have introduced an approach to simulate action-conditional dynamics and demonstrated that is highly adaptable to different environments, ranging from Atari games to 3D car racing environments and mazes. We showed state-of-the-art results on Atari, and demonstrated the feasibility of live human play in all three task families. The system is able to capture complex and long-term interactions, and displays a sense of spatial and temporal coherence that has, to our knowledge, not been demonstrated on high-dimensional time-series data such as these.

We have presented an in-deep analysis on the effect of different training approaches on short and long-term prediction capabilities, and showed that moving towards schemes in which the simulator relies less on past observations to form future predictions has the effect on focussing model resources on learning the global dynamics of the environment, leading to dramatic improvements in the long-term predictions. However, this requires a distribution of resources that impacts short-term performance, which can be harmful to the overall performance of the model for some games. This trade-off is also causing the model to be less robust to states of the environment not seen during training. To alleviate this problem would require the design of more sophisticated model architectures than the ones considered here. Whilst it is also expected that more ad-hoc architectures would be less sensitive

to different training approaches, we believe that guiding the noise as well as teaching the model to make use of past information through the objective function would still be beneficial for improving long-term prediction.

Complex environments have compositional structure, such as independently moving objects and other phenomena that only rarely interact. In order for our simulators to better capture this compositional structure, we may need to develop specialised functional forms and memory stores that are better suited to dealing with independent representations and their interlinked interactions and relationships. More homogeneous deep network architectures such as the one presented here are clearly not optimal for these domains, as can be seen in Atari environments such as Ms Pacman where the system has trouble keeping track of multiple independently moving ghosts. Whilst the LSTM memory and our training scheme have proven to capture long-term dependencies, alternative memory structures are required in order, for example, to learn spatial coherence at a more global level than the one displayed by our model in the 3D mazes in order to do navigation.

In the case of action-conditional dynamics, the policy-induced data distribution does not cover the state space and might in fact be nonstationary over an agent lifetime. This can cause some regions of the state space to be oversampled, whereas the regions we might actually care about the most – those just *around* the agent policy state distribution – to be underrepresented. In addition, this induces biases in the data that will ultimately not enable the model learn the environment dynamics correctly. As verified from the experiments in this paper, both on live human play and model-based exploration, this problem is not yet as pressing as might be expected in some environments. However, our simulators displayed limitations and faults due to the specificities of the training data, such as for example predicting an event based on the recognition of a particular sequence of actions always co-occurring with this event in the training data rather than on the recognition of the real causes.

Finally, a limitation of our approach is that, however capable it might be, it is a deterministic model designed for deterministic environments. Clearly most real world environments involve noisy state transitions, and future work will have to address the extension of the techniques developed in this paper to more generative temporal models.

#### ACKNOWLEDGMENTS

The authors would like to thank David Barber for helping with the graphical model interpretation, Alex Pritzel for preparing the DQN data, Yori Zwols and Frederic Besse for helping with the implementation of the model, and Oriol Vinyals, Yee Whye Teh, Junhyuk Oh, and the anonymous reviewers for useful discussions and feedback on the manuscript.

#### REFERENCES

- P. Auer, N. Cesa-Bianchi, and P. Fischer. Finite-time analysis of the multiarmed bandit problem. *Machine Learning*, 47:235–256, 2002.
- C. Beattie, J. Z. Leibo, D. Teplyashin, T. Ward, M. Wainwright, H. Küttler, A. Lefrancq, S. Green, V. Valdés, A. Sadik, J. Schrittwieser, K. Anderson, S. York, M. Cant, A. Cain, A. Bolton, S. Gaffney, H. King, D. Hassabis, S. Legg, and S. Petersen. Deepmind lab. *CoRR*, abs/1612.03801, 2016. URL <http://arxiv.org/abs/1612.03801>.
- 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*, 47:253–279, 2013.
- S. Bengio, O. Vinyals, N. Jaitly, and N. Shazeer. Scheduled sampling for sequence prediction with recurrent neural networks. In *Advances in Neural Information Processing Systems 28 (NIPS)*, pp. 1171–1179. 2015.
- A. Graves. Generating sequences with recurrent neural networks. 2013. URL <http://arxiv.org/abs/1308.0850>.
- S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- M. Lengyel and P. Dayan. Hippocampal contributions to control: The third way. In *Advances in Neural Information Processing Systems 20 (NIPS)*, pp. 889–896, 2008.
- M. L. Littman, R. S. Sutton, and S. Singh. Predictive representations of state. In *Advances in Neural Information Processing Systems 14 (NIPS)*, pp. 1555–1561. 2002.
- M. McCloskey. Intuitive physics. *Scientific American*, 248(4):122–130, 1983.
- V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis. Human-level control through deep reinforcement learning. *Nature*, 518(7540): 529–533, 02 2015. URL <http://dx.doi.org/10.1038/nature14236>.

- V. Mnih, A. Puigdomènech Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu. Asynchronous methods for deep reinforcement learning. In *Proceedings of the 33rd International Conference on Machine Learning (ICML)*, 2016.
- Y. Niv. Reinforcement learning in the brain. *Journal of Mathematical Psychology*, 53(3):139–154, 2009.
- J. Oh, X. Guo, H. Lee, R. L. Lewis, and S. P. Singh. Action-conditional video prediction using deep networks in Atari games. In *Advances in Neural Information Processing Systems 28 (NIPS)*, pp. 2863–2871. 2015. URL <http://arxiv.org/abs/1507.08750>.
- J. K. O’Regan and A. Noë. A sensorimotor account of vision and visual consciousness. *Behavioral and brain sciences*, 24(05):939–973, 2001.
- P.-Y. Oudeyer, F. Kaplan, and V. V. Hafner. Intrinsic motivation systems for autonomous mental development. *Evolutionary Computation, IEEE Transactions on*, 11(2):265–286, 2007.
- V. Patraucean, A. Handa, and R. Cipolla. Spatio-temporal video autoencoder with differentiable memory. *CoRR*, abs/1511.06309, 2015. URL <http://arxiv.org/abs/1511.06309>.
- J. Pearl. *Causality*. Cambridge University Press, 2009.
- N. Srivastava, E. Mansimov, and R. Salakhutdinov. Unsupervised learning of video representations using LSTMs. In *Proceedings of the 32nd International Conference on Machine Learning (ICML)*, pp. 843–852, 2015.
- W. Sun, A. Venkatraman, B. Boots, and J. A. Bagnell. Learning to filter with predictive state inference machines. *CoRR*, abs/1512.08836, 2015. URL <http://arxiv.org/abs/1512.08836>.
- R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT Press, 1998.
- E. Talvitie. Model regularization for stable sample rollouts. In *Proceedings of the Thirtieth Conference Annual Conference on Uncertainty in Artificial Intelligence (UAI-14)*, pp. 780–789, 2014.
- N. Wahlström, T. B. Schön, and M. P. Deisenroth. From pixels to torques: Policy learning with deep dynamical models. *CoRR*, abs/1502.02251, 2015. URL <http://arxiv.org/abs/1502.02251>.
- M. Watter, J. Springenberg, J. Boedecker, and M. Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. In *Advances in Neural Information Processing Systems 28 (NIPS)*, pp. 2728–2736, 2015.
- R. J. Williams and D. Zipser. Gradient-based learning algorithms for recurrent networks and their computational complexity. *Bibliometrics*, pp. 433–486, 1995.
- B. Wymann, E. Espié, C. Guionneau, C. Dimitrakakis, R. Coulom, and A. Sumner. Torcs: The open racing car simulator, v1.3.5. 2013. URL <http://www.torcs.org>.
- B. Xu, N. Wang, T. Chen, and M. Li. Empirical evaluation of rectified activations in convolutional network. 2015.
- W. Zaremba, I. Sutskever, and O. Vinyals. Recurrent neural network regularization. *CoRR*, abs/1409.2329, 2014. URL <http://arxiv.org/abs/1409.2329>.



## A DATA, PREPROCESSING AND TRAINING ALGORITHM

When generating the data, each selected action was repeated for 4 time-steps and only the 4th frame was recorded for the analysis. The RGB images were preprocessed by subtracting mean pixel values (calculated separately for each color channel and over an initial set of 2048 frames only) and by dividing each pixel value by 255.

As stochastic gradient algorithm, we used centered RMSProp (Graves, 2013) with learning rate<sup>4</sup> 1e-5, epsilon 0.01, momentum 0.9, decay 0.95, and mini-batch size 16. The model was implemented in Torch, using the default initialization of the parameters. The state  $\mathbf{s}_1$  was initialized to zero.

## B PREDICTION-DEPENDENT SIMULATORS

As baseline for the single-step simulators we used the following state transition:

$$\begin{aligned} \text{Encoding: } \mathbf{z}_{t-1} &= \mathcal{C}(\mathbb{I}(\hat{\mathbf{x}}_{t-1}, \mathbf{x}_{t-1})), \\ \text{Action fusion: } \mathbf{v}_t &= \mathbf{W}^h \mathbf{h}_{t-1} \otimes \mathbf{W}^a \mathbf{a}_{t-1}, \\ \text{Gate update: } \mathbf{i}_t &= \sigma(\mathbf{W}^{iv} \mathbf{v}_t + \mathbf{W}^{iz} \mathbf{z}_{t-1}), \quad \mathbf{f}_t = \sigma(\mathbf{W}^{fv} \mathbf{v}_t + \mathbf{W}^{fs} \mathbf{z}_{t-1}), \\ \mathbf{o}_t &= \sigma(\mathbf{W}^{ov} \mathbf{v}_t + \mathbf{W}^{oz} \mathbf{z}_{t-1}), \\ \text{Cell update: } \mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv} \mathbf{v}_t + \mathbf{W}^{cz} \mathbf{z}_{t-1}), \\ \text{State update: } \mathbf{h}_t &= \mathbf{o}_t \otimes \tanh(\mathbf{c}_t), \end{aligned}$$

with vectors  $\mathbf{h}_{t-1}$  and  $\mathbf{v}_t$  of dimension 1024 and 2048 respectively.

### B.1 ATARI

We used a trained DQN agent (the scores are given in the table on the right) to generate training and test datasets consisting of 5,000,000 and 1,000,000 ( $210 \times 160$ ) RGB images respectively, with actions chosen according to an  $\epsilon = 0.2$ -greedy policy. Such a large number of training frames was necessary to prevent our simulators from strongly overfitting to the training data. This would be the case with, for example, one million training frames, as shown in Fig. 13 (the corresponding video can be seen at [MSPacman](#)). The ghosts are in frightened mode at time-step 1 (first image), and have returned to chase mode at time-step 63 (second image). The simulator is able to predict the exact time of return to the chase mode without sufficient history, which suggests that the sequence was memorized.

Game Name	DQN Score
Bowling	51.84
Breakout	396.25
Fishing Derby	19.30
Freeway	33.38
Ms Pacman	2963.31
Pong	20.88
Qbert	14,865.43
Riverraid	13,593.49
Seaquest	17,250.31
Space Invaders	2952.09

The encoding consisted of 4 convolutional layers with 64, 32, 32 and 32 filters, of size  $8 \times 8$ ,  $6 \times 6$ ,  $6 \times 6$ , and  $4 \times 4$ , stride 2, and padding 0, 1, 1, 0 and 1, 1, 1, 0 for the height and width respectively. Every layer was followed by a randomized rectified linear function (RReLU) (Xu et al., 2015) with parameters  $l = 1/8$ ,  $u = 1/3$ . The output tensor of the convolutional layers of dimension  $32 \times 11 \times 8$  was then flattened into the vector  $\mathbf{z}_t$  of dimension 2816. The decoding consisted of one fully-connected layer with 2816 hidden units followed by 4 full convolutional layers with the inverse symmetric structure of the encoding transformation: 32, 32, 32 and 64 filters, of size  $4 \times 4$ ,  $6 \times 6$ ,  $6 \times 6$ , and  $8 \times 8$ , stride 2, and padding 0, 1, 1, 0 and 0, 1, 1, 1. Each full convolutional layer (except the last one) was followed by a RReLU.

In Fig. 14, we show one example of successful prediction at time-steps 100 and 200 for each game.

#### B.1.1 SHORT-TERM VERSUS LONG-TERM ACCURACY

In Figures 15-19, we show the prediction error obtained with the training schemes described in Sec. 3.1 for all games. Below we discuss the main findings for each game.

<sup>4</sup>We found that using a higher learning rate value of  $2e-5$  would generally increase convergence speed but cause major instability issues, suggesting that gradient clipping would need to be used.

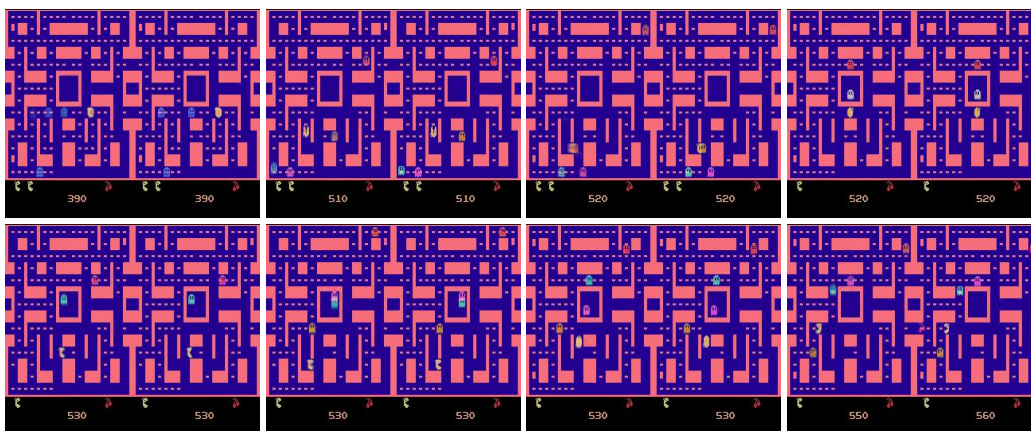


Figure 13: Prediction that demonstrates overfitting of the model when trained on one million frames.

**Bowling.** Bowling is one of the easiest games to model. A simulator trained using only observation-dependent transitions gives quite accurate predictions. However, using only prediction-dependent transitions reduces the error in updating the score and predicting the ball direction.

**Breakout.** Breakout is a difficult game to model. A simulator trained with only prediction-dependent transitions predicts the paddle movement very accurately but almost always fails to represent the ball. A simulator trained with only observation-dependent transitions struggles much less to represent the ball but does not predict the paddle and ball positions as accurately, and the ball also often disappears after hitting the paddle. Interestingly, the long-term prediction error (bottom-right of Fig. 15(b)) for the 100%PDT training scheme is the lowest, as when not representing the ball the predicted frames look closer to the real frames than when representing the ball incorrectly. A big improvement in the model ability to represent the ball could be obtained by pre-processing the frames with max-pooling as done for DQN, as this increases the ball size. We believe that a more sophisticated convolutional structure would be even more effective, but did not succeed in discovering such a structure.

**Fishing Derby.** In Fishing Derby, long-term accuracy is disastrous with the 0%PDT training scheme and good with the 100%PDT training scheme. Short-term accuracy is better with schemes using more observation-dependent transitions than in the 100% or 0%-100%PDT training schemes, especially at low numbers of parameter updates.

**Freeway.** With Bowling, Freeway is one of the easiest games to model, but more parameter updates are required for convergence than for Bowling. The 0%PDT training scheme gives good accuracy, although sometimes the chicken disappears or its position is incorrectly predicted – this happens extremely rarely with the 100%PDT training scheme. In both schemes, the score is often wrongly updated in the warning phase.

**Ms Pacman.** Ms Pacman is a very difficult game to model and accurate prediction can only be obtained for a few time-steps into the future. The movement of the ghosts, especially when in frightened mode, is regulated by the position of Ms Pacman according to complex rules. Furthermore, the DQN  $\epsilon = 0.2$ -greedy policy does not enable the agent to explore certain regions of the state space. As a result, the simulator can predict well the movement of Ms Pacman, but fails to predict long-term the movement of the ghosts when in frightened mode or when in chase mode later in the episodes.

**Pong.** With the 0%PDT training scheme, the model often incorrectly predicts the direction of the ball when hit by the agent or by the opponent. Quite rarely, the ball disappears when hit by the agent. With the 100%PDT training scheme, the direction the ball is much more accurately predicted, but the ball more often disappears when hit by the agent, and the ball and paddles are generally less sharp.

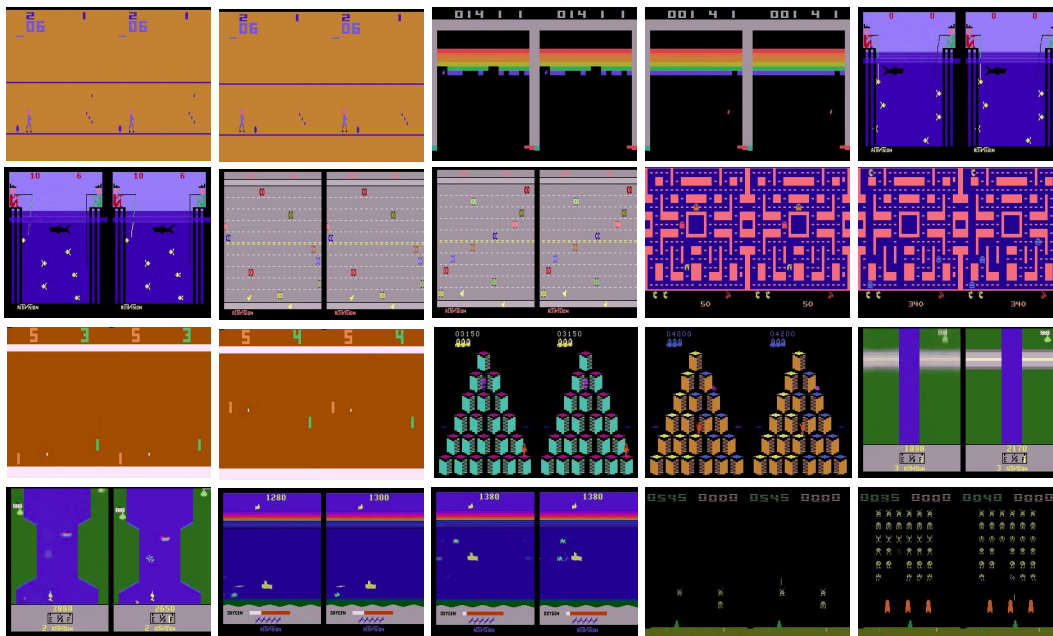


Figure 14: One example of 200 time-step ahead prediction for each of the 10 Atari games. Displayed are predicted (left) and real (right) frames at time-steps 100 and 200.

**Qbert.** Qbert is a game for which the 0%PDT training scheme is unable to predict accurately beyond very short-term, as after a few frames only the background is predicted. The more prediction-dependent transitions are used, the less sharply the agent and the moving objects are represented.

**Riverraid.** In Riverraid, prediction with the 0%PDT training scheme is very poor, as this scheme causes no generation of new objects or background. With all schemes, the model fails to predict the frames that follow a jet loss – that’s why the prediction error increases sharply after around time-step 13 in Fig. 18(b). The long-term prediction error is lower with the 100%PDT training scheme, as with this scheme the simulator is more accurate before, and sometimes after, a jet loss. The problem of incorrect prediction after a jet loss disappears when using BBTT(15,2) with prediction-dependent transitions.

**Seaquest.** In Seaquest, with the 0%PDT training scheme, the existing fish disappears after a few time-steps and no new fish ever appears from the sides of the frame. The higher the number of prediction-dependent transitions the less sharply the fish is represented, but the more accurately its dynamics and appearance from the sides of the frame can be predicted.

**Space Invaders** Space Invaders is a very difficult game to model and accurate prediction can only be obtained for a few time-steps into the future. The 0%PDT training scheme is unable to predict accurately beyond very short-term. The 100%PDT training scheme struggles to represent the bullets.

In Figs. 20-24 we show the effect of using different prediction lengths  $T \leq 20$  with the training schemes 0%PDT, 67%PDT, and 100%PDT for all games.

In Figs. 25-29 we show the effect of using different prediction lengths  $T > 20$  through truncated backpropagation with the training schemes 0%PDT, 33%PDT, and 100%PDT for all games.

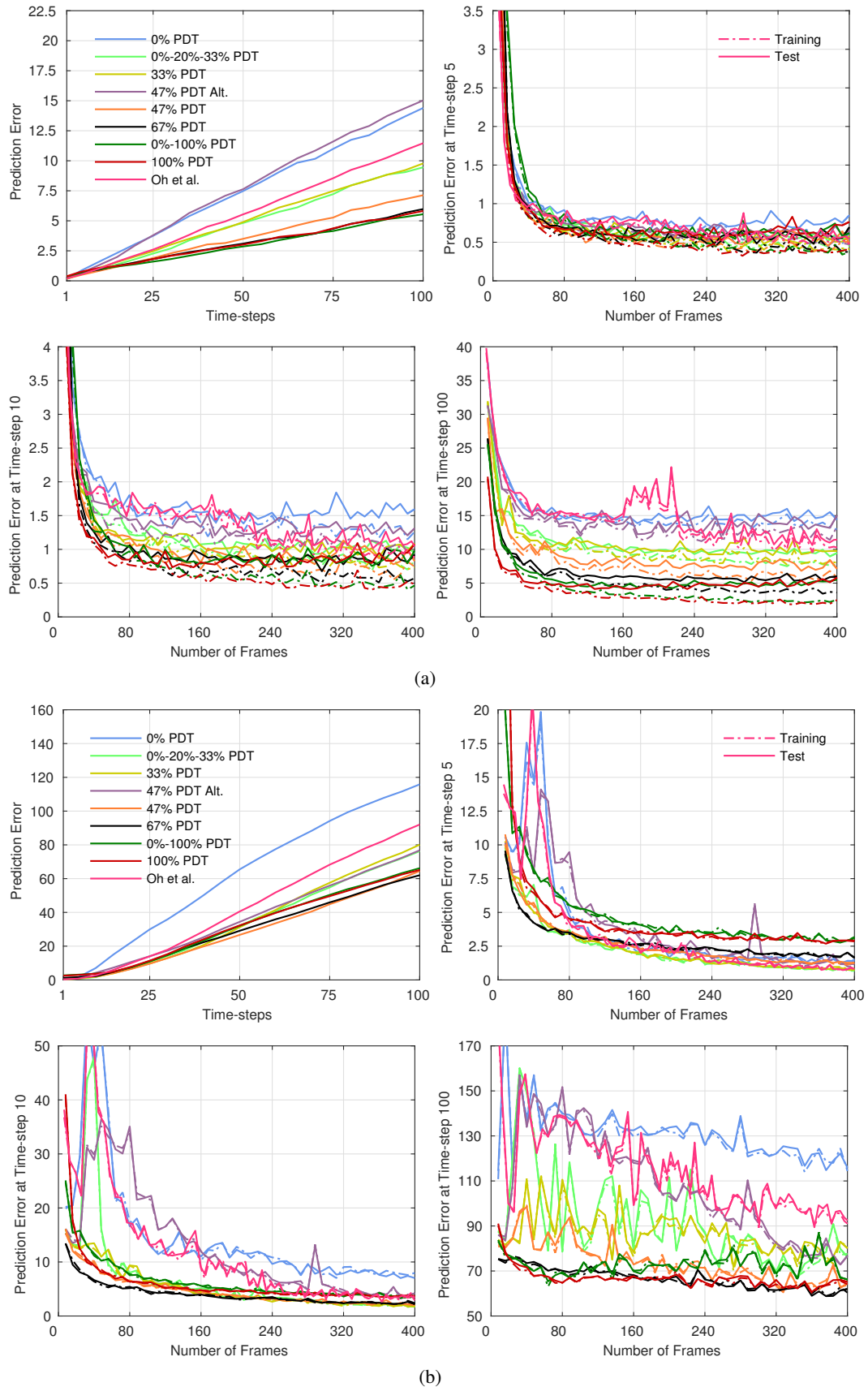


Figure 15: Prediction error (average over 10,000 sequences) for different training schemes on (a) Bowling and (b) Breakout. Number of frames is in millions.

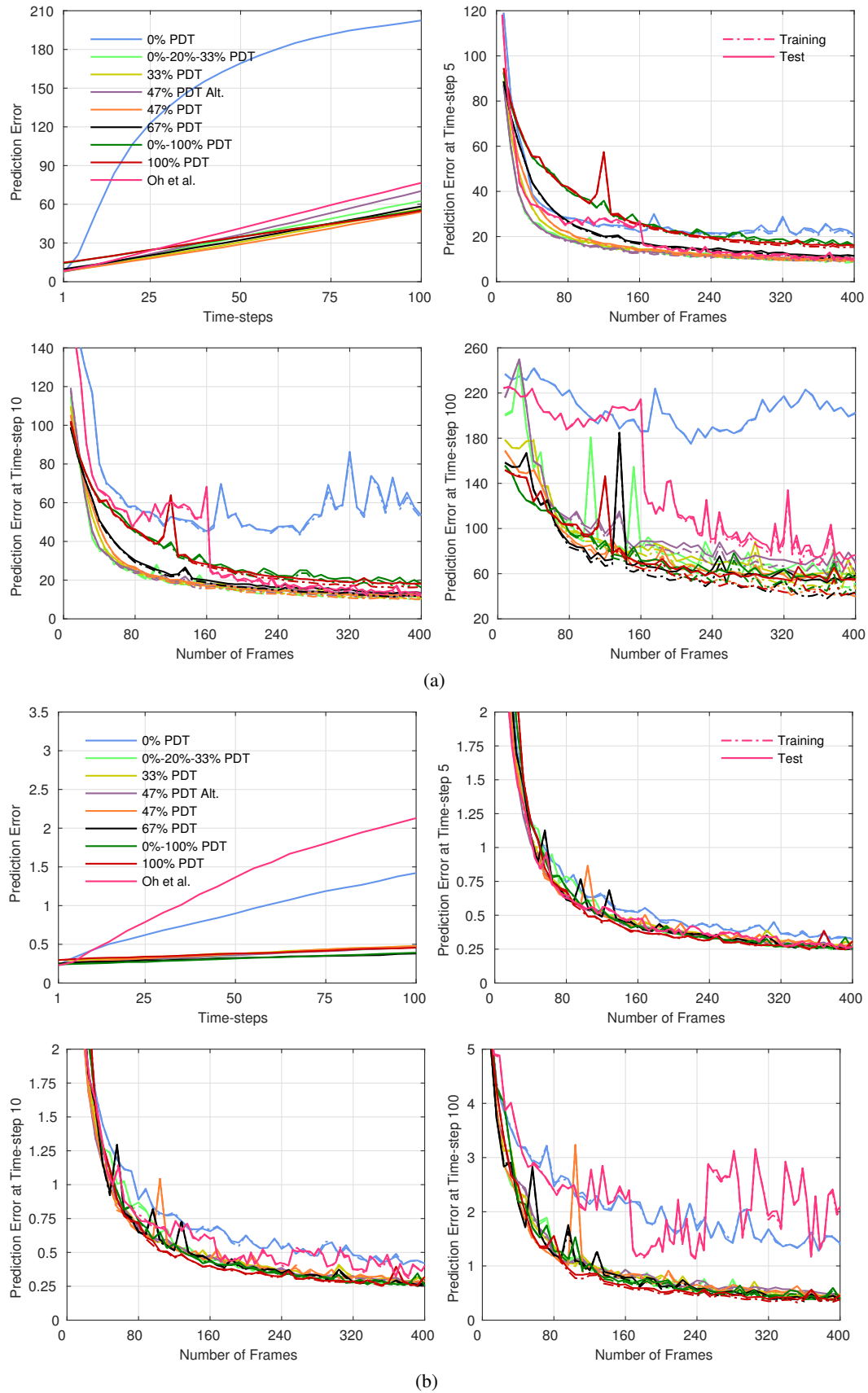


Figure 16: Prediction error for different training schemes on (a) Fishing Derby and (b) Freeway.

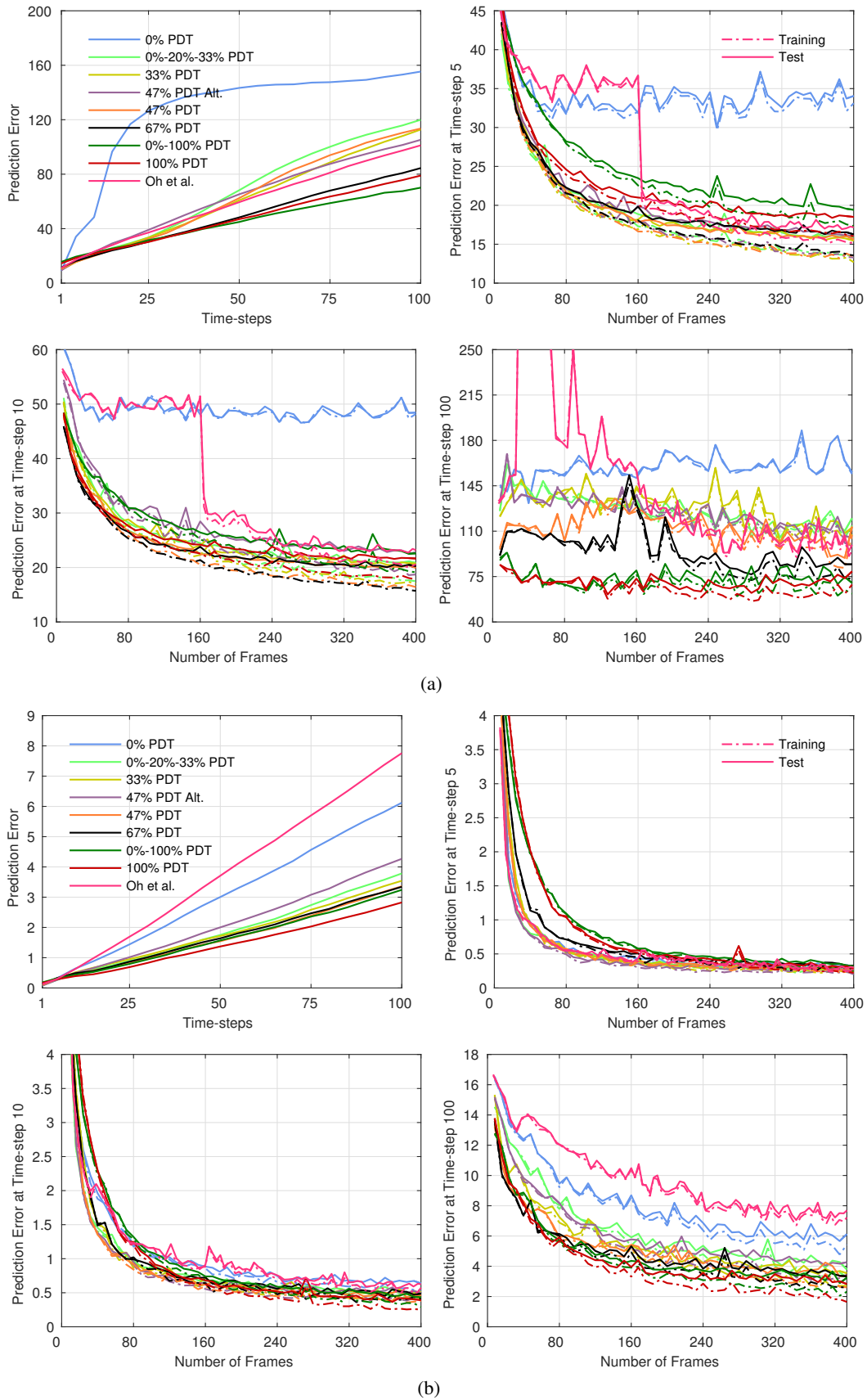
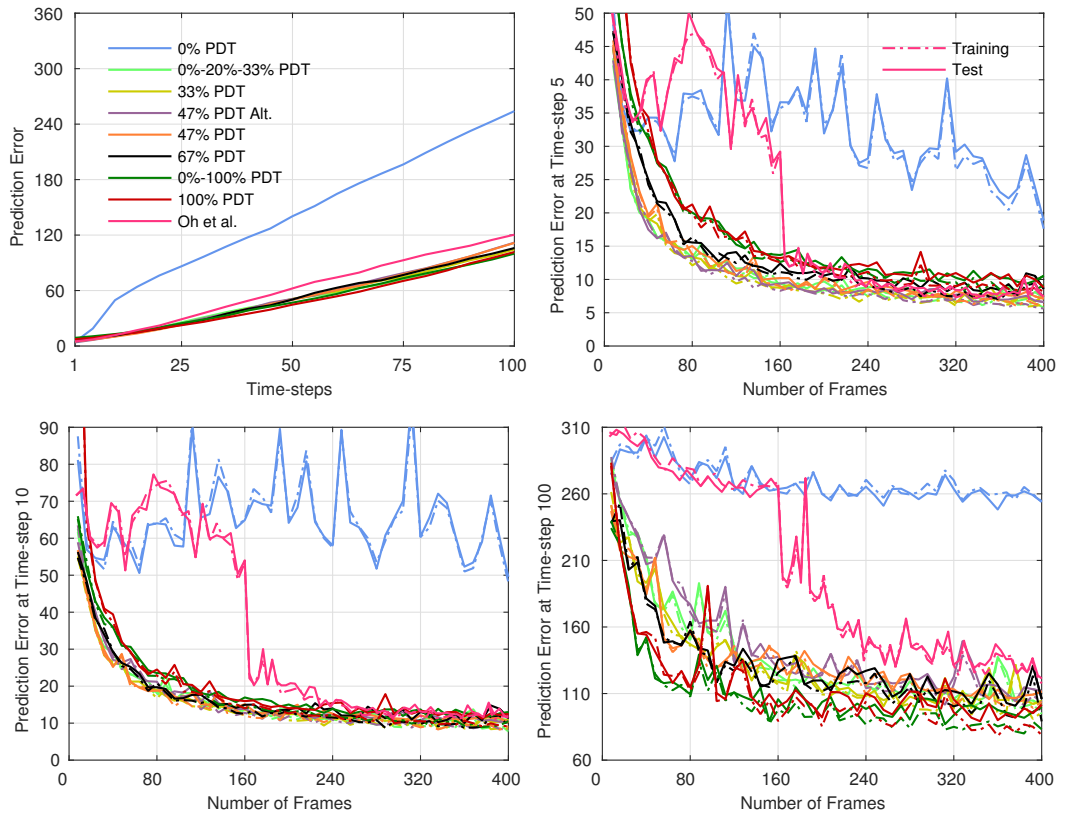
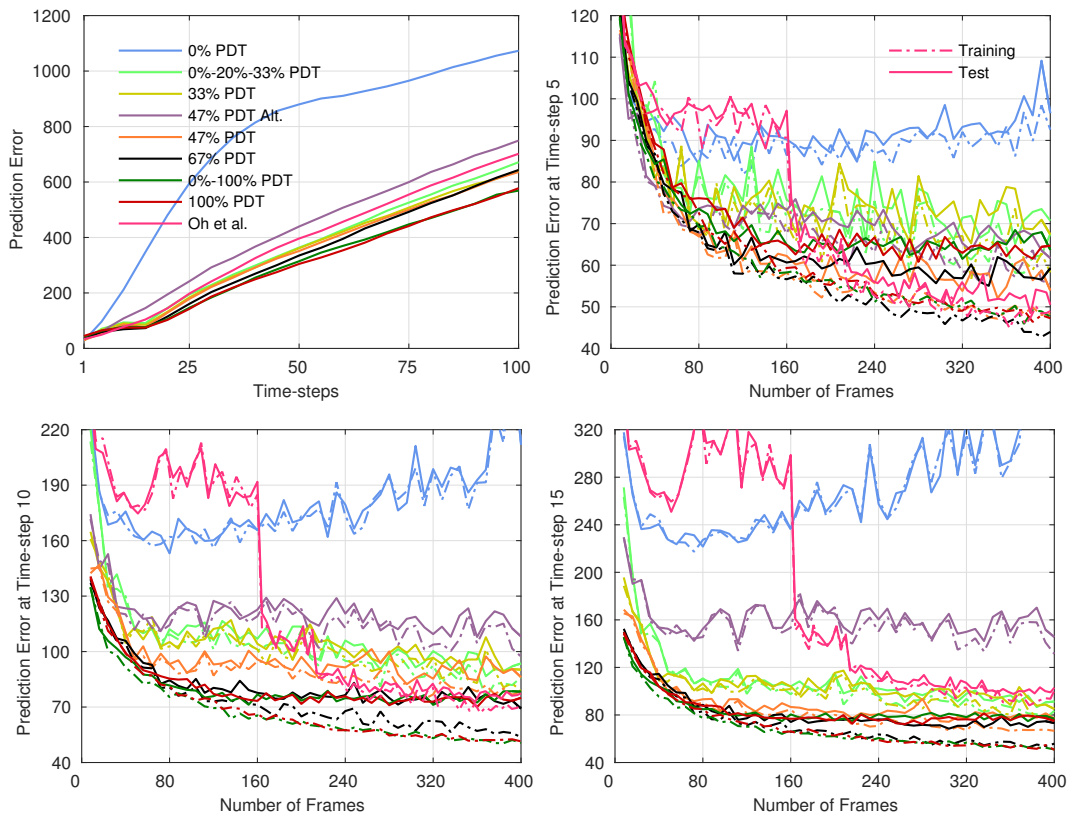


Figure 17: Prediction error for different training schemes on (a) Ms Pacman and (b) Pong.



(a)



(b)

Figure 18: Prediction error for different training schemes on (a) Qbert and (b) Riverraid.

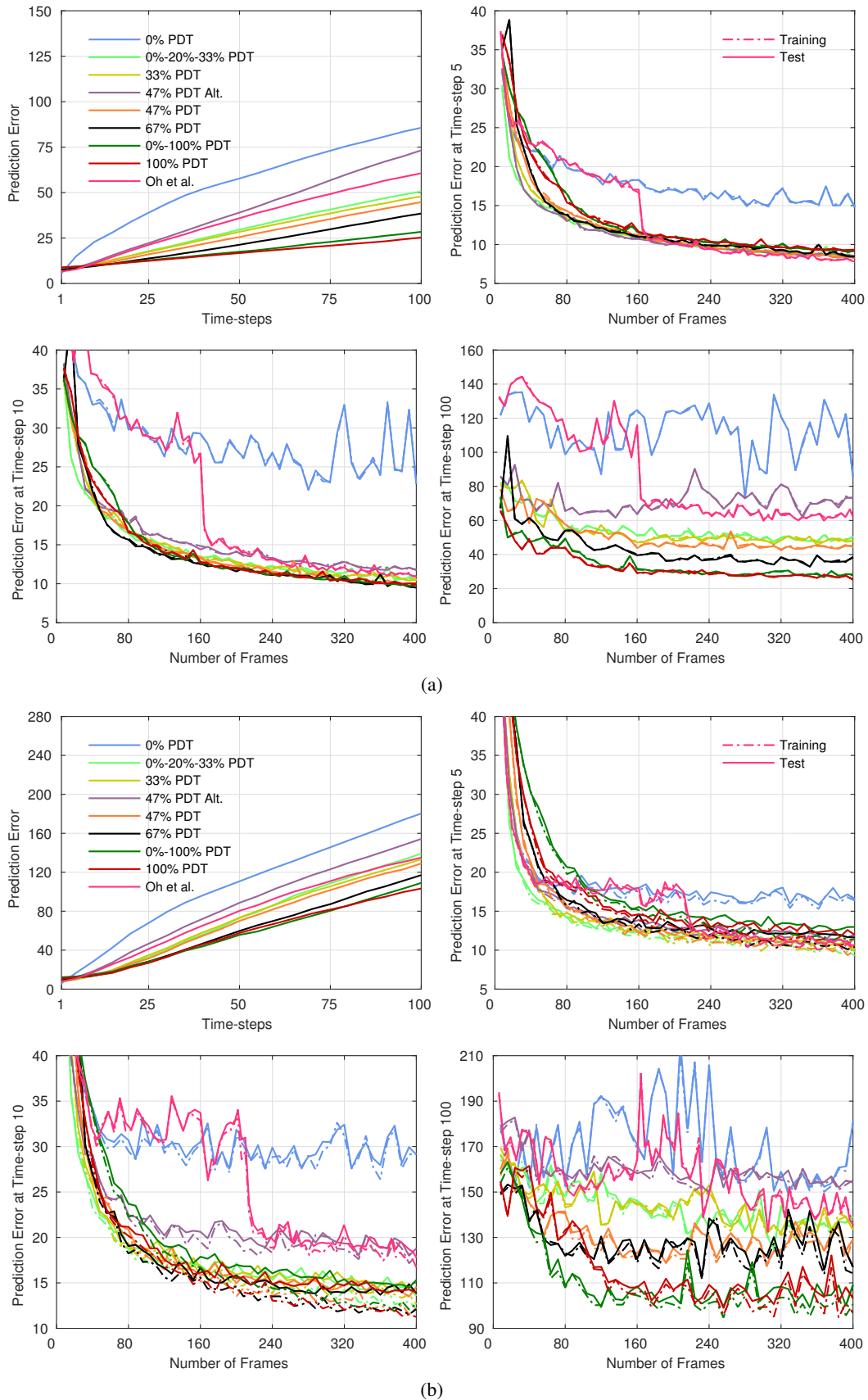


Figure 19: Prediction error for different training schemes on (a) Seaquest and (b) Space Invaders.



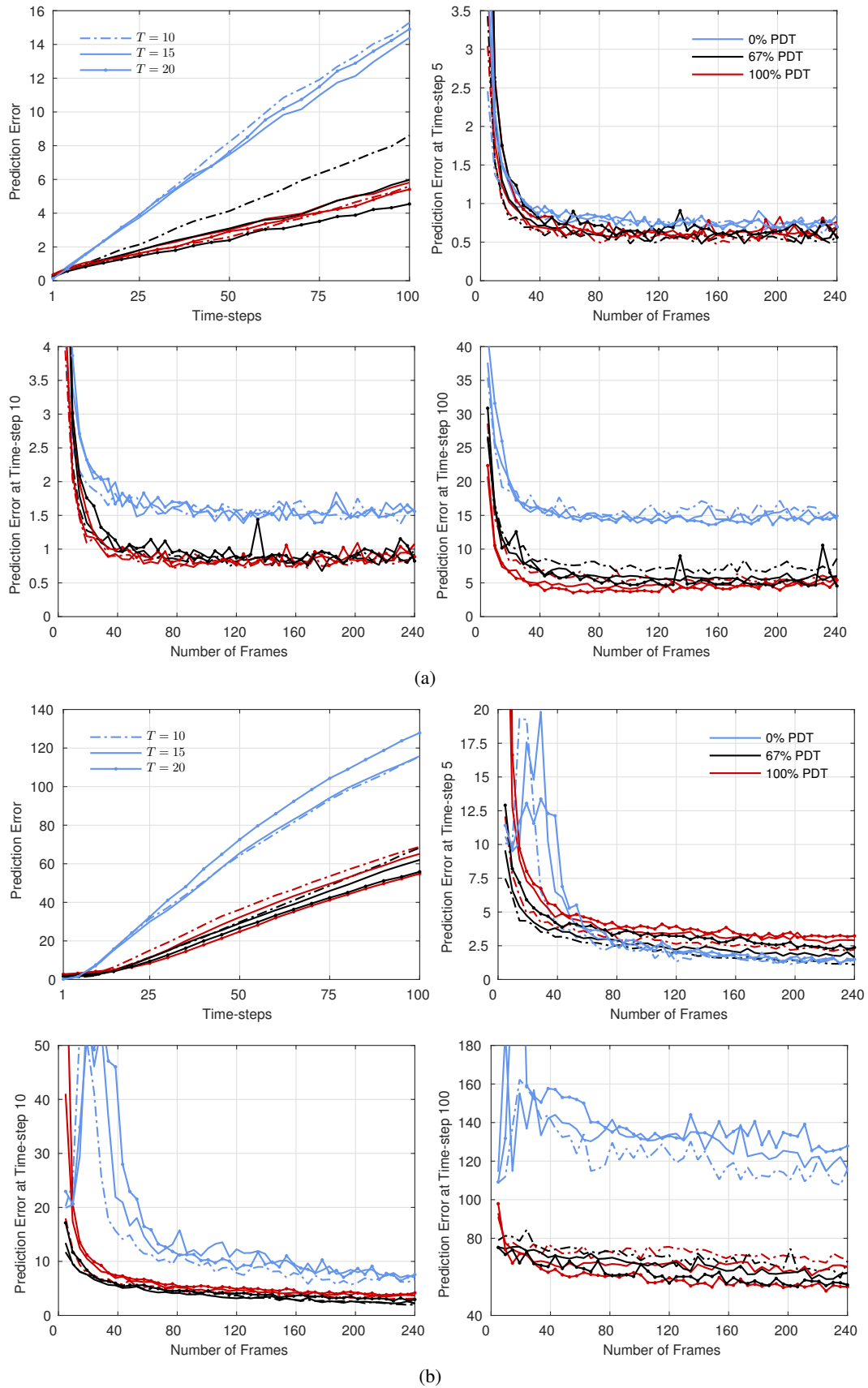


Figure 20: Prediction error (average over 10,000 sequences) for different prediction lengths  $T \leq 20$  on (a) Bowling and (b) Breakout. Number of frames is in millions and excludes warm-up frames.

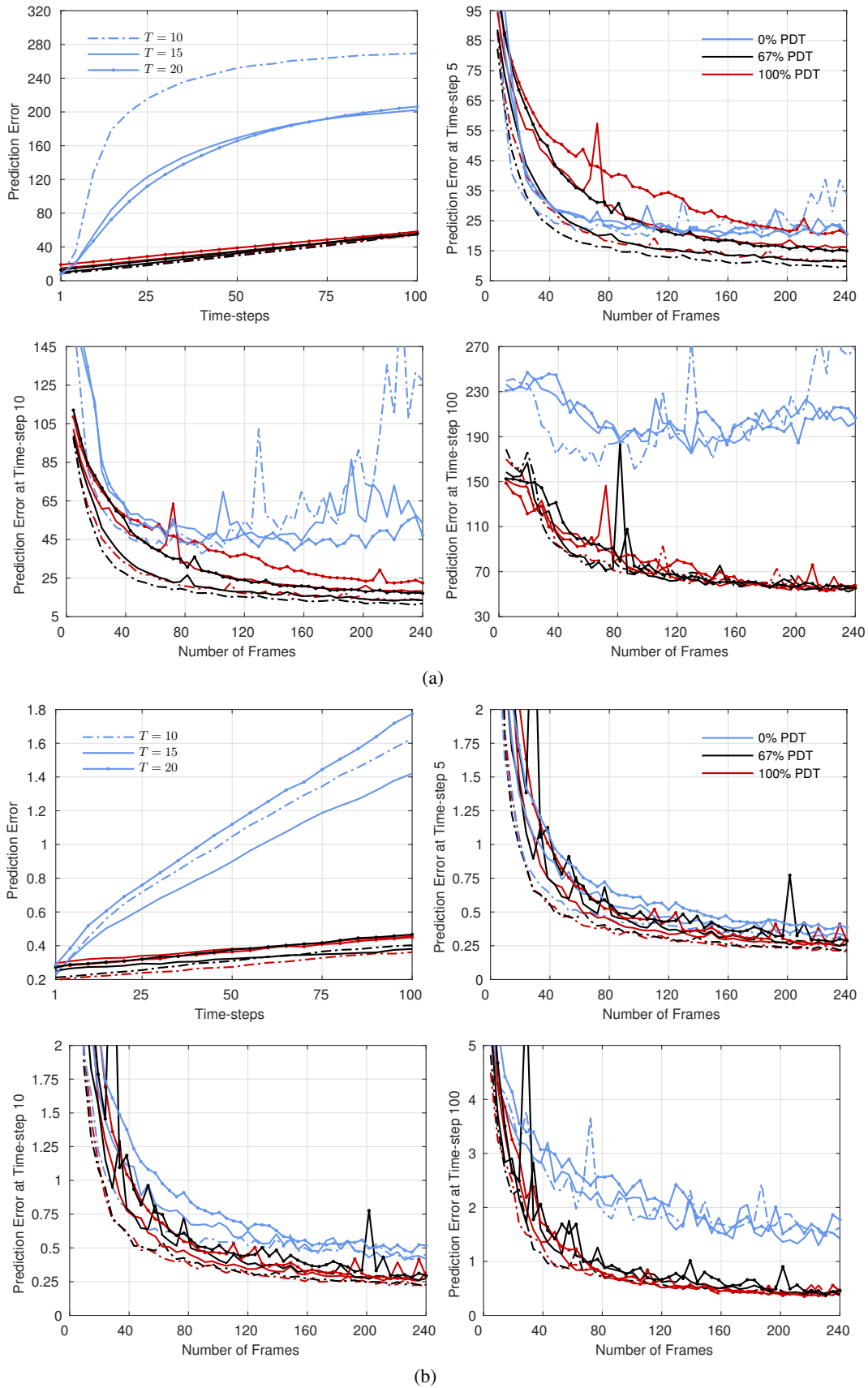


Figure 21: Prediction error for different prediction lengths  $T \leq 20$  on (a) Fishing Derby and (b) Freeway.

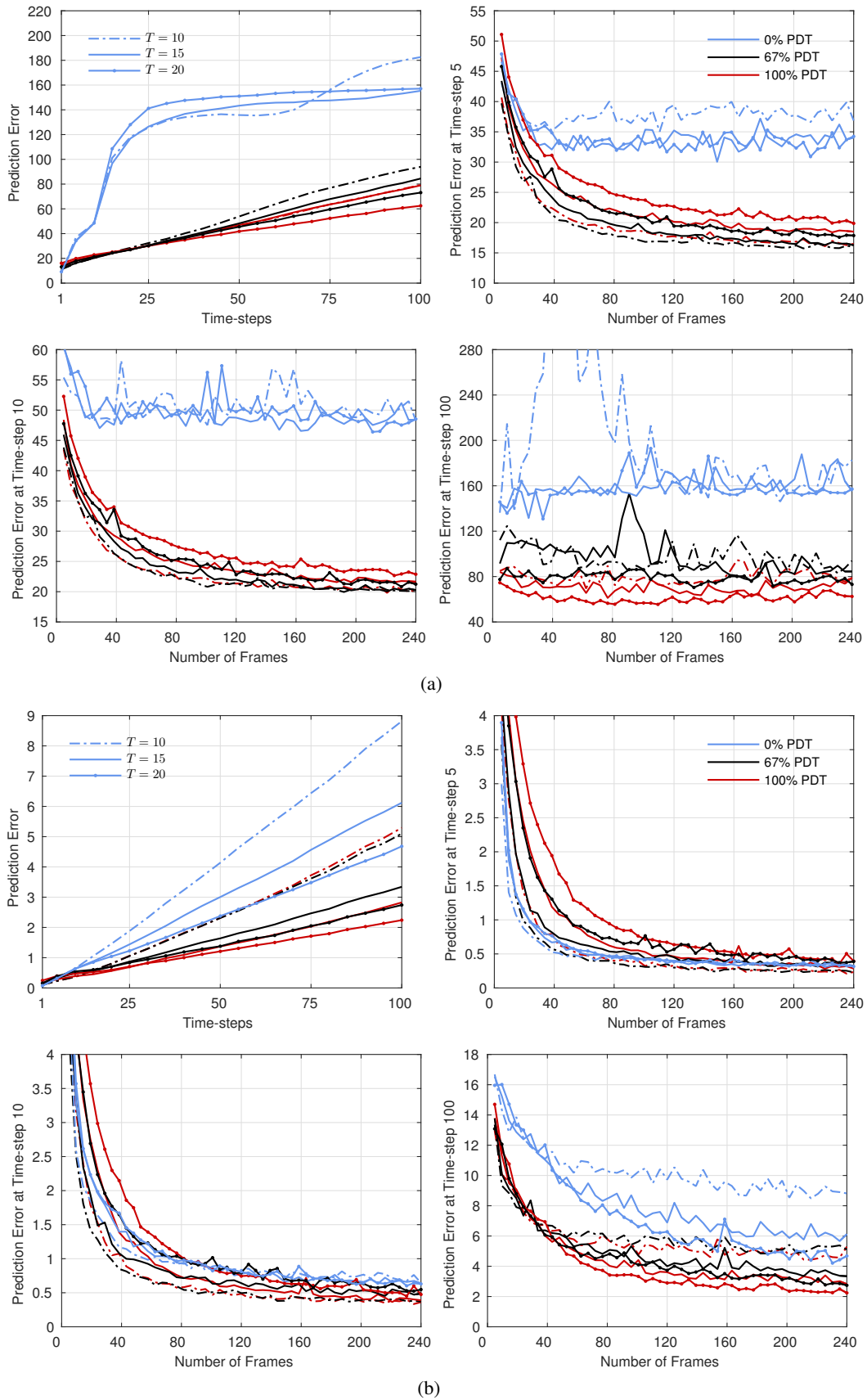
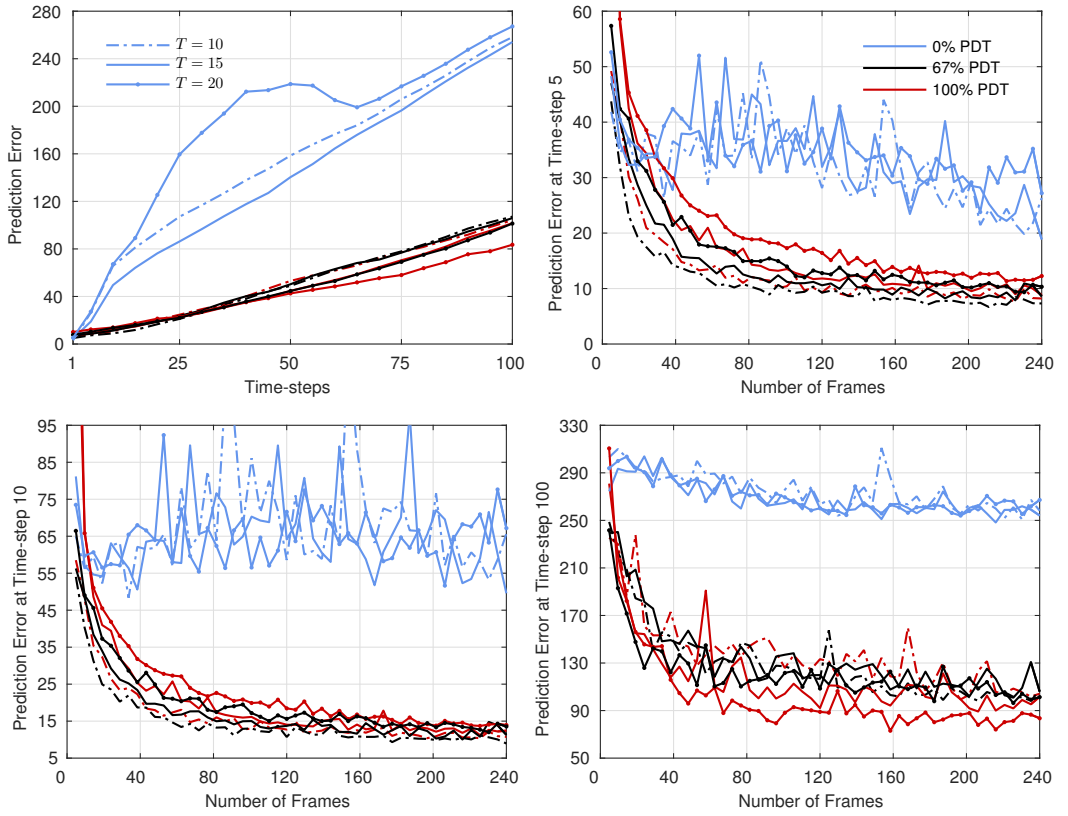
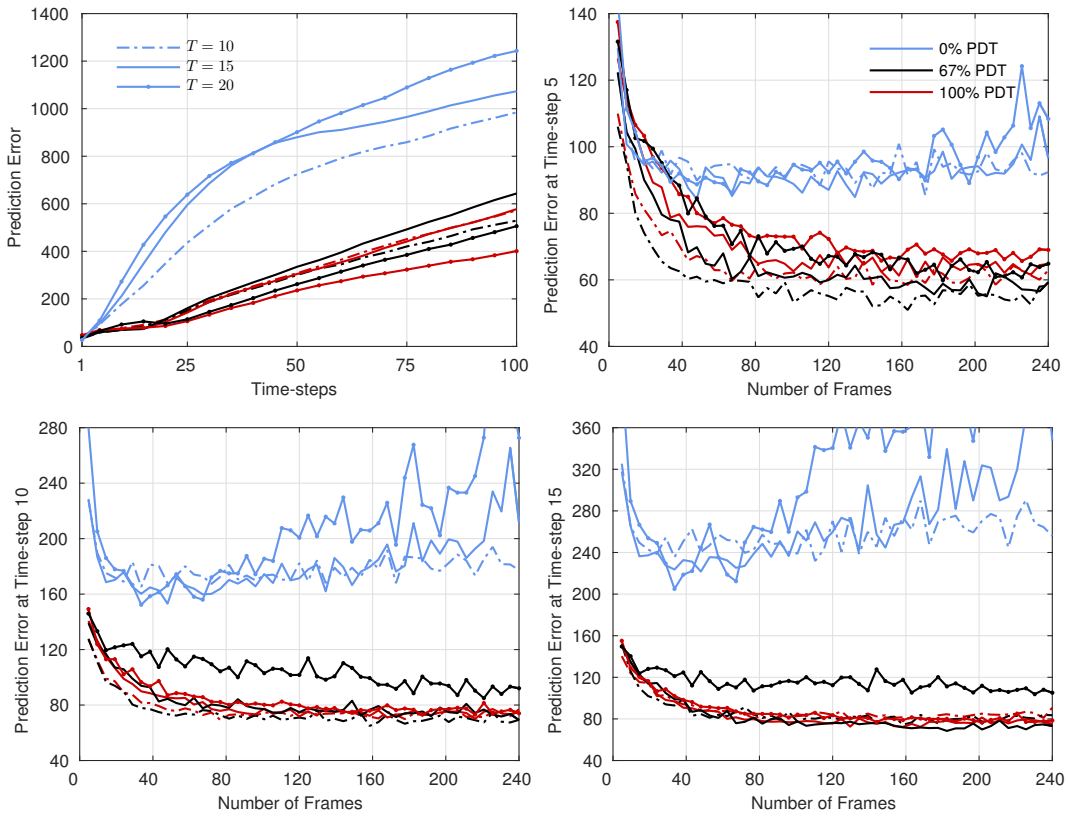


Figure 22: Prediction error for different prediction lengths  $T \leq 20$  on (a) Ms Pacman and (b) Pong.



(a)



(b)

Figure 23: Prediction error for different prediction lengths  $T \leq 20$  on (a) Qbert and (b) Riverraid.

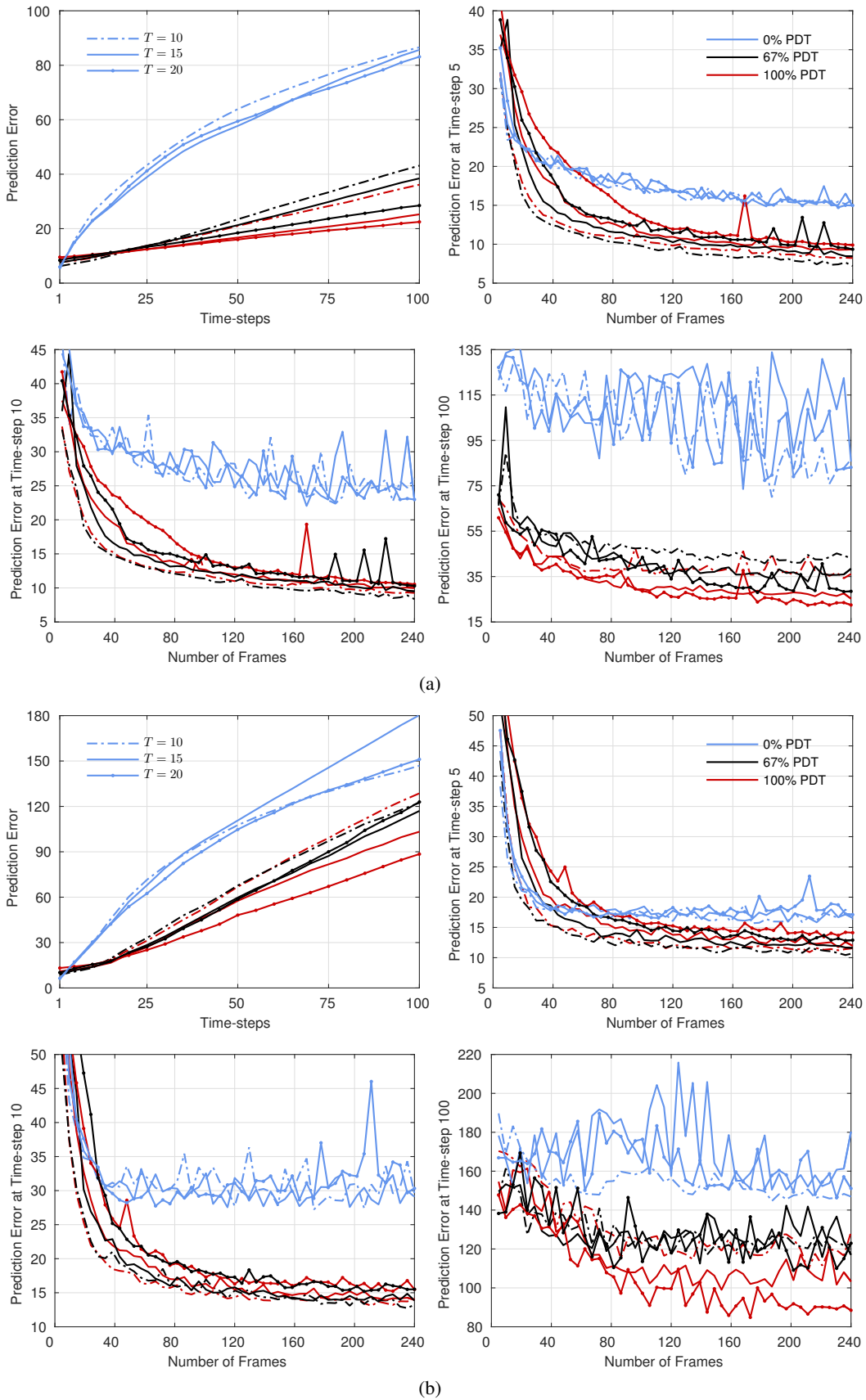


Figure 24: Prediction error for different prediction lengths  $T \leq 20$  on (a) Seaquest and (b) Space Invaders.

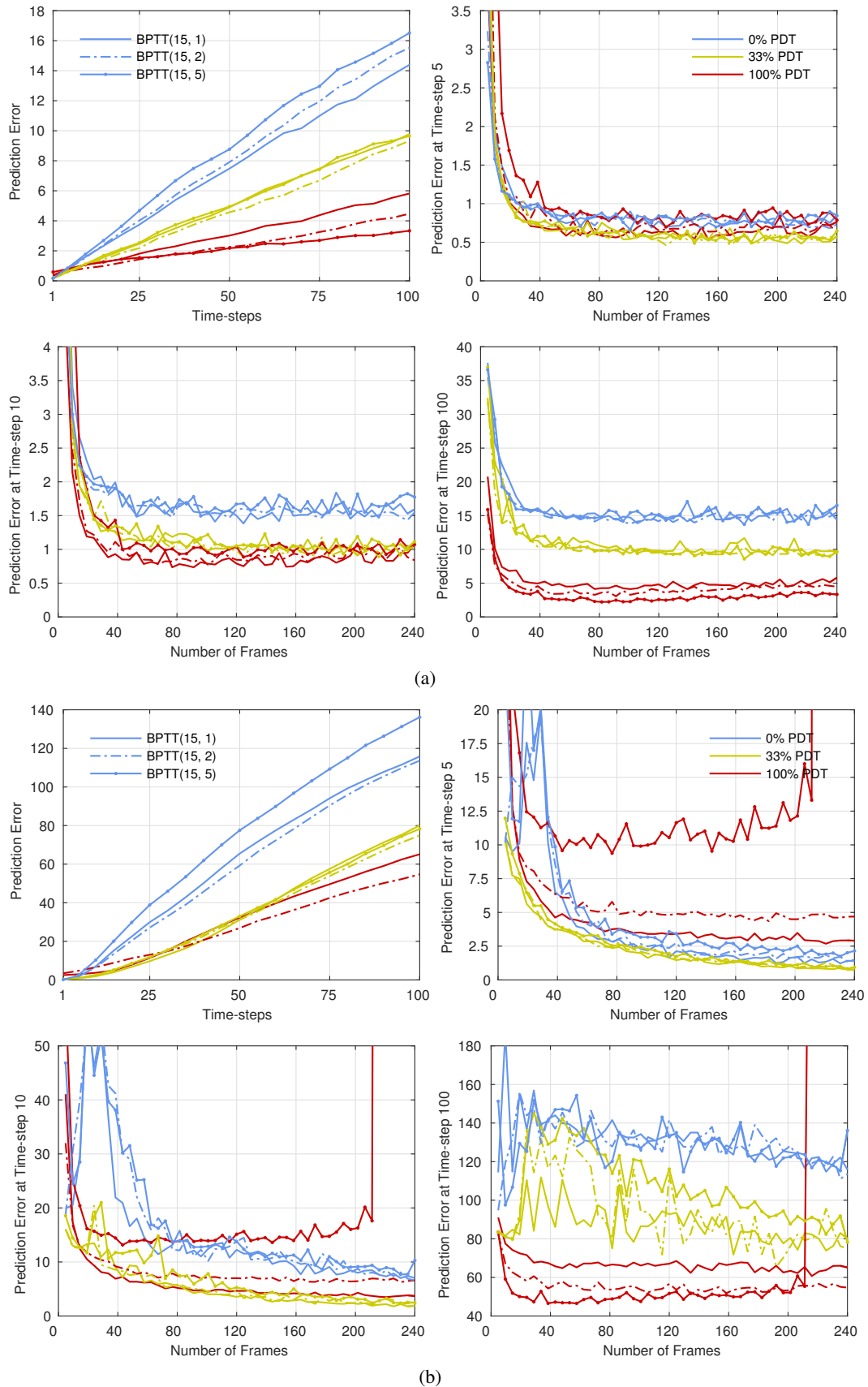


Figure 25: Prediction error (average over 10,000 sequences) for different prediction lengths through truncated BPTT on (a) Bowling and (b) Breakout. Number of frames is in millions and excludes warm-up frames.

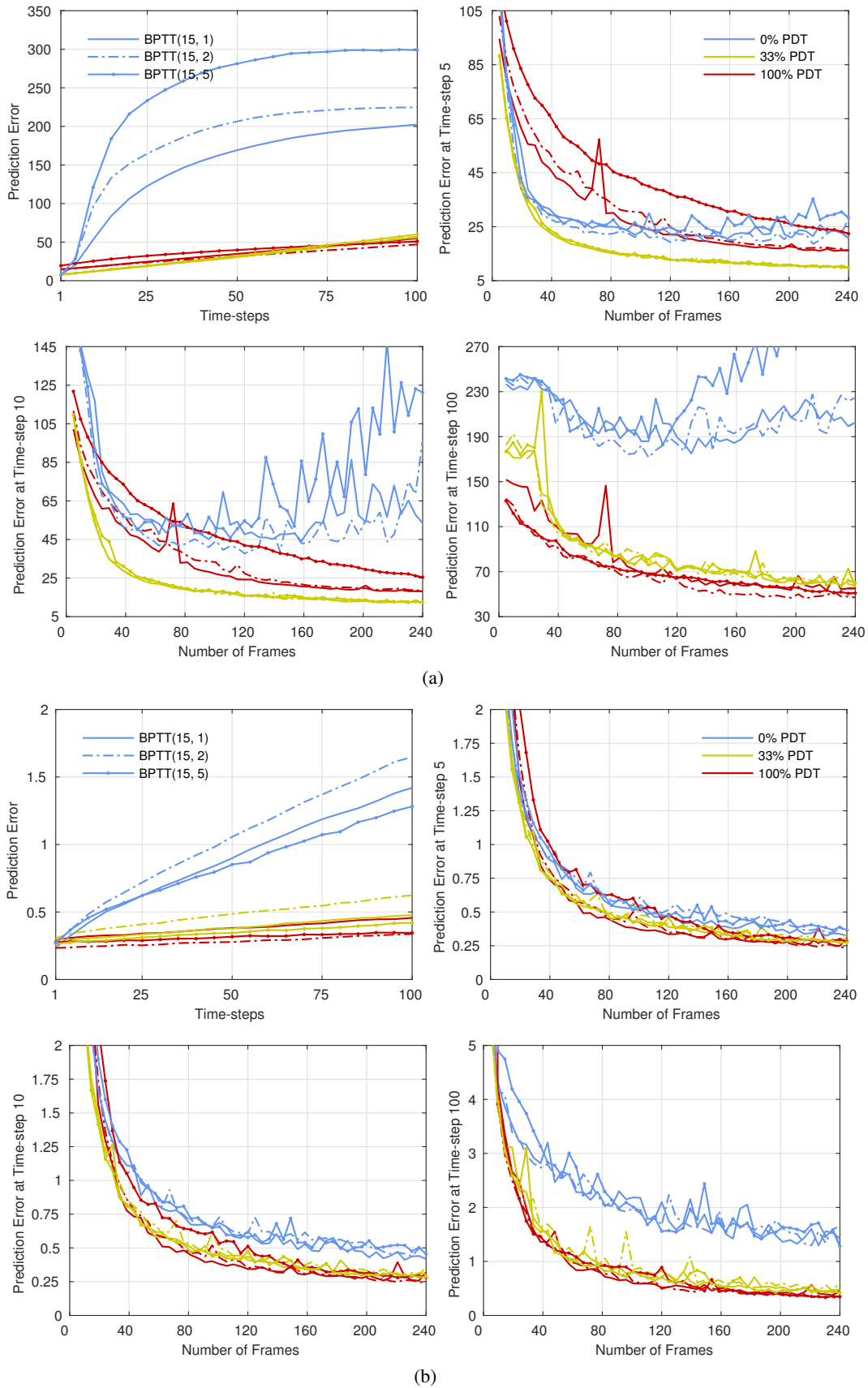


Figure 26: Prediction error for different prediction lengths through truncated BPTT on (a) Fishing Derby and (b) Freeway.

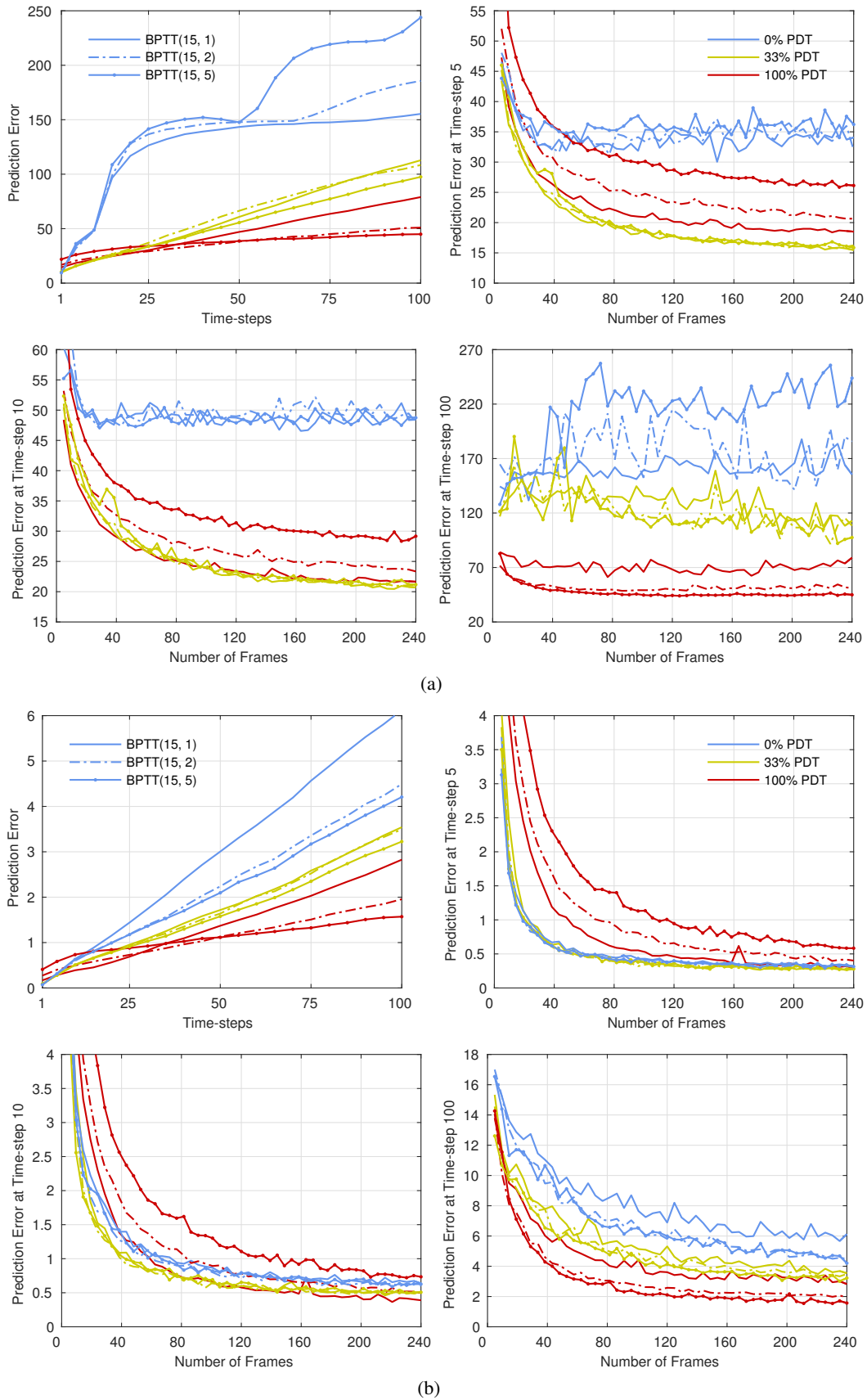


Figure 27: Prediction error for different prediction lengths through truncated BPTT on (a) Ms Pacman and (b) Pong.



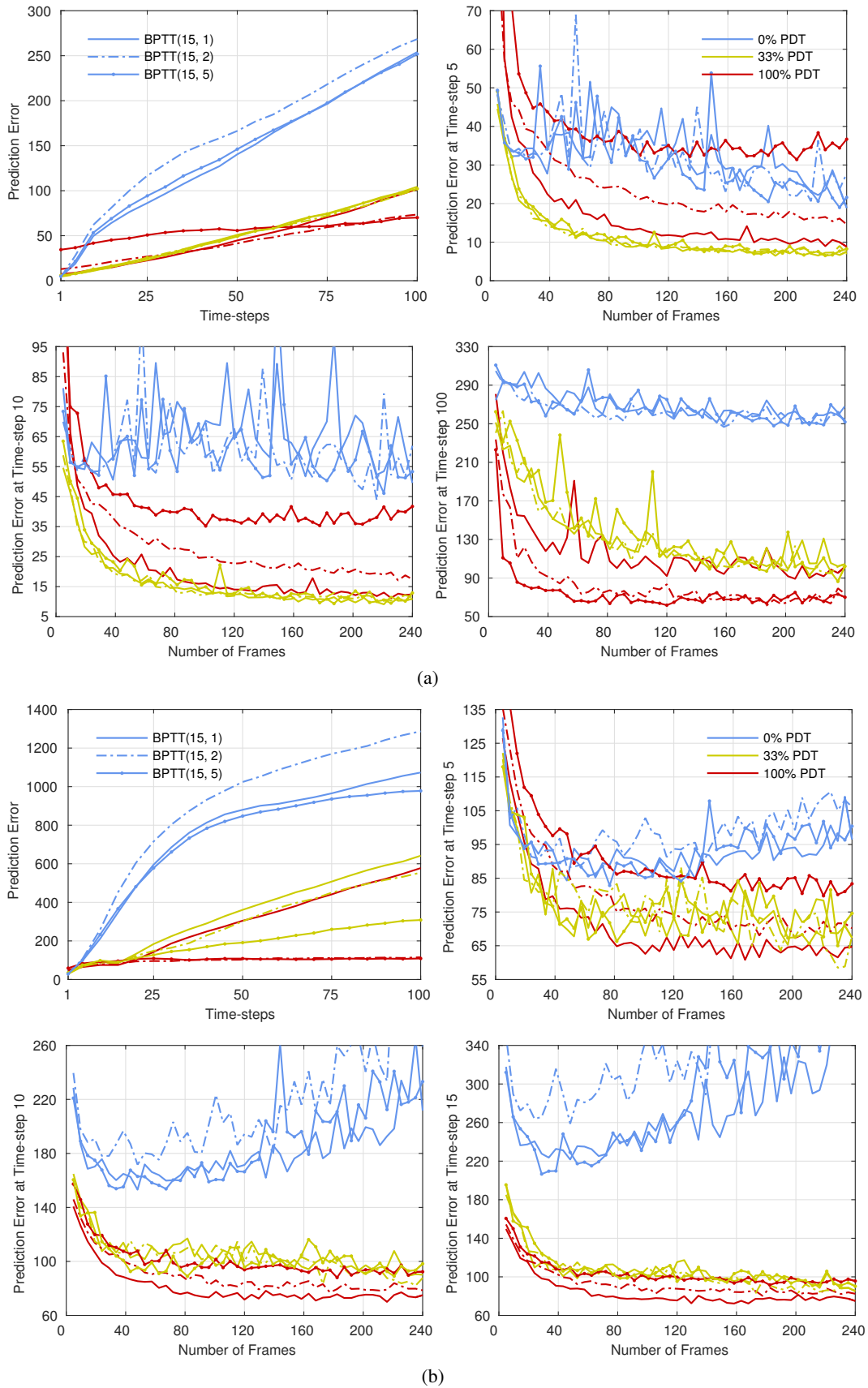


Figure 28: Prediction error for different prediction lengths through truncated BPTT on (a) Qbert and (b) Riverraid.

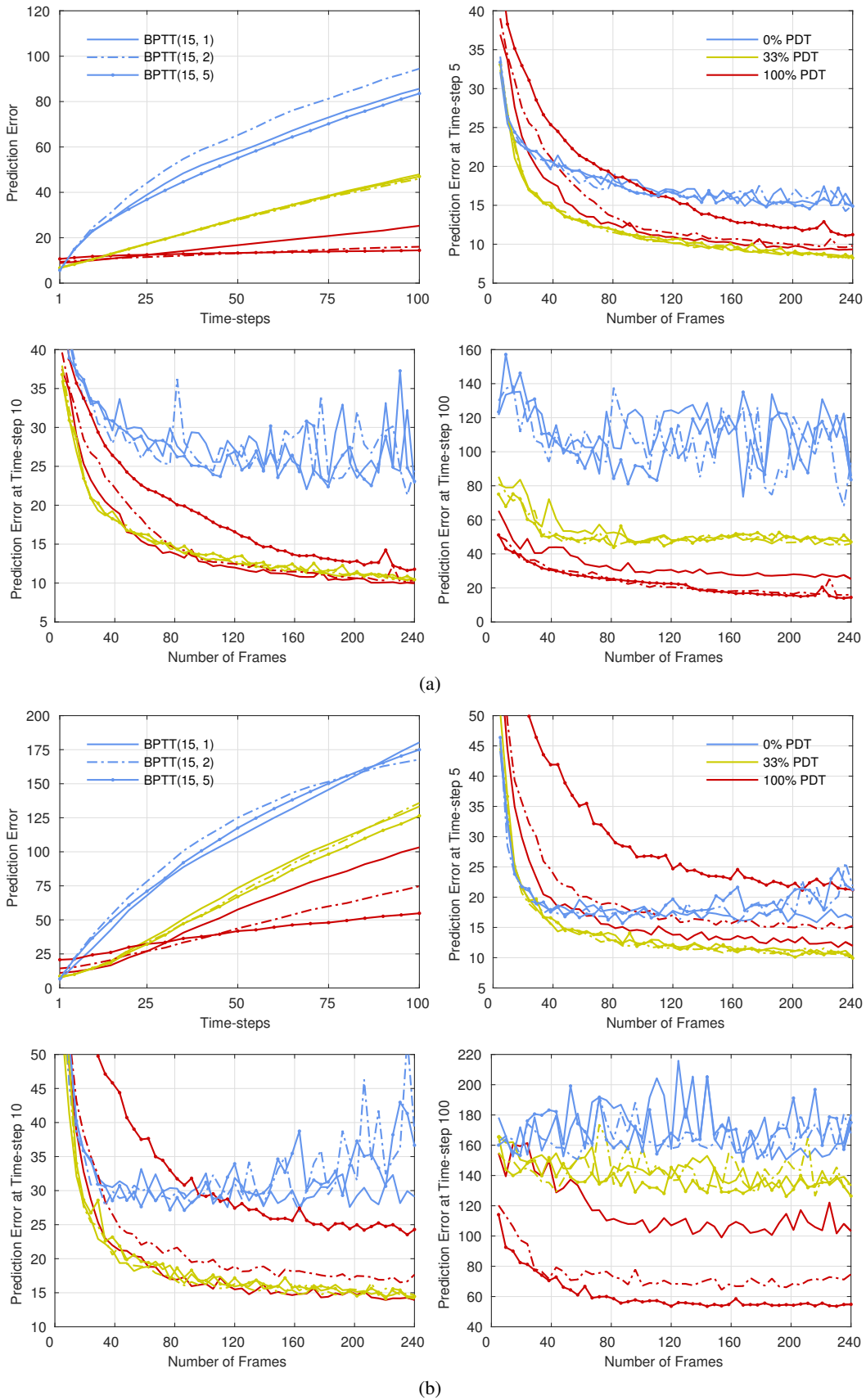


Figure 29: Prediction error for different prediction lengths through truncated BPTT on (a) Seaquest and (b) Space Invaders.

### B.1.2 DIFFERENT ACTION-DEPENDENT STATE TRANSITIONS

In this section we compare the baseline state transition

$$\begin{aligned}
&\text{Encoding: } \mathbf{z}_{t-1} = \mathcal{C}(\mathbb{I}(\hat{\mathbf{x}}_{t-1}, \mathbf{x}_{t-1})), \\
&\text{Action fusion: } \mathbf{v}_t = \mathbf{W}^h \mathbf{h}_{t-1} \otimes \mathbf{W}^a \mathbf{a}_{t-1}, \\
&\text{Gate update: } \mathbf{i}_t = \sigma(\mathbf{W}^{iv} \mathbf{v}_t + \mathbf{W}^{iz} \mathbf{z}_{t-1}), \quad \mathbf{f}_t = \sigma(\mathbf{W}^{fv} \mathbf{v}_t + \mathbf{W}^{fz} \mathbf{z}_{t-1}), \\
&\quad \mathbf{o}_t = \sigma(\mathbf{W}^{ov} \mathbf{v}_t + \mathbf{W}^{oz} \mathbf{z}_{t-1}), \\
&\text{Cell update: } \mathbf{c}_t = \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv} \mathbf{v}_t + \mathbf{W}^{cz} \mathbf{z}_{t-1}), \\
&\text{State update: } \mathbf{h}_t = \mathbf{o}_t \otimes \tanh(\mathbf{c}_t),
\end{aligned}$$

where the vectors  $\mathbf{h}_{t-1}$  and  $\mathbf{v}_t$  have dimension 1024 and 2048 respectively (this model has around 25 millions (25M) parameters), with alternatives using unconstrained or convolutional transformations, for prediction length  $T = 15$  and the 0%-100%PDT training scheme.

More specifically, in Figs. 30-34 we compare the baseline transition with the following alternatives:

- **Base2816:** The vectors  $\mathbf{h}_{t-1}$  and  $\mathbf{v}_t$  have the same dimension as  $\mathbf{z}_{t-1}$ , namely 2816. This model has around 80M parameters.
- **$\mathbf{i}_t^z$  and  $\mathbf{i}_t^z$ 2816:** Have a separate gating for  $\mathbf{z}_{t-1}$  in the cell update, *i.e.*

$$\mathbf{c}_t = \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv} \mathbf{v}_t) + \mathbf{i}_t^z \otimes \tanh(\mathbf{W}^{cz} \mathbf{z}_{t-1}).$$

This model has around 30 million parameters. We also considered removing the linear projection of  $\mathbf{z}_{t-1}$ , *i.e.*

$$\mathbf{c}_t = \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv} \mathbf{v}_t) + \mathbf{i}_t^z \otimes \tanh(\mathbf{z}_{t-1}),$$

without RReLU after the last convolution and with vectors  $\mathbf{h}_{t-1}$  and  $\mathbf{v}_t$  of dimensionality 2816. This model has around 88M parameters.

- **$\neg \mathbf{z}_{t-1}$  and  $\neg \mathbf{z}_{t-1} - \mathbf{i}_t^z$ 2816:** Remove  $\mathbf{z}_{t-1}$  in the gate updates, *i.e.*

$$\mathbf{i}_t = \sigma(\mathbf{W}^{iv} \mathbf{v}_t), \quad \mathbf{f}_t = \sigma(\mathbf{W}^{fv} \mathbf{v}_t), \quad \mathbf{o}_t = \sigma(\mathbf{W}^{ov} \mathbf{v}_t),$$

with one of the following cell updates

$$\begin{aligned}
\mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv} \mathbf{v}_t + \mathbf{W}^{cz} \mathbf{z}_{t-1}), \quad 17\text{M parameters}, \\
\mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv} \mathbf{v}_t) + \mathbf{i}_t^z \otimes \tanh(\mathbf{z}_{t-1}), \quad 56\text{M parameters}.
\end{aligned}$$

- **$\mathbf{h}_{t-1}$ ,  $\mathbf{h}_{t-1} - \mathbf{i}_t^z$ , and  $\mathbf{h}_{t-1} - \mathbf{i}_t^z$ 2816:** Substitute  $\mathbf{z}_{t-1}$  with  $\mathbf{h}_{t-1}$  in the gate updates, *i.e.*

$$\begin{aligned}
\mathbf{i}_t &= \sigma(\mathbf{W}^{iv} \mathbf{v}_t + \mathbf{W}^{ih} \mathbf{h}_{t-1}), \quad \mathbf{f}_t = \sigma(\mathbf{W}^{fv} \mathbf{v}_t + \mathbf{W}^{fh} \mathbf{h}_{t-1}), \\
\mathbf{o}_t &= \sigma(\mathbf{W}^{ov} \mathbf{v}_t + \mathbf{W}^{oh} \mathbf{h}_{t-1}),
\end{aligned}$$

with one of the following cell updates

$$\begin{aligned}
\mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv} \mathbf{v}_t + \mathbf{W}^{ch} \mathbf{h}_{t-1} + \mathbf{W}^{cz} \mathbf{z}_{t-1}), \quad 21\text{M parameters}, \\
\mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv} \mathbf{v}_t + \mathbf{W}^{ch} \mathbf{h}_{t-1}) + \mathbf{i}_t^z \otimes \tanh(\mathbf{W}^{cz} \mathbf{z}_{t-1}), \quad 24\text{M parameters}, \\
\mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv} \mathbf{v}_t + \mathbf{W}^{ch} \mathbf{h}_{t-1}) + \mathbf{i}_t^z \otimes \tanh(\mathbf{z}_{t-1}), \quad 95\text{M parameters}.
\end{aligned}$$

As we can see from the figures, there is no other transition that is clearly preferable to the baseline, with the exception of Fishing Derby, for which transitions with 2816 hidden dimensionality perform better and converge earlier in terms of number of parameter updates.

In Figs. 35-39 we compare the baseline transition with the following convolutional alternatives (where to apply the convolutional transformations the vectors  $\mathbf{z}_{t-1}$  and  $\mathbf{v}_t$  of dimensionality 2816 are reshaped into tensors of dimension  $32 \times 11 \times 8$ )

- **$\mathcal{C}$  and  $2\mathcal{C}$** : Convolutional gate and cell updates, *i.e.*

$$\begin{aligned}\mathbf{i}_t &= \sigma(\mathcal{C}^{iv}(\mathbf{v}_t) + \mathcal{C}^{iz}(\mathbf{z}_{t-1})), & \mathbf{f}_t &= \sigma(\mathcal{C}^{fv}(\mathbf{v}_t) + \mathcal{C}^{fz}(\mathbf{z}_{t-1})), \\ \mathbf{o}_t &= \sigma(\mathcal{C}^{ov}(\mathbf{v}_t) + \mathcal{C}^{oz}(\mathbf{z}_{t-1})), \\ \mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathcal{C}^{cv}(\mathbf{v}_t) + \mathcal{C}^{cz}(\mathbf{z}_{t-1})),\end{aligned}$$

where  $\mathcal{C}$  denotes either one convolution with 32 filters of size  $3 \times 3$ , with stride 1 and padding 1 (as to preserve the input size), or two such convolutions with RReLU nonlinearity in between. These two models have around 16M parameters.

- **$\mathcal{CDA}$  and  $2\mathcal{CDA}$** : As above but with different action fusion parameters for the gate and cell updates, *i.e.*

$$\begin{aligned}\mathbf{v}_t^i &= \mathbf{W}^{ih} \mathbf{h}_{t-1} \otimes \mathbf{W}^{ia} \mathbf{a}_{t-1}, & \mathbf{v}_t^f &= \mathbf{W}^{fh} \mathbf{h}_{t-1} \otimes \mathbf{W}^{fa} \mathbf{a}_{t-1}, \\ \mathbf{v}_t^o &= \mathbf{W}^{oh} \mathbf{h}_{t-1} \otimes \mathbf{W}^{oa} \mathbf{a}_{t-1}, & \mathbf{v}_t^c &= \mathbf{W}^{ch} \mathbf{h}_{t-1} \otimes \mathbf{W}^{ca} \mathbf{a}_{t-1}, \\ \mathbf{i}_t &= \sigma(\mathcal{C}^{iv}(\mathbf{v}_t^i) + \mathcal{C}^{iz}(\mathbf{z}_{t-1})), & \mathbf{f}_t &= \sigma(\mathcal{C}^{fv}(\mathbf{v}_t^f) + \mathcal{C}^{fz}(\mathbf{z}_{t-1})), \\ \mathbf{o}_t &= \sigma(\mathcal{C}^{ov}(\mathbf{v}_t^o) + \mathcal{C}^{oz}(\mathbf{z}_{t-1})), \\ \mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathcal{C}^{cv}(\mathbf{v}_t^c) + \mathcal{C}^{cz}(\mathbf{z}_{t-1})).\end{aligned}$$

These two models have around 40M parameters.

- **$\mathbf{h}_{t-1}-\mathbf{i}_t^z$ 2816- $2\mathcal{C}$** : As ' $\mathbf{h}_{t-1}-\mathbf{i}_t^z$ 2816' with convolutional gate and cell updates, *i.e.*

$$\begin{aligned}\mathbf{i}_t &= \sigma(\mathcal{C}^{iv}(\mathbf{v}_t) + \mathcal{C}^{ih}(\mathbf{h}_{t-1})), & \mathbf{f}_t &= \sigma(\mathcal{C}^{fv}(\mathbf{v}_t) + \mathcal{C}^{fh}(\mathbf{h}_{t-1})), \\ \mathbf{o}_t &= \sigma(\mathcal{C}^{ov}(\mathbf{v}_t) + \mathcal{C}^{oh}(\mathbf{h}_{t-1})), \\ \mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathcal{C}^{cv}(\mathbf{v}_t) + \mathcal{C}^{ch}(\mathbf{h}_{t-1})) + \mathbf{i}_t^z \otimes \tanh(\mathbf{z}_{t-1}),\end{aligned}$$

where  $\mathcal{C}$  denotes two convolutions as above. This model has around 16M parameters.

- **$\mathbf{h}_{t-1}-\mathbf{i}_t^z$ 2816- $\mathcal{CDA}$  and  $\mathbf{h}_{t-1}-\mathbf{i}_t^z$ 2816- $2\mathcal{CDA}$** : As above but with different parameters for the gate and cell updates, and one or two convolutions. These two models have around 48M parameters.

- **$\mathbf{h}_{t-1}-\mathbf{i}_t^z$ 2816- $2\mathcal{CA}$** : As ' $\mathbf{h}_{t-1}-\mathbf{i}_t^z$ 2816' with convolutional action fusion, gate and cell updates, *i.e.*

$$\begin{aligned}\mathbf{v}_t &= \mathcal{C}^h(\mathbf{h}_{t-1}) \otimes \mathbf{W}^a \mathbf{a}_{t-1}, \\ \mathbf{i}_t &= \sigma(\mathcal{C}^{iv}(\mathbf{v}_t) + \mathcal{C}^{ih}(\mathbf{h}_{t-1})), & \mathbf{f}_t &= \sigma(\mathcal{C}^{fv}(\mathbf{v}_t) + \mathcal{C}^{fh}(\mathbf{h}_{t-1})), \\ \mathbf{o}_t &= \sigma(\mathcal{C}^{ov}(\mathbf{v}_t) + \mathcal{C}^{oh}(\mathbf{h}_{t-1})), \\ \mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathcal{C}^{cv}(\mathbf{v}_t) + \mathcal{C}^{ch}(\mathbf{h}_{t-1})) + \mathbf{i}_t^z \otimes \tanh(\mathbf{z}_{t-1}),\end{aligned}$$

where  $\mathcal{C}$  indicates two convolutions as above. This model has around 8M parameters.

### B.1.3 ACTION INCORPORATION

In Figs. 40-44 we compare different ways of incorporating the action for action-dependent state transitions, using prediction length  $T = 15$  and the 0%-100%PDT training scheme. More specifically, we compare the baseline structure (denoted as ' $\mathbf{W}^h \mathbf{h}_{t-1} \otimes \mathbf{W}^a \mathbf{a}_{t-1}$ ' in the figures) with the following alternatives:

- **$\mathbf{W}^h \mathbf{h}_{t-1} \otimes \mathbf{W}^{a_1} \mathbf{a}_{t-1} + \mathbf{W}^{a_2} \mathbf{a}_{t-1}$** : Multiplicative/additive interaction of the action with  $\mathbf{h}_{t-1}$ , *i.e.*  $\mathbf{v}_t = \mathbf{W}^h \mathbf{h}_{t-1} \otimes \mathbf{W}^{a_1} \mathbf{a}_{t-1} + \mathbf{W}^{a_2} \mathbf{a}_{t-1}$ . This model has around 25M parameters.
- **$\mathbf{W}^s \mathbf{s}_{t-1} \otimes \mathbf{W}^a \mathbf{a}_{t-1}$** : Multiplicative interaction of the action with the encoded frame  $\mathbf{z}_{t-1}$ , *i.e.*

$$\begin{aligned}\mathbf{v}_t &= \mathbf{W}^z \mathbf{z}_{t-1} \otimes \mathbf{W}^a \mathbf{a}_{t-1}, \\ \mathbf{i}_t &= \sigma(\mathbf{W}^{ih} \mathbf{h}_{t-1} + \mathbf{W}^{iv} \mathbf{v}_t), & \mathbf{f}_t &= \sigma(\mathbf{W}^{fh} \mathbf{h}_{t-1} + \mathbf{W}^{fv} \mathbf{v}_t), \\ \mathbf{o}_t &= \sigma(\mathbf{W}^{oh} \mathbf{h}_{t-1} + \mathbf{W}^{ov} \mathbf{v}_t), \\ \mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{ch} \mathbf{h}_{t-1} + \mathbf{W}^{cv} \mathbf{v}_t).\end{aligned}$$

This model has around 22M parameters.

- $\mathbf{W}^h \mathbf{h}_{t-1} \otimes \mathbf{W}^z \mathbf{z}_{t-1} \otimes \mathbf{W}^a \mathbf{a}_{t-1}$ : Multiplicative interaction of the action with both  $\mathbf{h}_{t-1}$  and  $\mathbf{z}_{t-1}$  in the following way

$$\begin{aligned} \mathbf{v}_t &= \mathbf{W}^h \mathbf{h}_{t-1} \otimes \mathbf{W}^z \mathbf{z}_{t-1} \otimes \mathbf{W}^a \mathbf{a}_{t-1}, \\ \mathbf{i}_t &= \sigma(\mathbf{W}^{iv} \mathbf{v}_t), \quad \mathbf{f}_t = \sigma(\mathbf{W}^{fv} \mathbf{v}_t), \quad \mathbf{o}_t = \sigma(\mathbf{W}^{ov} \mathbf{v}_t), \\ \mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv} \mathbf{v}_t). \end{aligned}$$

This model has around 19M parameters. We also considered having different matrices for the gate and cell updates (denoted in the figures as ' $\mathbf{W}^{*h} \mathbf{h}_{t-1} \otimes \mathbf{W}^{*z} \mathbf{z}_{t-1} \otimes \mathbf{W}^{*a} \mathbf{a}_{t-1}$ '). This model has around 43M parameters.

- $\mathbf{W}^h \mathbf{h}_{t-1} \otimes \mathbf{W}^{a1} \mathbf{a}_{t-1}$ : Alternative multiplicative interaction of the action with  $\mathbf{h}_{t-1}$  and  $\mathbf{z}_{t-1}$

$$\begin{aligned} \mathbf{v}_t^1 &= \mathbf{W}^h \mathbf{h}_{t-1} \otimes \mathbf{W}^{a1} \mathbf{a}_{t-1}, \quad \mathbf{v}_t^2 = \mathbf{W}^z \mathbf{z}_{t-1} \otimes \mathbf{W}^{a2} \mathbf{a}_{t-1}, \\ \mathbf{i}_t &= \sigma(\mathbf{W}^{iv1} \mathbf{v}_t^1 + \mathbf{W}^{iv2} \mathbf{v}_t^2), \quad \mathbf{f}_t = \sigma(\mathbf{W}^{fv1} \mathbf{v}_t^1 + \mathbf{W}^{fv2} \mathbf{v}_t^2), \\ \mathbf{o}_t &= \sigma(\mathbf{W}^{ov1} \mathbf{v}_t^1 + \mathbf{W}^{ov2} \mathbf{v}_t^2), \\ \mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv1} \mathbf{v}_t^1 + \mathbf{W}^{cv2} \mathbf{v}_t^2). \end{aligned}$$

This model has around 28M parameters. We also considered having different matrices for the gate and cell updates (denoted in the figures as ' $\mathbf{W}^{*h} \mathbf{h}_{t-1} \otimes \mathbf{W}^{*a1} \mathbf{a}_{t-1}$ '). This model has around 51M parameters.

- **As Input:** Consider the action as an additional input, *i.e.*

$$\begin{aligned} \mathbf{i}_t &= \sigma(\mathbf{W}^{ih} \mathbf{h}_{t-1} + \mathbf{W}^{iz} \mathbf{z}_{t-1} + \mathbf{W}^{ia} \mathbf{a}_{t-1}), \\ \mathbf{f}_t &= \sigma(\mathbf{W}^{fh} \mathbf{h}_{t-1} + \mathbf{W}^{fz} \mathbf{z}_{t-1} + \mathbf{W}^{fa} \mathbf{a}_{t-1}), \\ \mathbf{o}_t &= \sigma(\mathbf{W}^{oh} \mathbf{h}_{t-1} + \mathbf{W}^{oz} \mathbf{z}_{t-1} + \mathbf{W}^{oa} \mathbf{a}_{t-1}), \\ \mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{ch} \mathbf{h}_{t-1} + \mathbf{W}^{cz} \mathbf{z}_{t-1} + \mathbf{W}^{ca} \mathbf{a}_{t-1}). \end{aligned}$$

This model has around 19M parameters.

- $\mathcal{CA}$ : Combine the action with the frame, by replacing the encoding with

$$\mathbf{z}_{t-1} = \mathcal{C}(\mathcal{A}(\mathbb{I}(\hat{\mathbf{x}}_{t-1}, \mathbf{x}_{t-1}), \mathbf{a}_{t-1})),$$

where  $\mathcal{A}$  indicates an augmenting operation: the frame of dimension  $nC = 3 \times nH = 210 \times nW = 160$  is augmented with  $nA$  (number of actions) full-zero or full-one matrices of dimension  $nH \times nW$ , producing a tensor of dimension  $(nC + nA) \times nH \times nW$ . As the output of the first convolution can be written as

$$\mathbf{y}_{j,k,l} = \sum_{h=1}^{nH} \sum_{w=1}^{nW} \left\{ \sum_{i=1}^{nC} \mathbf{W}^{i,j}_{h,w} \mathbf{x}_{i,h+dH(k-1),w+dW(l-1)} + \mathbf{x}_{nC+a,h+dH(k-1),w+dW(l-1)} \right\},$$

where  $dH$  and  $dW$  indicate the filter strides, with this augmentation the action has a local linear interaction. This model has around 19M parameters.

As we can see from the figures, ' $\mathcal{CA}$ ' is generally considerably worse than the other structures.

#### ACTION-INDEPENDENT VERSUS ACTION-DEPENDENT STATE TRANSITION

In Fig. 45, we compare the baseline structure with one that is action-independent as in Oh et al. (2015), using prediction length  $T = 15$  and the 0%-100%PDT training scheme.

As we can see, having an action-independent state transition generally gives worse performance in the games with higher error. An interesting disadvantage of such a structure is its inability to predict the moving objects around the agent in Seaquest. This can be noticed in the videos in [Seaquest](#), which show poor modelling of the fish. This structure also makes it more difficult to correctly update the score in some games such as Seaquest and Fishing Derby.

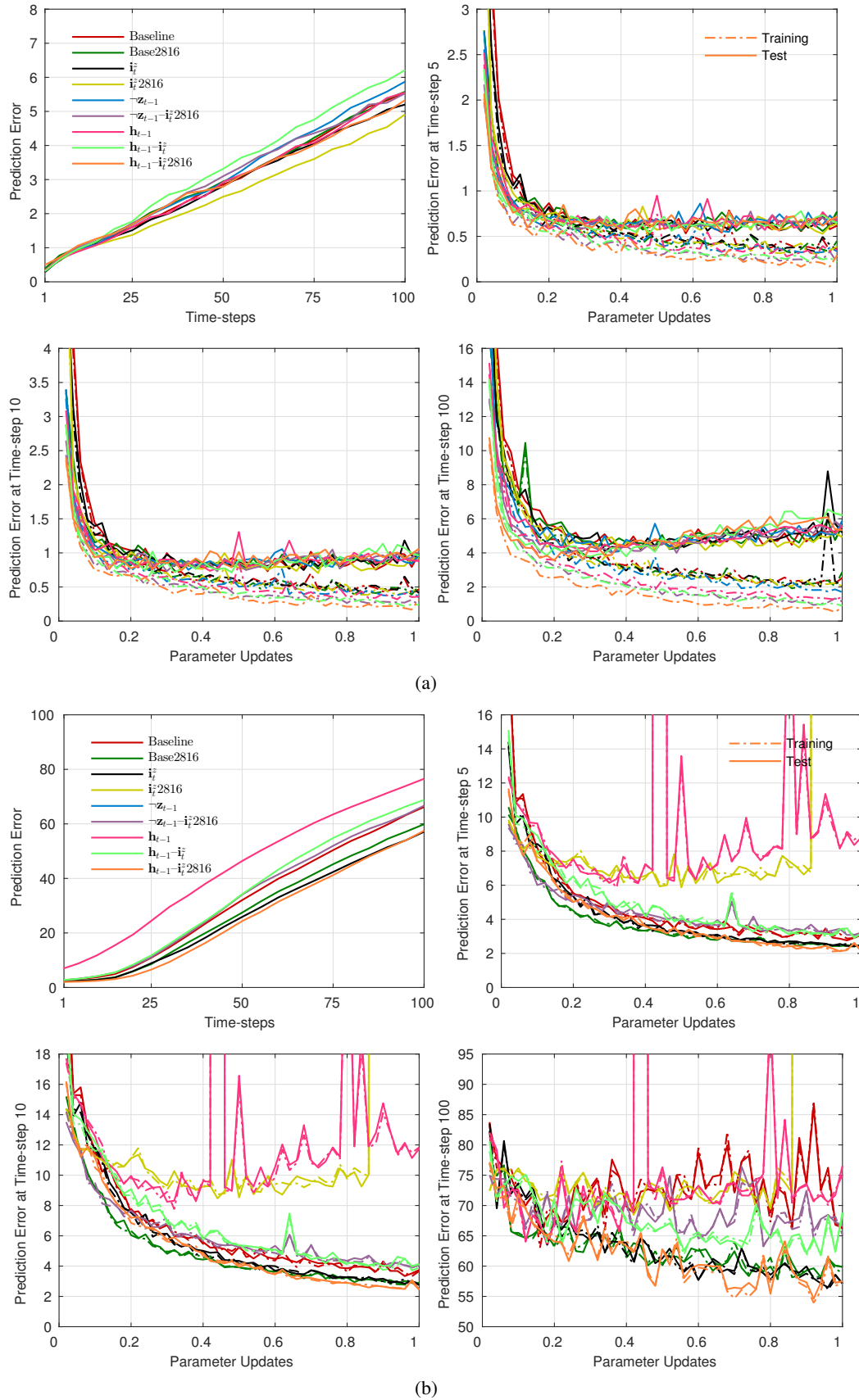


Figure 30: Prediction error (average over 10,000 sequences) for different action-dependent state transitions on (a) Bowling and (b) Breakout. Parameter updates are in millions.

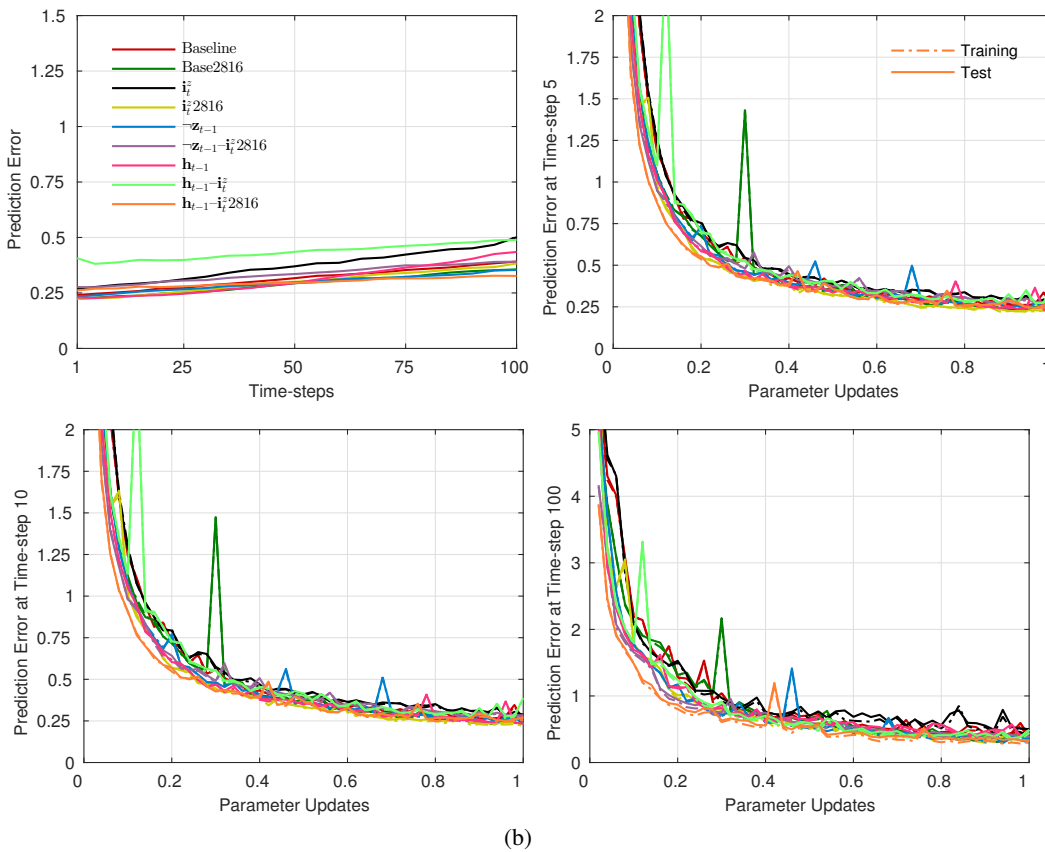
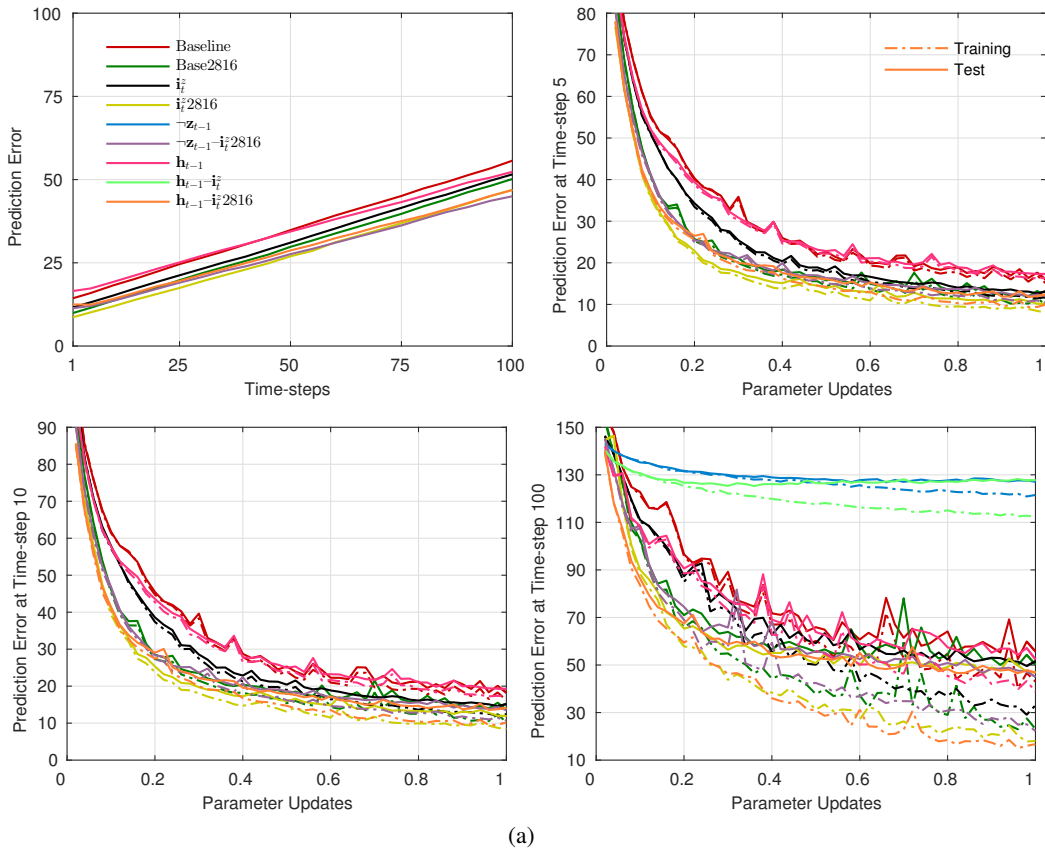


Figure 31: Prediction error for different action-dependent state transitions on (a) Fishing Derby and (b) Freeway.

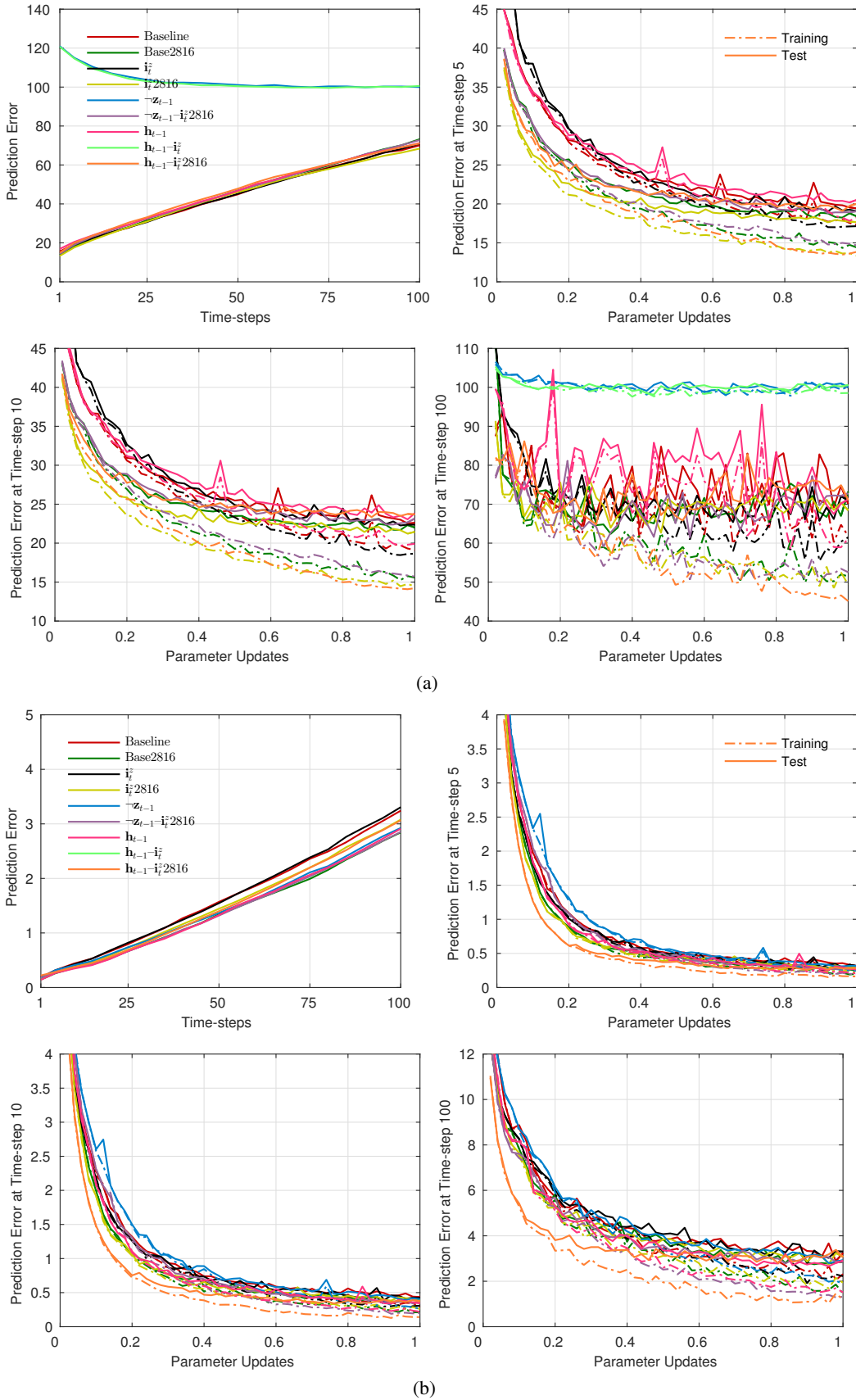


Figure 32: Prediction error for different action-dependent state transitions on (a) Ms Pacman and (b) Pong.



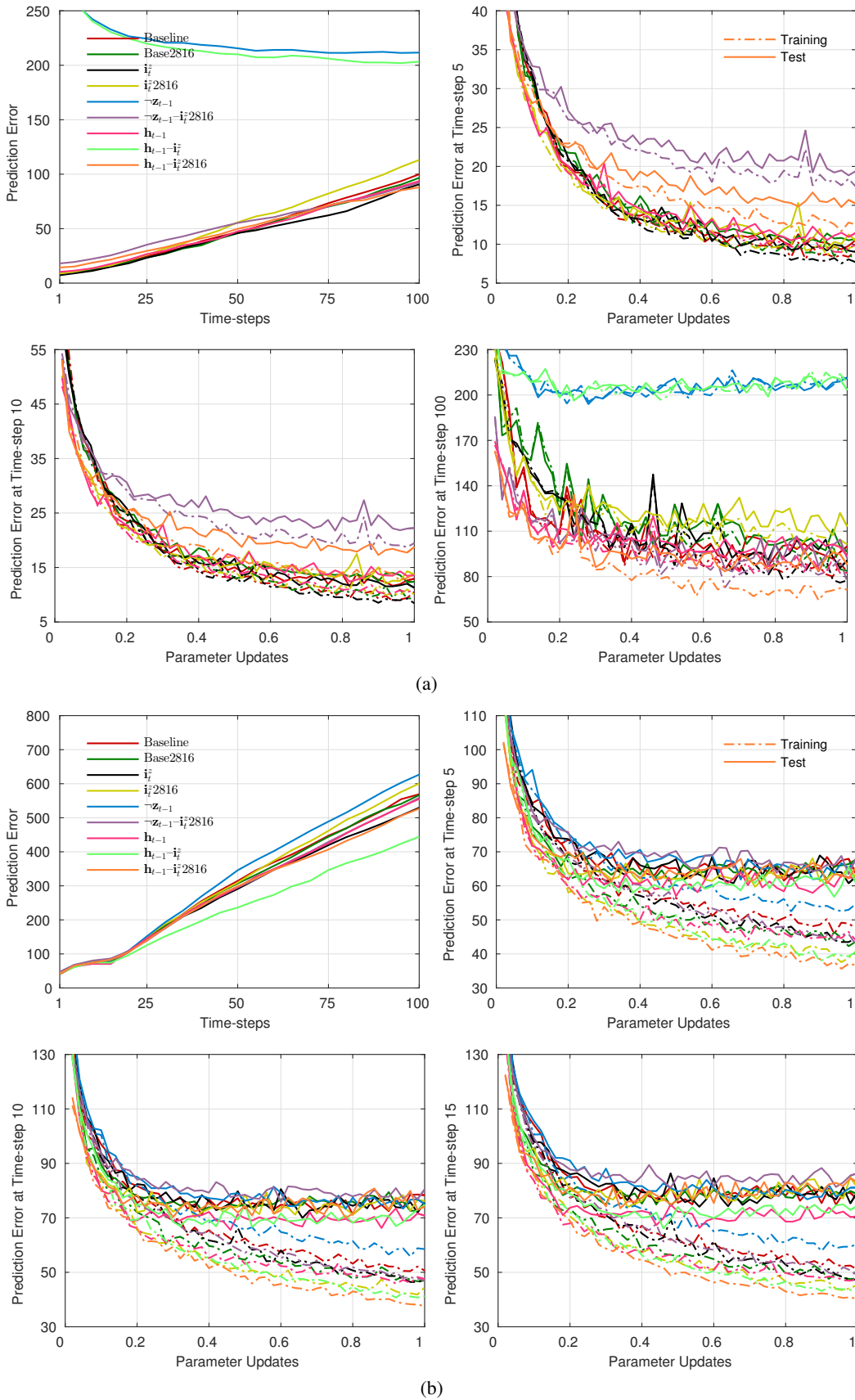


Figure 33: Prediction error for different action-dependent state transitions on (a) Qbert and (b) Riverraid.

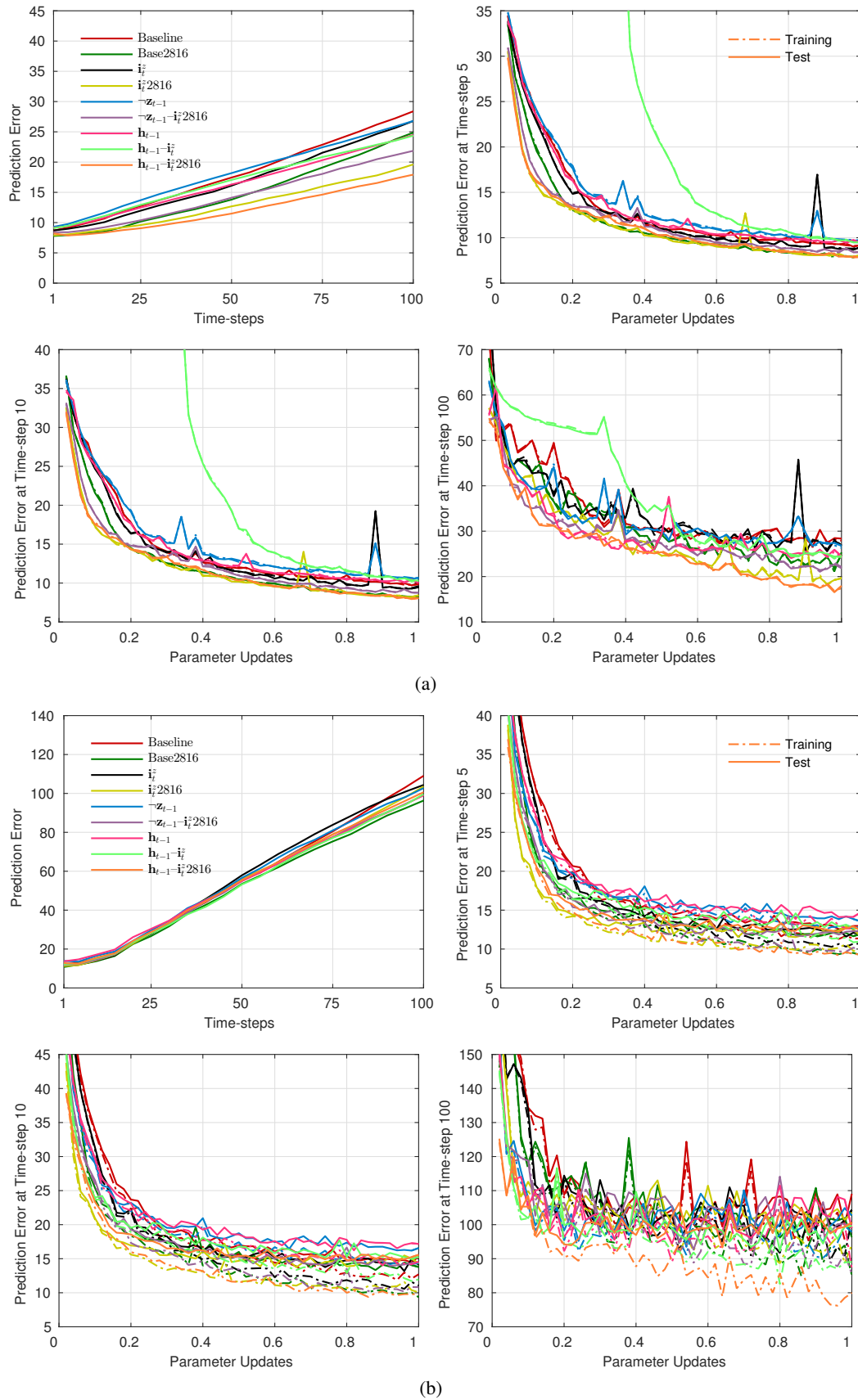


Figure 34: Prediction error for different action-dependent state transitions on (a) Seaquest and (b) Space Invaders.

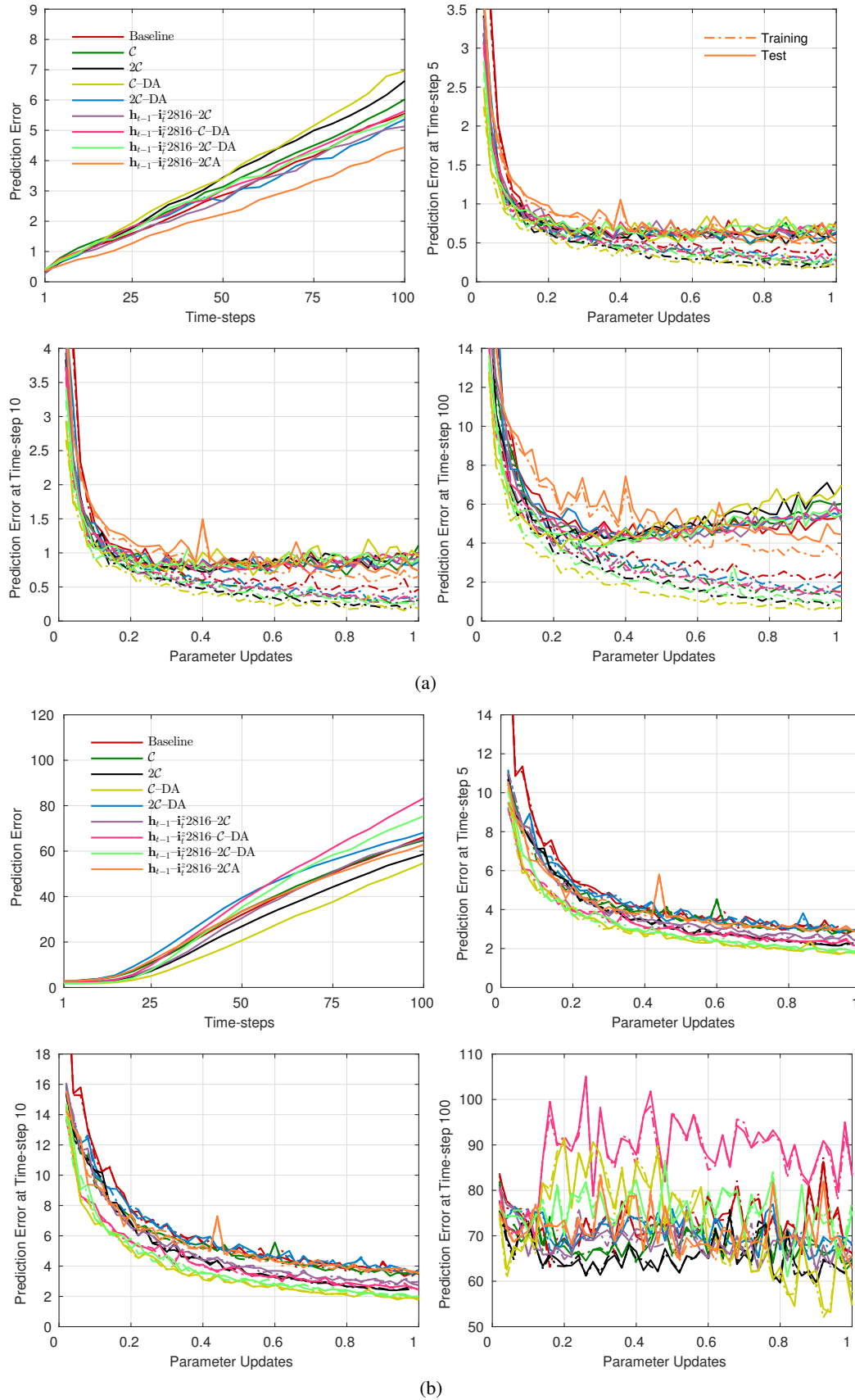


Figure 35: Prediction error (average over 10,000 sequences) for different convolutional action-dependent state transitions on (a) Bowling and (b) Breakout. Parameter updates are in millions.

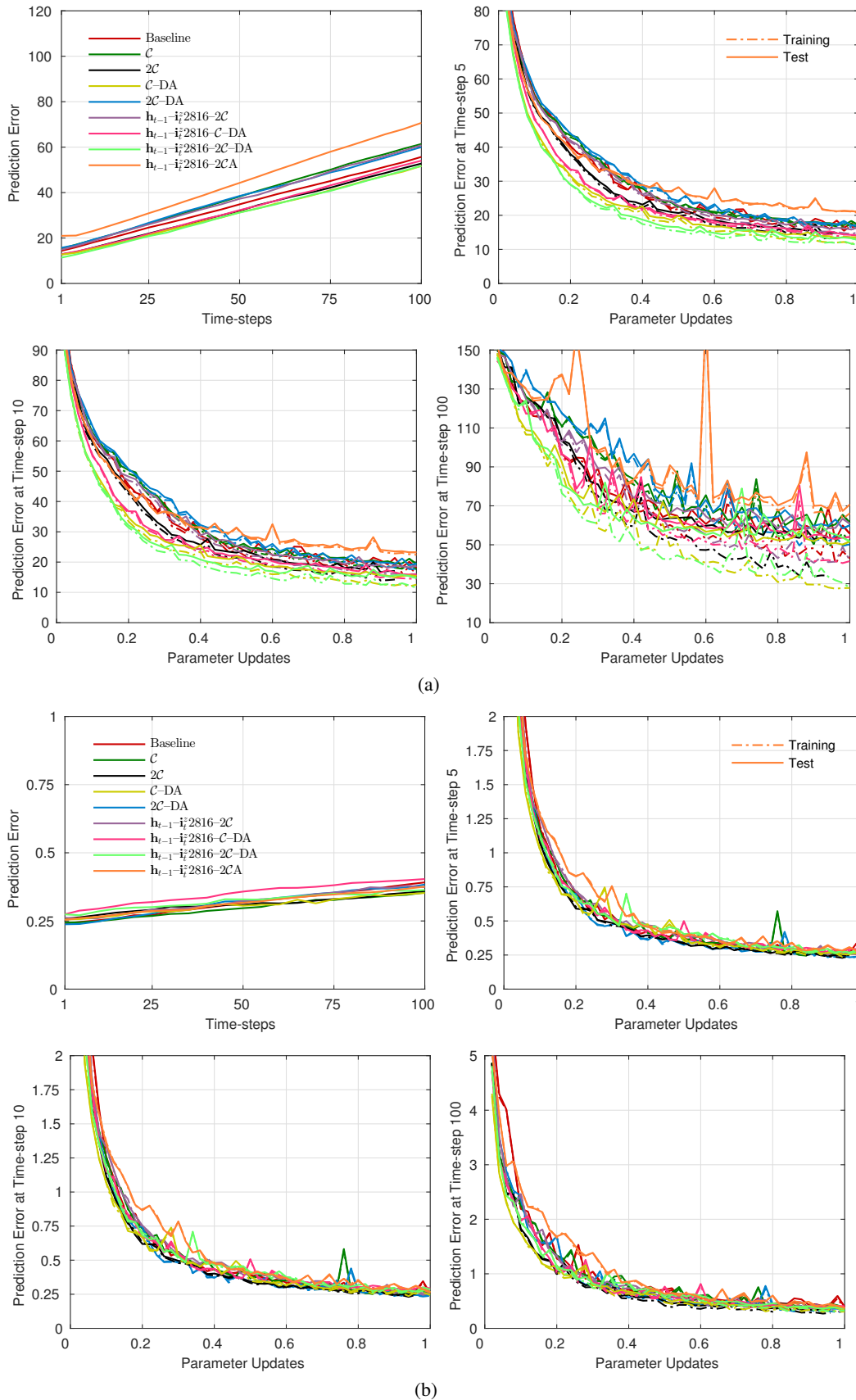


Figure 36: Prediction error for different convolutional action-dependent state transitions on (a) Fishing Derby and (b) Freeway.

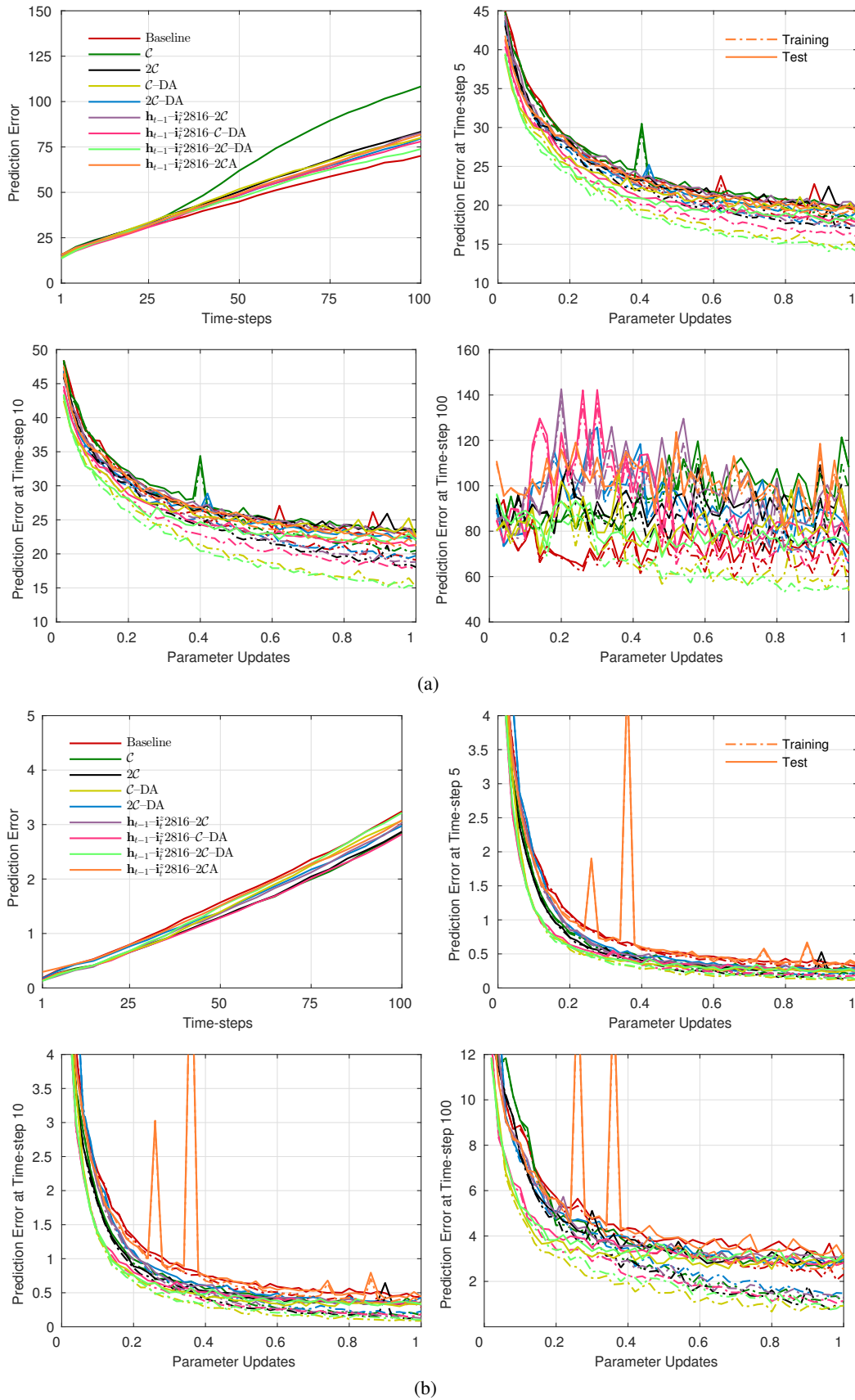


Figure 37: Prediction error for different convolutional action-dependent state transitions on (a) Ms Pacman and (b) Pong.

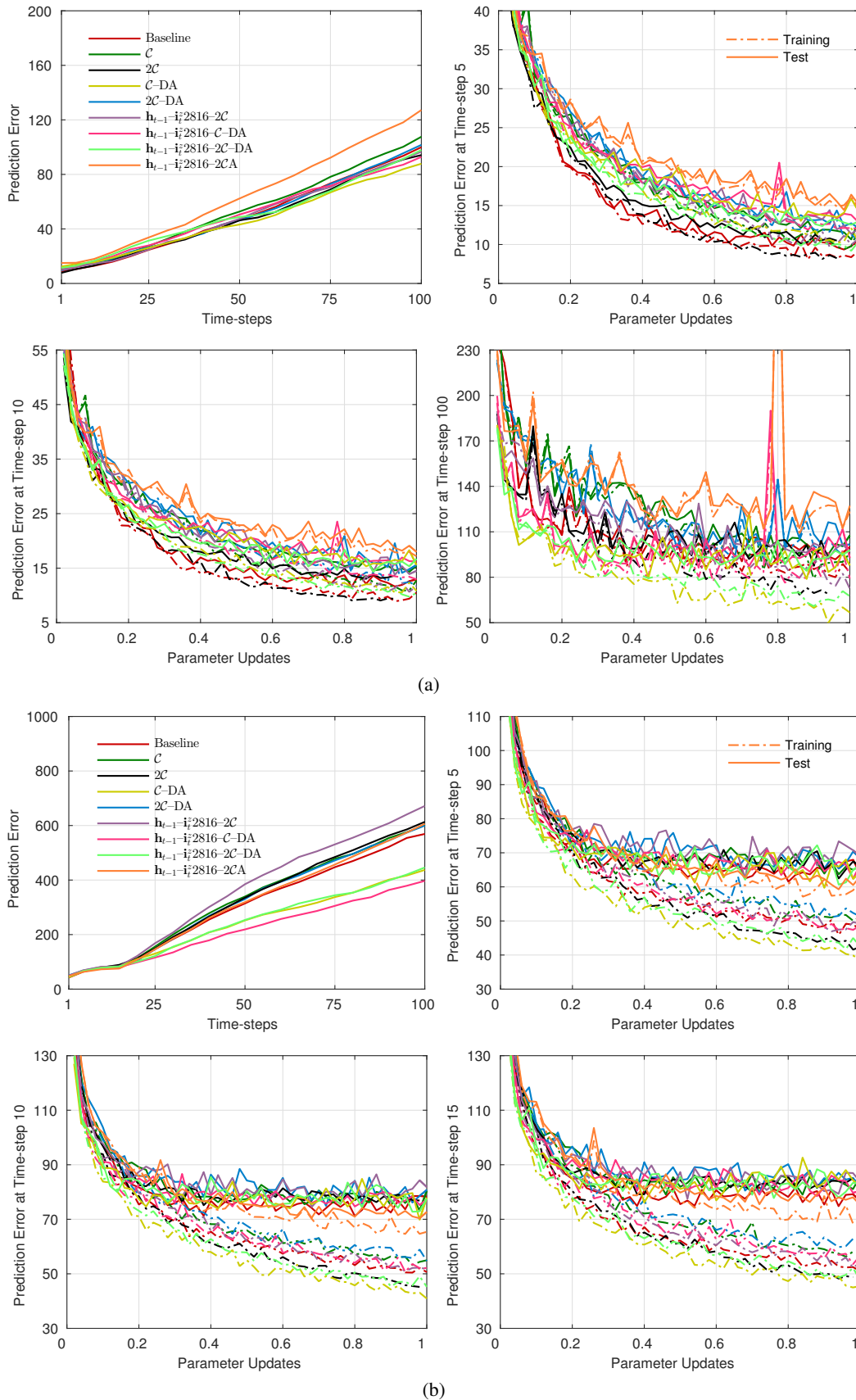


Figure 38: Prediction error for different convolutional action-dependent state transitions on (a) Qbert and (b) Riverraid.

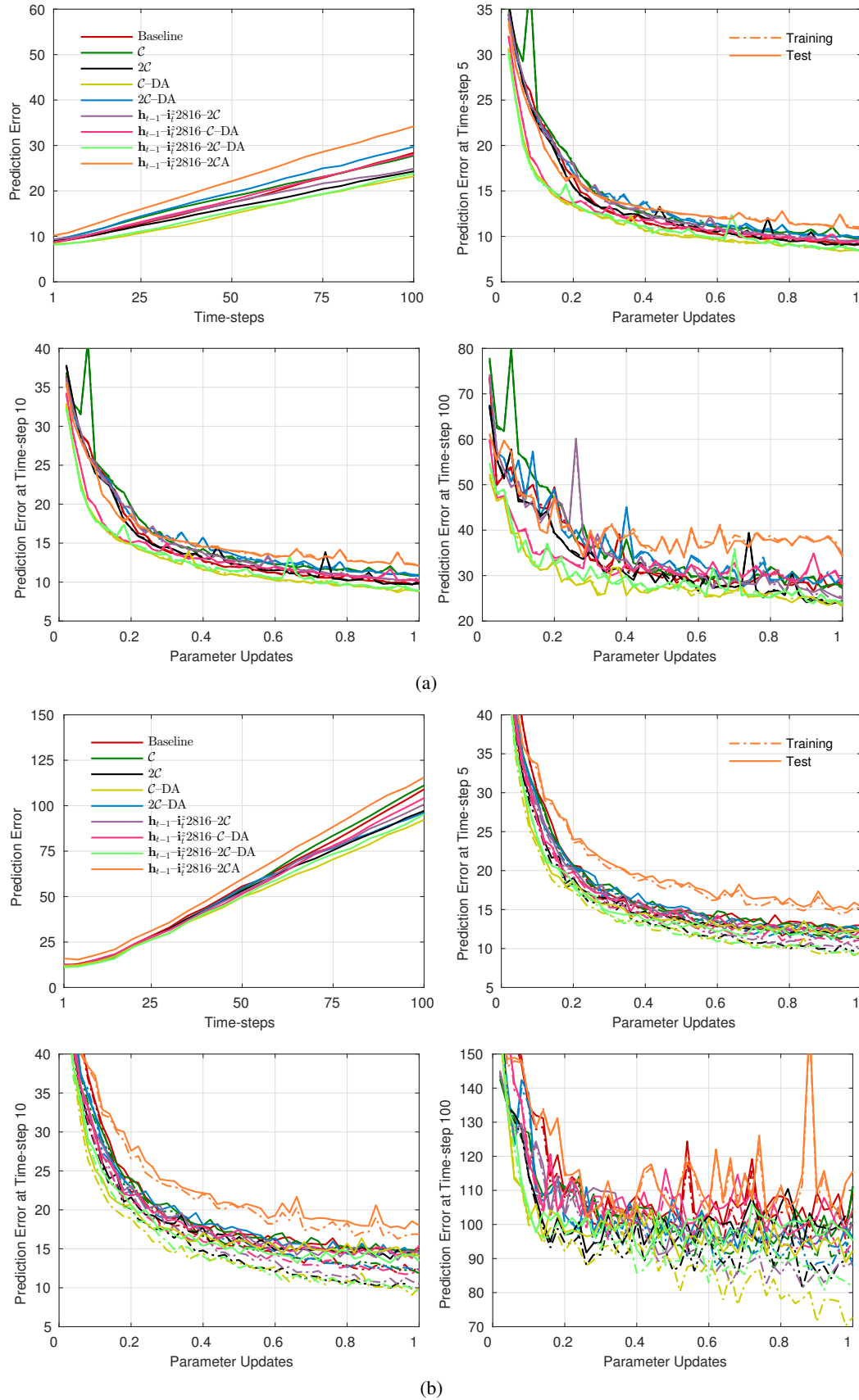


Figure 39: Prediction error for different convolutional action-dependent state transitions on (a) Seaquest and (b) Space Invaders.

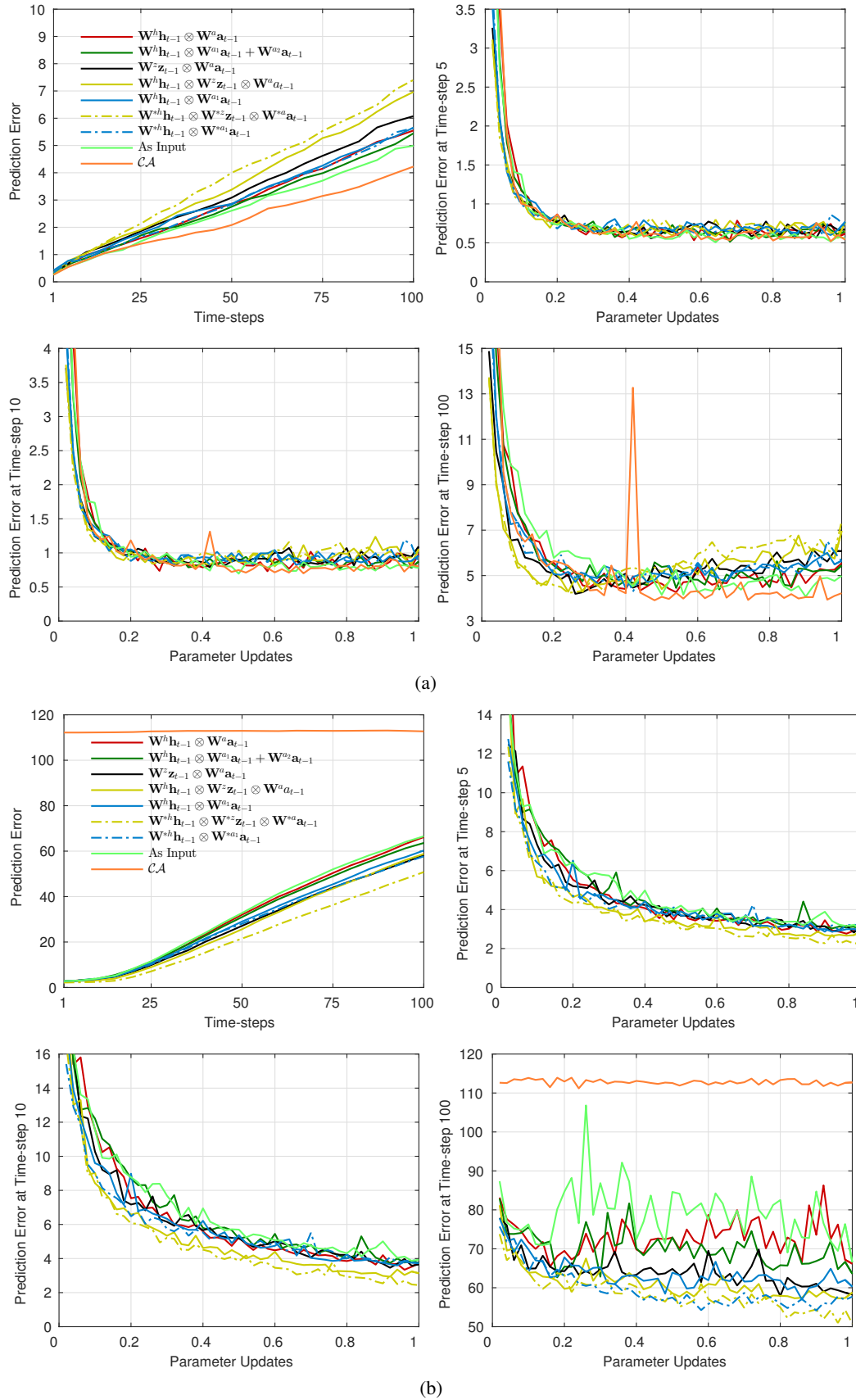


Figure 40: Prediction error (average over 10,000 sequences) for different ways of incorporating the action on (a) Bowling and (b) Breakout. Parameter updates are in millions.



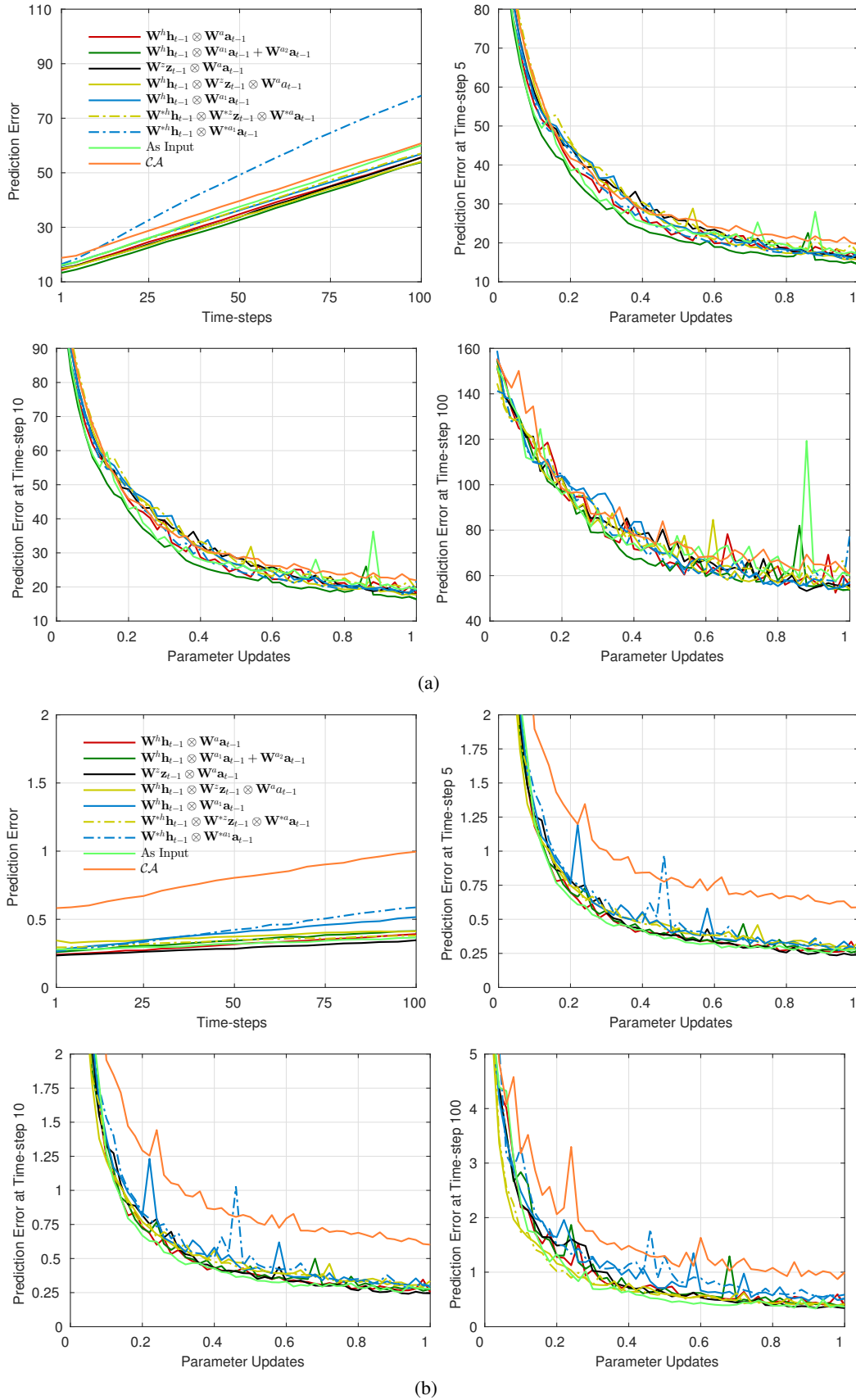


Figure 41: Prediction error for different ways of incorporating the action on (a) Fishing Derby and (b) Freeway.

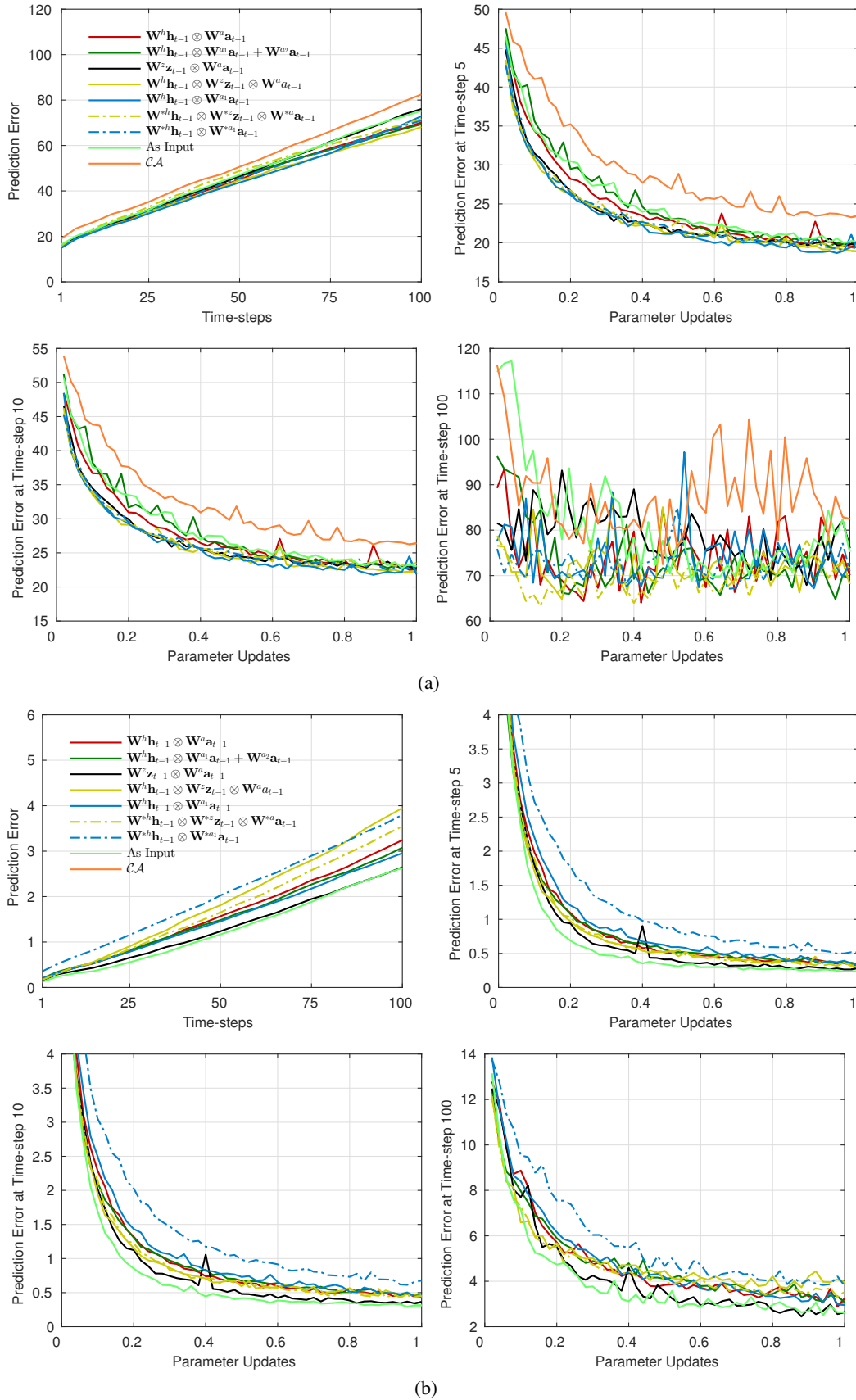


Figure 42: Prediction error for different ways of incorporating the action on (a) Ms Pacman and (b) Pong.

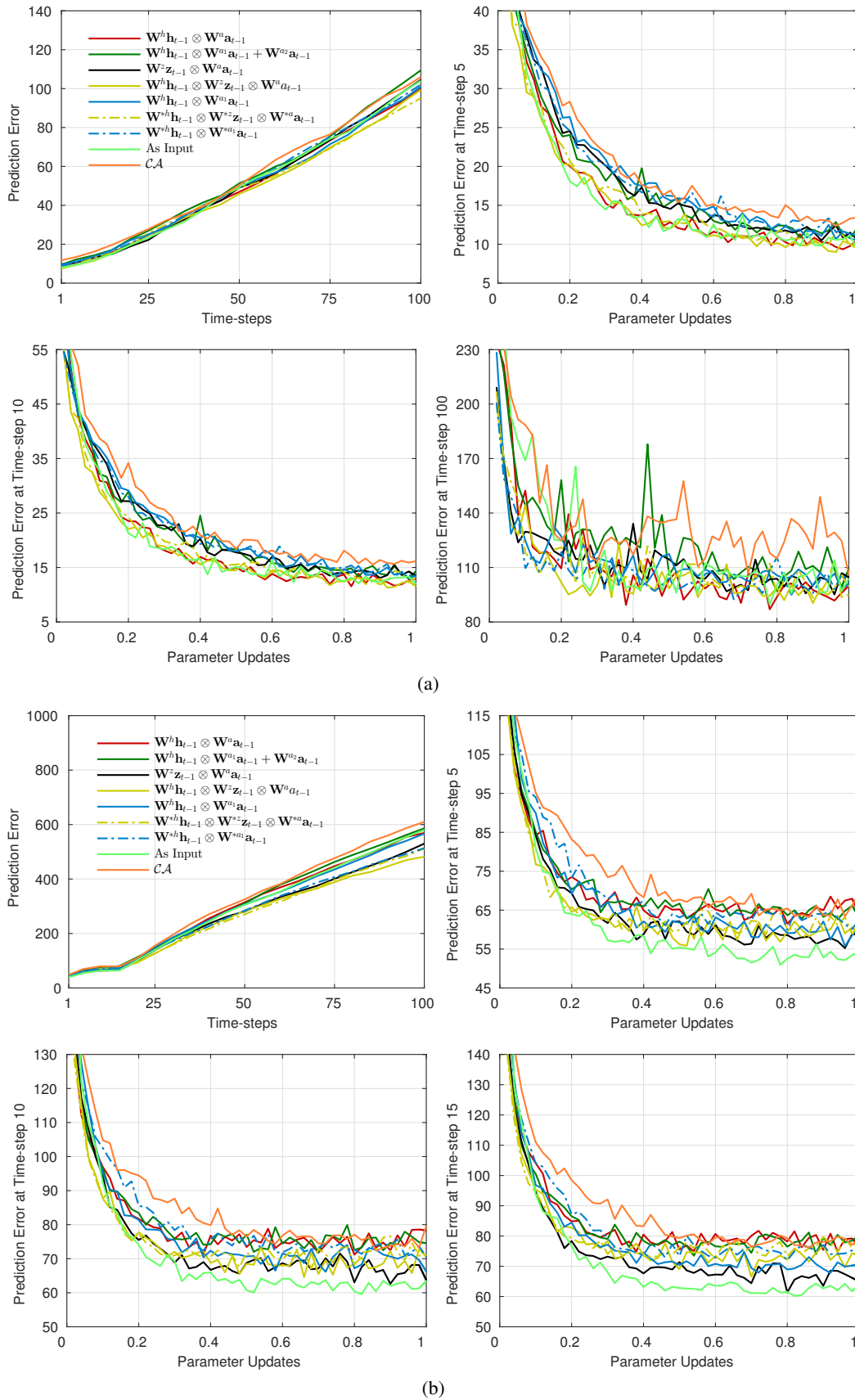


Figure 43: Prediction error for different ways of incorporating the action on (a) Qbert and (b) Riverraid.

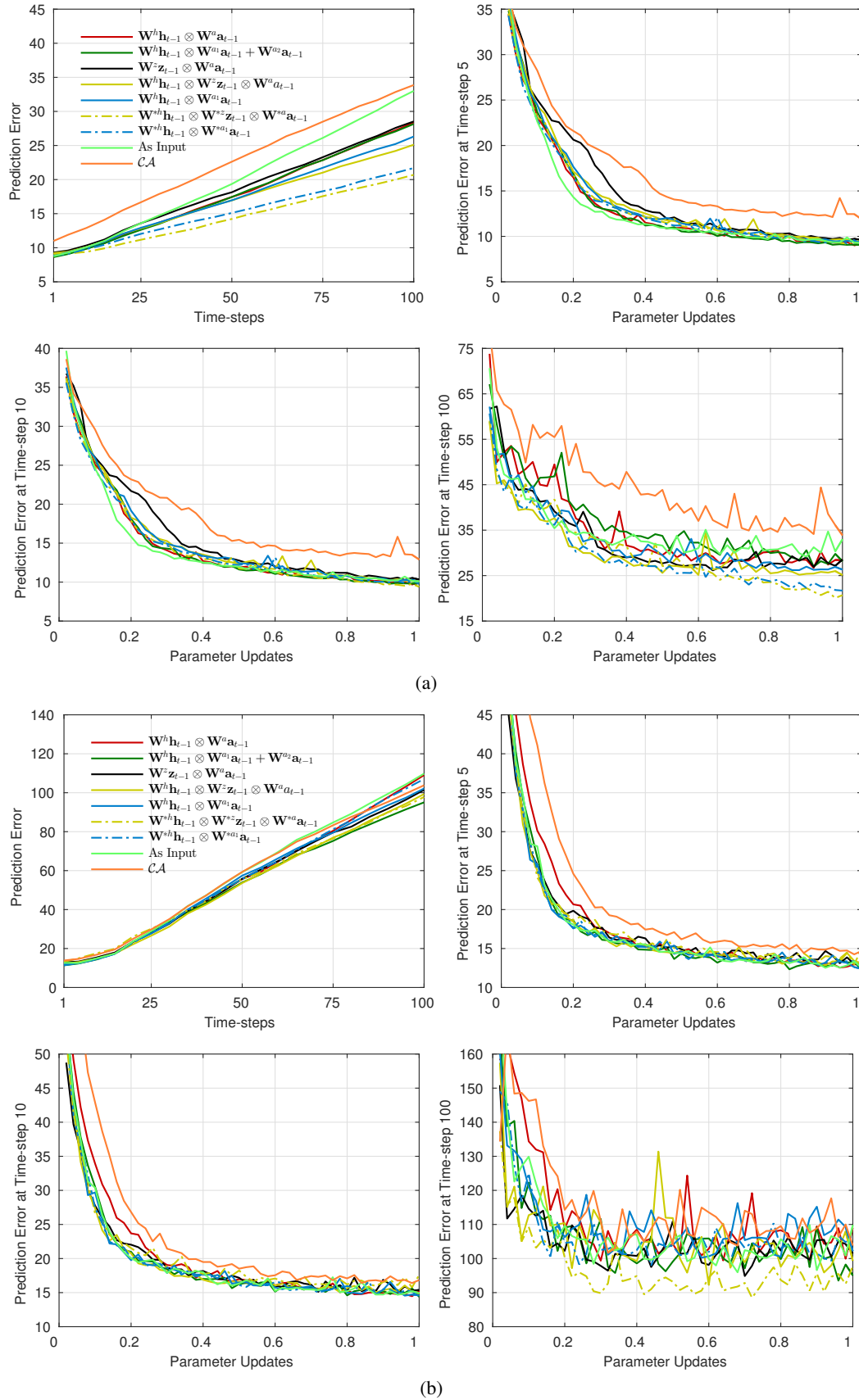


Figure 44: Prediction error for different ways of incorporating the action on (a) Seaquest and (b) Space Invaders.

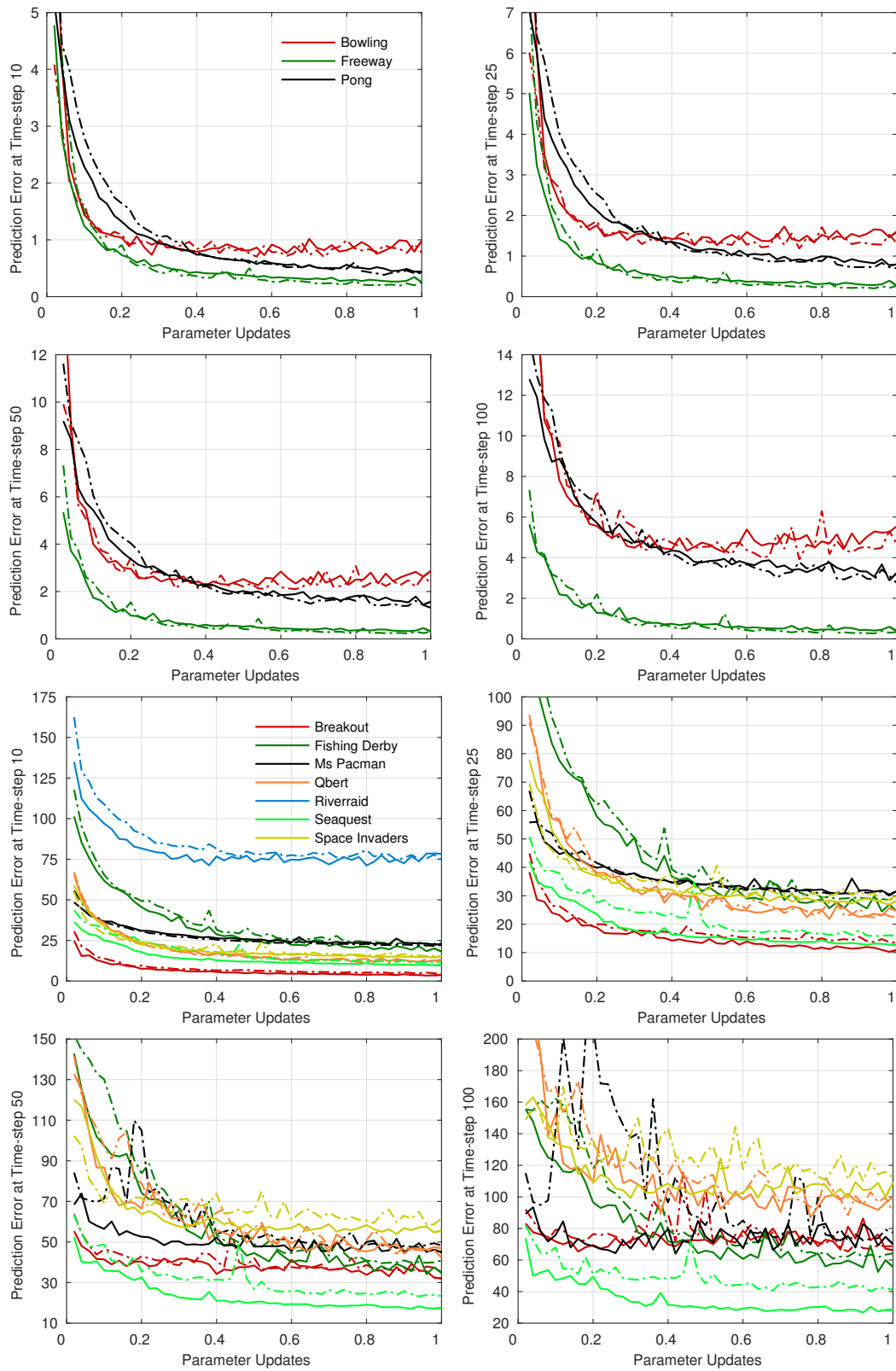


Figure 45: Prediction error (average over 10,000 sequences) with (continuous lines) action-dependent and (dashed lines) action-independent state transition. Parameter updates are in millions.

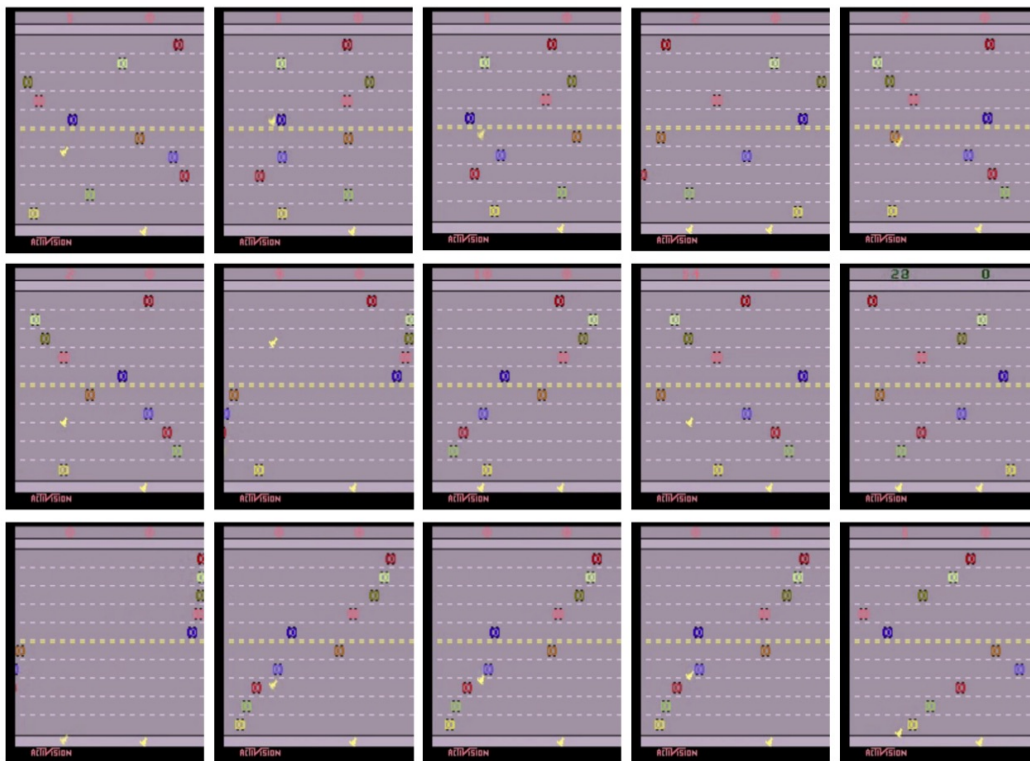


Figure 46: Salient frames extracted from 2000 frames of Freeway generated using our simulator with actions chosen by a human player.

#### B.1.4 HUMAN PLAY

In Fig. 46, we show the results of a human playing Freeway for 2000 time-steps (the corresponding video is available at [Freeway-HPlay](#)). The model is able to update the score correctly up to (14,0). At that point the score starts flashing and to change color as a warn to the resetting of the game. The model is not able to predict the score correctly in this warning phase, due to the bias in the data (DQN always achieves score above 20 at this point in the game), but flashing starts at the right time as does the resetting of the game.

Figs. 47 and 48 are larger views of the same frames shown in Fig. 7.

## B.2 3D CAR RACING

We generated 10 and one million ( $180 \times 180$ ) RGB images for training and testing respectively, with an agent trained with the asynchronous advantage actor critic algorithm (Fig. 2 in (Mnih et al., 2016)). The agent could choose among the three actions accelerate straight, accelerate left, and accelerate right, according to an  $\epsilon$ -greedy policy, with  $\epsilon$  selected at random between 0 and 0.5, independently for each episode. We added a 4th ‘do nothing’ action when generating actions at random. Smaller  $\epsilon$  lead to longer episodes ( $\sim 1500$  frames), while larger  $\epsilon$  lead to shorter episodes ( $\sim 200$  frames).

We could use the same number of convolutional layers, filters and kernel sizes as in Atari, with no padding.

Fig. 49 shows side by side predicted and real frames for up to 200 actions. We found that this quality of predictions was very common.

When using our model as an interactive simulator, we observed that the car would slightly slow down when selecting no action, but fail to stop. Since the model had never seen occurrences of the agent

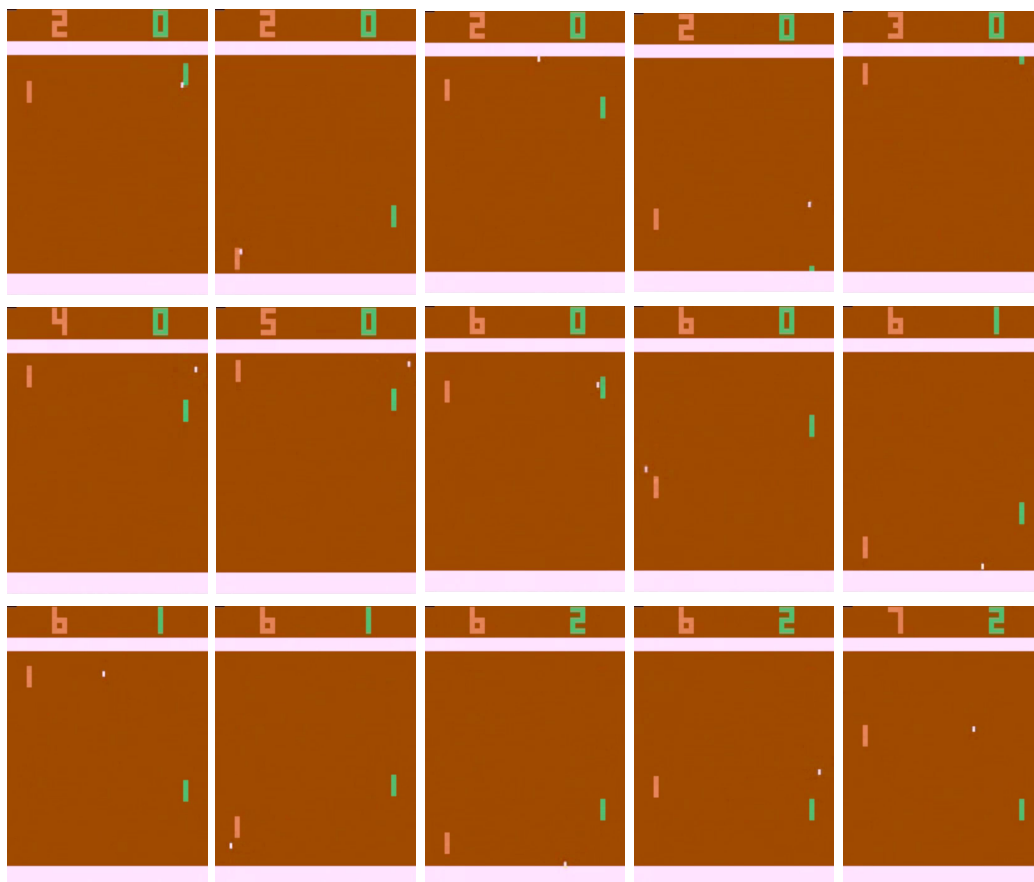


Figure 47: Salient frames extracted from 500 frames of Pong generated using our simulator with actions chosen by a human player.

completely releasing the accelerator for more than a few consecutive actions, it makes sense it would fail to deal with this case appropriately.

### B.3 3D MAZES

Unlike Atari and TORCS, we could rely on agents with random policies to generate interesting sequences. The agent could choose one of five actions: forward, backward, rotate left, rotate right or do nothing. During an episode, the agent alternated between a random walk for 15 steps, and spinning on itself for 15 steps (roughly, a complete  $360^\circ$  spin). This encourages coherent learning of the predicted frames after a spin. The random walk was with dithering of 0.7, meaning that new actions were chosen with a probability of 0.7 at every time-step. The training and test datasets were made of 7,600 and 1,100 episodes, respectively. All episodes were of length 900 frames, resulting in 6,840,000 and 990,000 ( $48 \times 48$ ) RGB images for training and testing respectively.

We adapted the encoding by having only 3 convolutions with 64 filters of size  $6 \times 6$ , stride 2, and padding 0, 1, and 2. The decoding transformation was adapted accordingly.

### B.4 MODEL-BASED EXPLORATION

We observed that increasing the number of Monte-Carlo simulations beyond 100 made little to no difference, probably because with  $n_a$  possible actions the number of possible Monte-Carlo simulations  $n_a^d$  is so large that we quickly get diminishing returns with every new simulation.

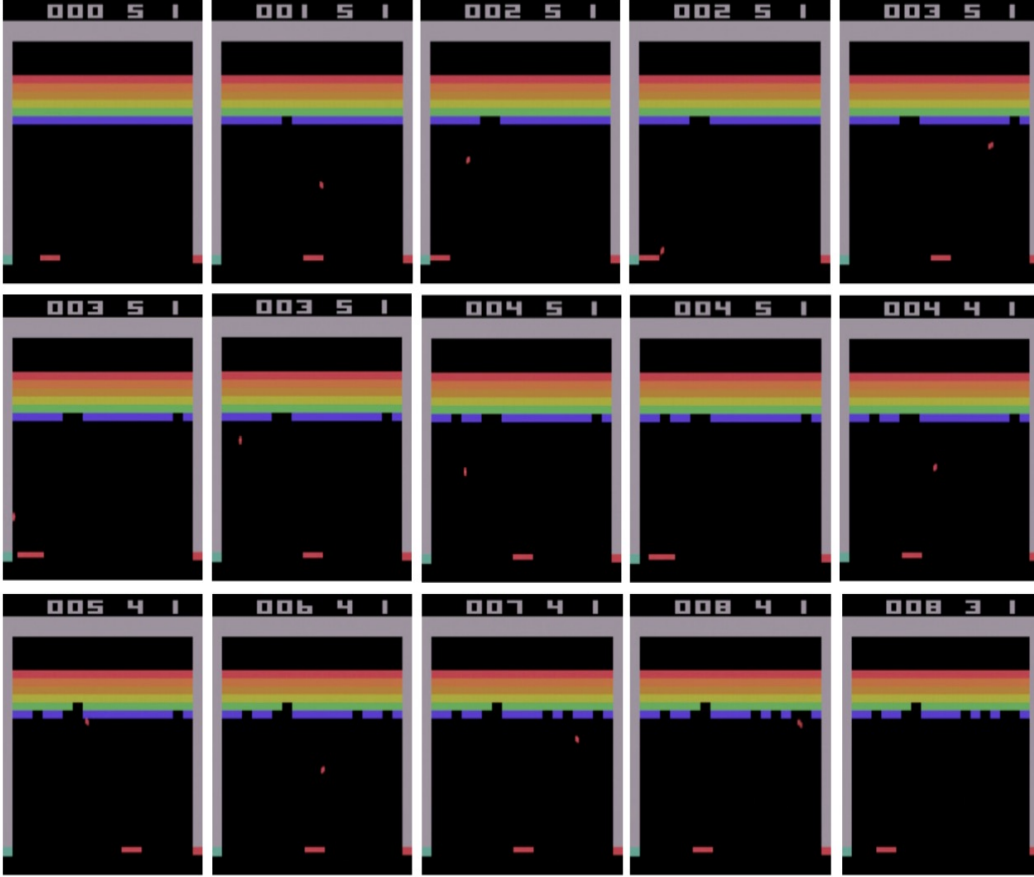


Figure 48: Salient frames extracted from 350 frames of Breakout generated using our simulator with actions taken by a human player.

Increasing significantly the sequence length of actions beyond  $d = 6$  lead to a large decrease in performance. To explain this, we observed that after 6 steps, our average prediction error was less than half the average prediction error after 30 steps (0.16 and 0.37 respectively). Since the average minimum and maximum distances did not vary significantly (from 0.23 to 0.36, and from 0.24 to 0.4 respectively), for deep simulations we ended up with more noise than signal in our predictions and our decisions were no better than random.

Fig. 50 shows some examples of trajectories chosen by our explorer. Note that all these trajectories are much smoother than for our baseline agent.

## C PREDICTION-INDEPENDENT SIMULATORS

In this section we compare different action-dependent state transitions and prediction lengths  $T$  for the prediction-independent simulator.

More specifically, in Fig. 51 we compare (with  $T = 15$ ) the state transition

$$\text{Encoding: } \mathbf{z}_{t-1} = \begin{cases} \mathcal{C}(\mathbf{x}_{t-1}) & \text{Up to } t-1 = \tau-1, \\ \mathbf{h}_{t-1} & \text{From } t-1 = \tau, \end{cases}$$

$$\text{Action fusion: } \mathbf{v}_t = \mathbf{W}^h \mathbf{h}_{t-1} \otimes \mathbf{W}^a \mathbf{a}_{t-1},$$

$$\text{Gate update: } \mathbf{i}_t = \sigma(\mathbf{W}^{iv} \mathbf{v}_t + \mathbf{W}^{iz} \mathbf{z}_{t-1}), \quad \mathbf{f}_t = \sigma(\mathbf{W}^{fv} \mathbf{v}_t + \mathbf{W}^{fz} \mathbf{z}_{t-1}),$$

$$\mathbf{o}_t = \sigma(\mathbf{W}^{ov} \mathbf{v}_t + \mathbf{W}^{oz} \mathbf{z}_{t-1}),$$

$$\text{Cell update: } \mathbf{c}_t = \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv} \mathbf{v}_t + \mathbf{W}^{cz} \mathbf{z}_{t-1}),$$



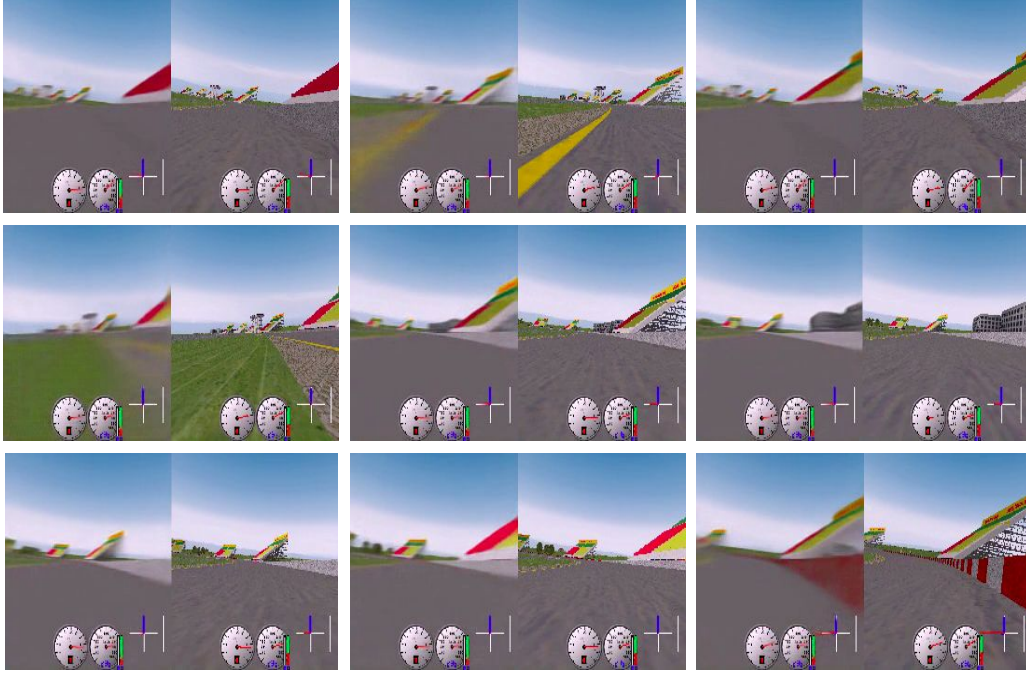


Figure 49: Salient frames, predicted (left) and real (right), for TORCS from a 200 time-steps video.

where the vectors  $\mathbf{h}_{t-1}$  and  $\mathbf{v}_t$  have dimension 1024 and 2048 respectively and with different matrices  $\mathbf{W}$  for the warm-up and the prediction phases (the resulting model has around 40M parameters – we refer to this structure as ‘Base- $\mathbf{z}_{t-1} = \mathbf{h}_{t-1}$ ’ in the figures), with the following alternatives:

- **Base- $\mathbf{z}_{t-1} = \mathbf{0}$ :** Remove the action-independent transformation of  $\mathbf{h}_{t-1}$ , *i.e.*

$$\mathbf{z}_{t-1} = \begin{cases} \mathcal{C}(\mathbf{x}_{t-1}) & \text{Up to } t-1 = \tau-1, \\ \mathbf{0} & \text{From } t-1 = \tau, \end{cases}$$

where  $\mathbf{0}$  represents a zero-vector and with different matrices  $\mathbf{W}$  for the warm-up and the prediction phases. This model has around 40M parameters.

- **$\mathbf{h}_{t-1} - \mathbf{i}_t^z 2816 - \mathbf{z}_{t-1} = \mathbf{0}$ :** Substitute  $\mathbf{z}_{t-1}$  with  $\mathbf{h}_{t-1}$  in the gate updates and have a separate gating for the encoded frame, *i.e.*

$$\begin{aligned} \mathbf{z}_{t-1} &= \begin{cases} \mathcal{C}(\mathbf{x}_{t-1}) & \text{Up to } t-1 = \tau-1, \\ \mathbf{0} & \text{From } t-1 = \tau, \end{cases} \\ \mathbf{v}_t &= \mathbf{W}^h \mathbf{h}_{t-1} \otimes \mathbf{W}^a \mathbf{a}_{t-1}, \\ \mathbf{i}_t &= \sigma(\mathbf{W}^{iv} \mathbf{v}_t + \mathbf{W}^{ih} \mathbf{h}_{t-1}), \quad \mathbf{f}_t = \sigma(\mathbf{W}^{fv} \mathbf{v}_t + \mathbf{W}^{fh} \mathbf{h}_{t-1}), \\ \mathbf{o}_t &= \sigma(\mathbf{W}^{ov} \mathbf{v}_t + \mathbf{W}^{oh} \mathbf{h}_{t-1}), \\ \mathbf{c}_t &= \mathbf{f}_t \otimes \mathbf{c}_{t-1} + \mathbf{i}_t \otimes \tanh(\mathbf{W}^{cv} \mathbf{v}_t + \mathbf{W}^{cs} \mathbf{h}_{t-1}) + \mathbf{i}_t^z \otimes \tanh(\mathbf{z}_{t-1}), \end{aligned}$$

with shared  $\mathbf{W}$  matrices for the warm-up and the prediction phases, without RReLU after the last convolution of the encoding, and with vectors  $\mathbf{h}_{t-1}$  and  $\mathbf{v}_t$  of dimensionality 2816. This model has around 95M parameters.

As we can see, the ‘Base- $\mathbf{z}_{t-1} = \mathbf{0}$ ’ state transition performs quite poorly for long-term prediction compared to the other transitions. With this transition, the prediction-independent simulator performs much worse than the prediction-dependent simulator with the baseline state transition (Appendix B.1.1). The best performance is obtained with the ‘ $\mathbf{h}_{t-1} - \mathbf{i}_t^z 2816 - \mathbf{z}_{t-1} = \mathbf{0}$ ’ structure, which however has a large number of parameters.

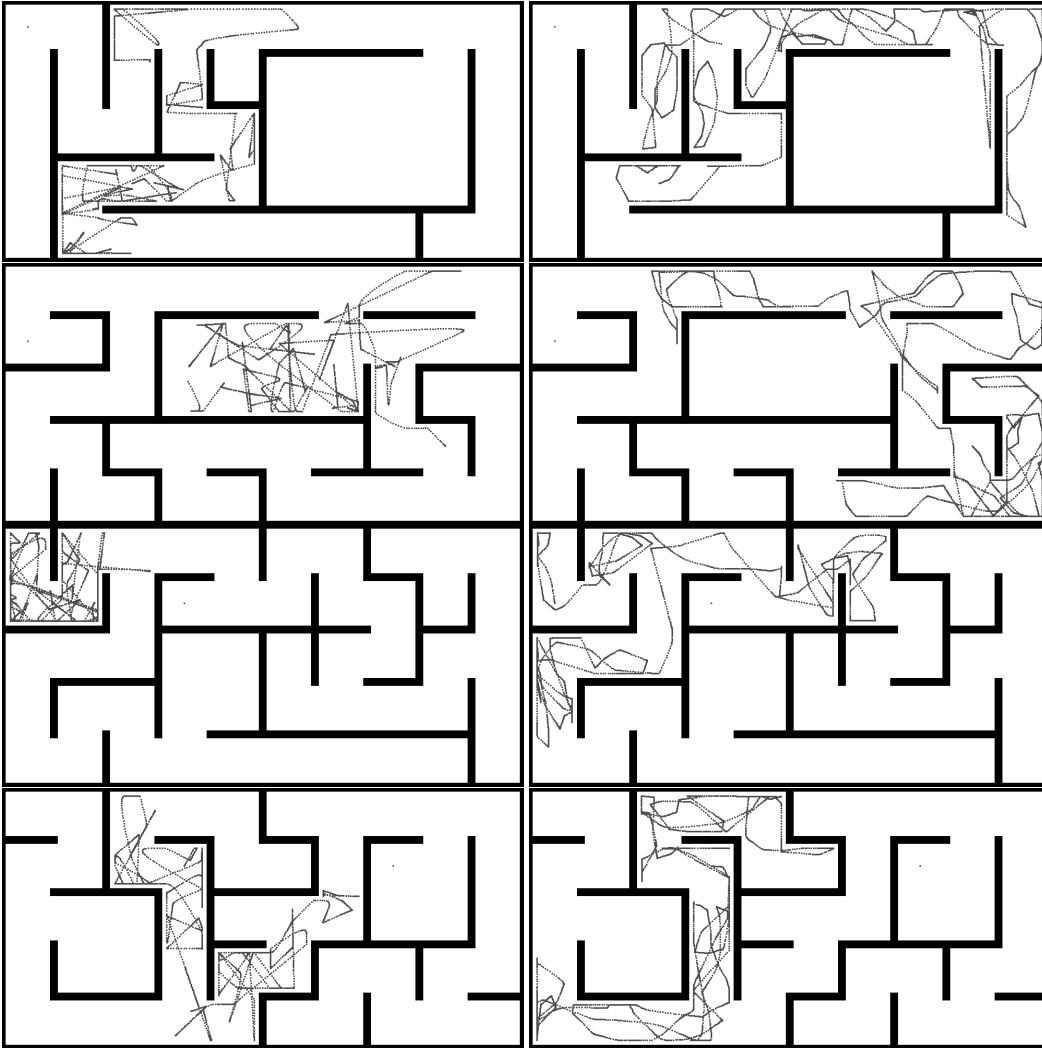


Figure 50: Examples of paths followed by random baseline (left), and explorers based on our simulator (right).

In Figs. 52 and 53, we show the effect of using different prediction lengths  $T$  on the structure 'Base- $\mathbf{z}_{t-1} = \mathbf{h}_{t-1}$ '. As we can see, using longer prediction lengths dramatically improves long-term. Overall, the best performance is obtained using two subsequences of length  $T = 15$ .

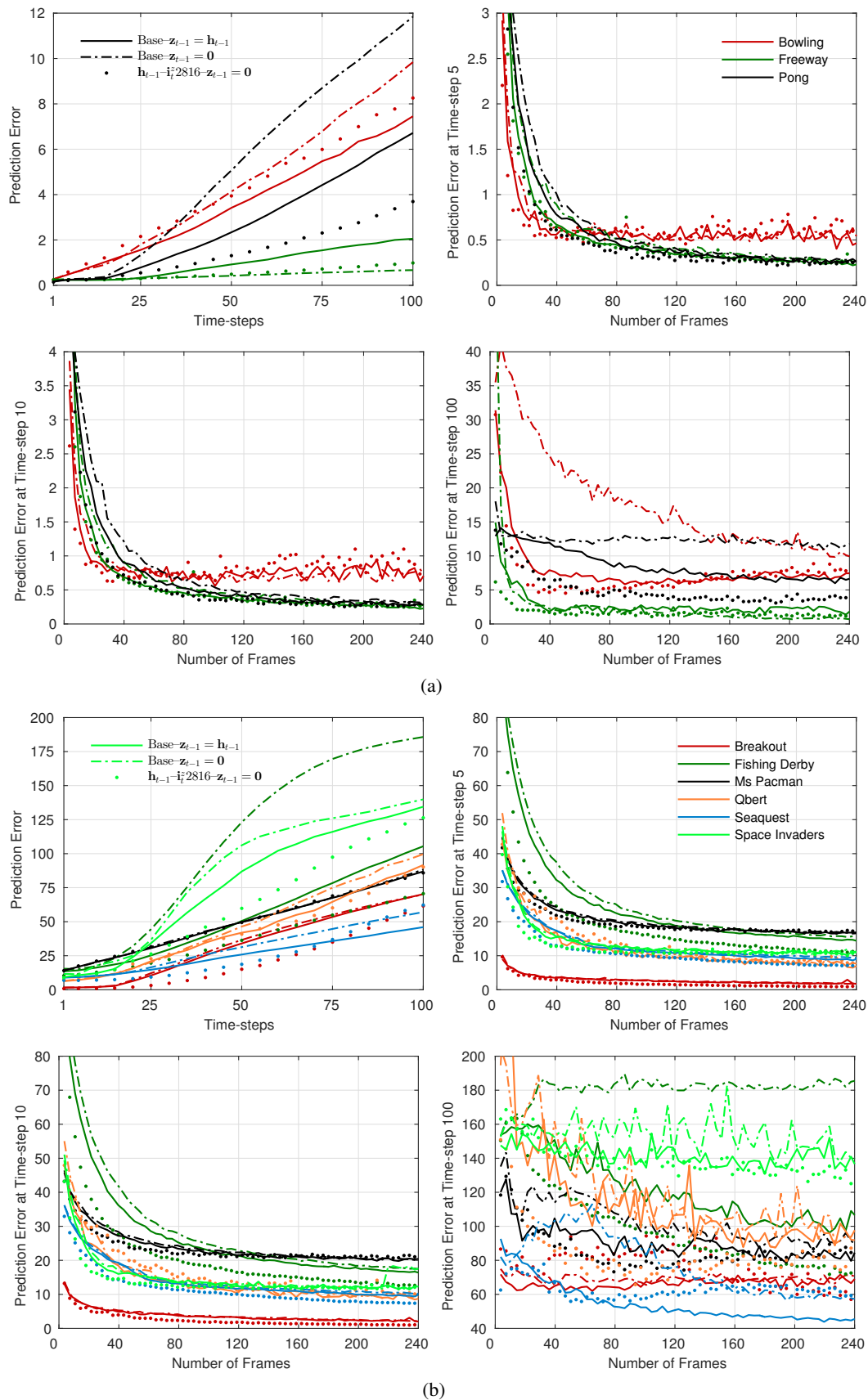
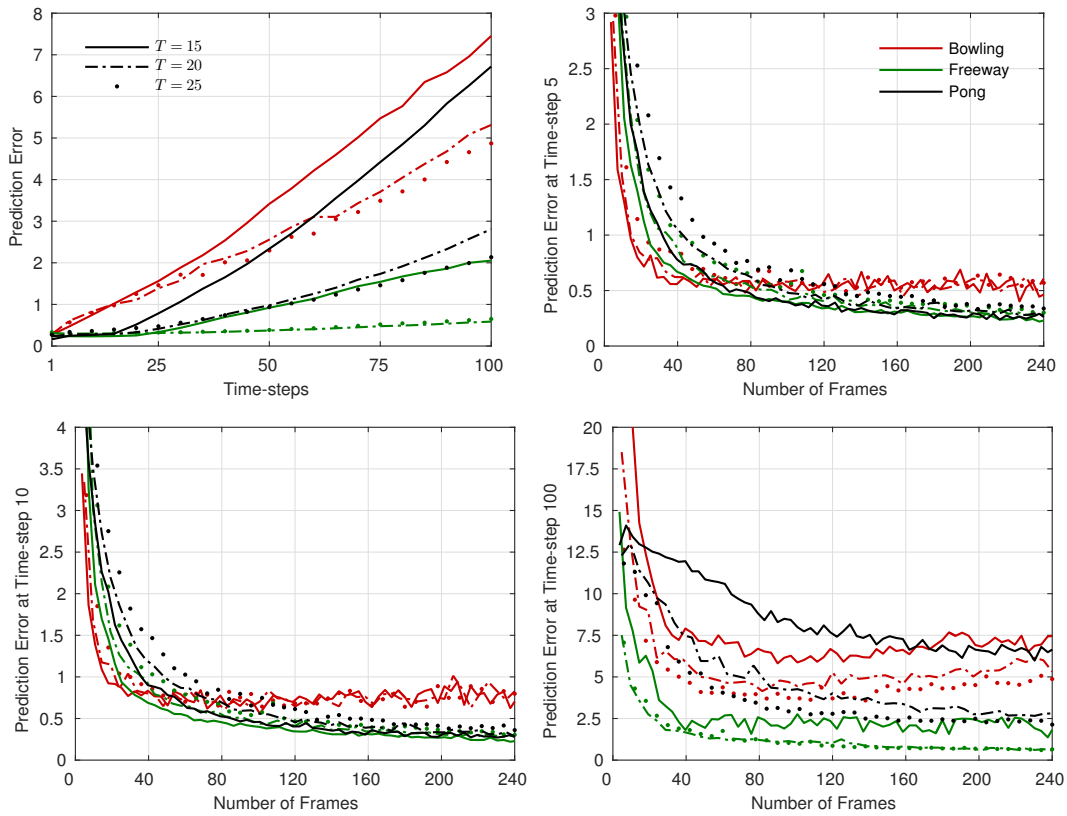
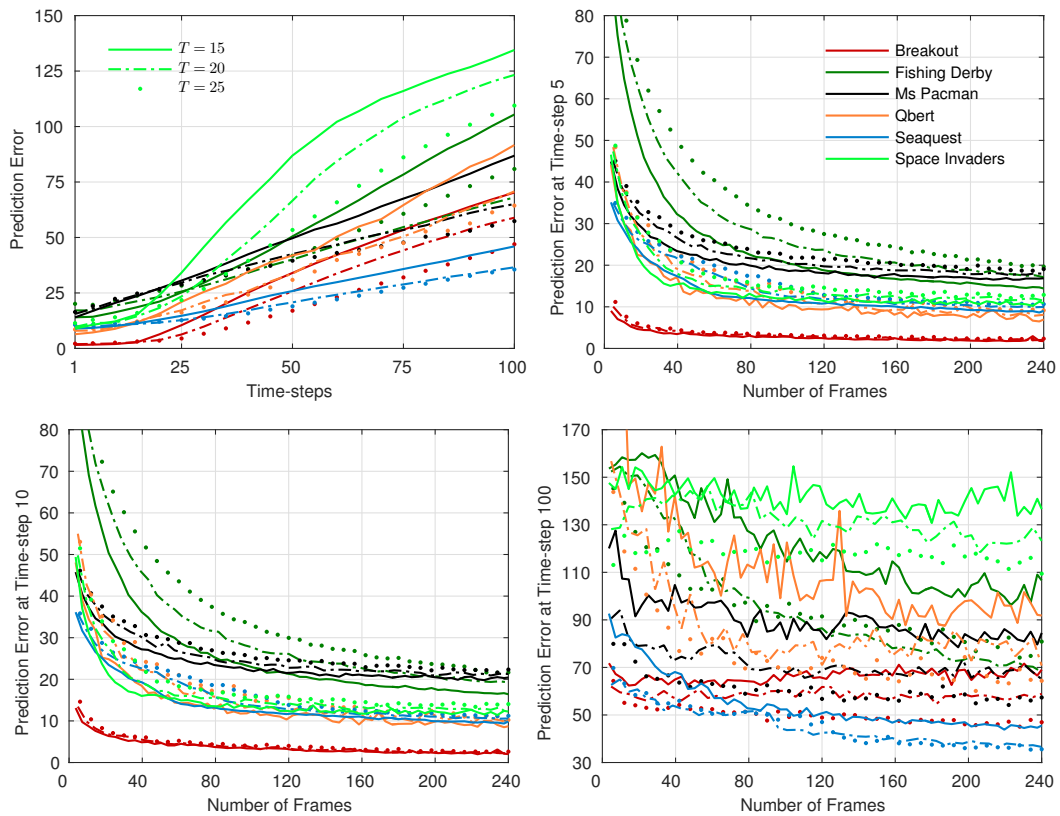


Figure 51: Prediction error (average over 10,000 sequences) for the prediction-independent simulator with different action-dependent state transitions for (a) Bowling, Freeway, Pong, and (b) Breakout, Fishing Derby, Ms Pacman, Qbert, Seaquest, Space Invaders. Number of frames is in million and excludes warm-up frames.



(a)



(b)

Figure 52: Prediction error for the prediction-independent simulator with different prediction lengths  $T \leq 25$  for (a) Bowling, Freeway, Pong, and (b) Breakout, Fishing Derby, Ms Pacman, Qbert, Seaquest, Space Invaders.

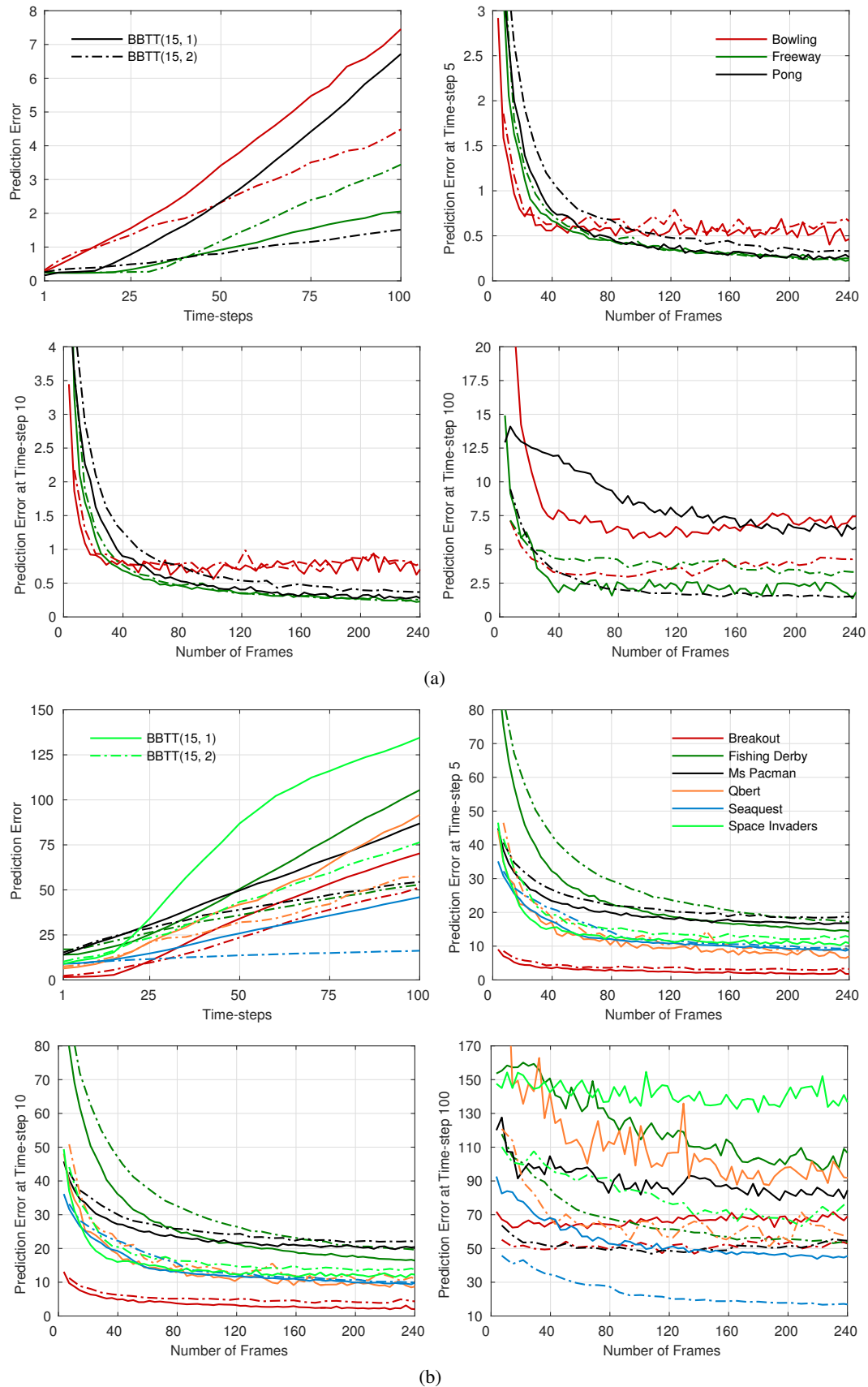


Figure 53: Prediction error for the prediction-independent simulator with BPTT(15, 1) and BPTT(15, 2) for (a) Bowling, Freeway, Pong, and (b) Breakout, Fishing Derby, Ms Pacman, Qbert, Seaquest, Space Invaders.