

---

# Data-Efficient Learning of Feedback Policies from Image Pixels using Deep Dynamical Models

---

**John-Alexander M. Assael\***  
 Department of Computing  
 Imperial College London, UK  
 i.assael@imperial.ac.uk

**Niklas Wahlström\***  
 Division of Automatic Control  
 Linköping University, Sweden  
 nikwa@isy.liu.se

**Thomas B. Schön**  
 Department of Information Technology  
 Uppsala University, Sweden  
 thomas.schon@it.uu.se

**Marc Peter Deisenroth**  
 Department of Computing  
 Imperial College London, UK  
 m.deisenroth@imperial.ac.uk

## Abstract

Data-efficient reinforcement learning (RL) in continuous state-action spaces using very high-dimensional observations remains a key challenge in developing fully autonomous systems. We consider a particularly important instance of this challenge, the pixels-to-torques problem, where an RL agent learns a closed-loop control policy (“torques”) from pixel information only. We introduce a data-efficient, model-based reinforcement learning algorithm that learns such a closed-loop policy directly from pixel information. The key ingredient is a deep dynamical model for learning a low-dimensional feature embedding of images jointly with a predictive model in this low-dimensional feature space. Joint learning is crucial for long-term predictions, which lie at the core of the adaptive nonlinear model predictive control strategy that we use for closed-loop control. Compared to state-of-the-art RL methods for continuous states and actions, our approach learns quickly, scales to high-dimensional state spaces, is lightweight and an important step toward fully autonomous end-to-end learning from pixels to torques.

## 1 Introduction

The vision of fully autonomous and intelligent systems that learn by themselves has inspired artificial intelligence (AI) and robotics research for many decades. The *pixels to torques problem* identifies key aspects of such an autonomous system: autonomous thinking and decision making using (general-purpose) sensor measurements only, intelligent exploration and learning from mistakes. We consider the problem of efficiently learning closed-loop policies (“torques”) from pixel information end-to-end. Although, this problem falls into the general class of reinforcement learning (RL) [1], it is challenging for the following reasons: (1) The state-space (here defined by pixel values) is enormous (e.g., for a  $50 \times 50$  image, we are looking at 2500 continuous-valued dimensions); (2) In many practical applications, we need to find solutions data efficiently: When working with real systems, e.g., robots, we cannot perform millions of experiments because of time and hardware constraints.

One way of using data efficiently, and, therefore, reducing the number of experiments, is to learn predictive forward models of the underlying dynamical system, which are then used for internal simulations and policy learning. These ideas have been successfully applied to RL, control and

---

\*These authors contributed equally to this work.

robotics [2–7], for instance. However, they often rely on heuristics, demonstrations or engineered low-dimensional features, and do not easily scale to data-efficient RL using pixel information only.

A common way of dealing with high-dimensional data is to learn low-dimensional feature representations. Deep learning architectures, such as deep neural networks [8], stacked auto-encoders [9, 10], or convolutional neural networks [11], are the current state-of-the-art in learning parsimonious representations of high-dimensional data. Since 2006, deep learning has produced outstanding empirical results in image, text and audio tasks [12].

**Related Work** Over the last months, there has been significant progress in the context of the pixels-to-torques problem. A first working solution was presented in 2015 [13], where an RL agent automatically learned to play Atari games purely based on pixel information. The key idea was to embed the high-dimensional pixel space into a lower-dimensional space using deep neural networks and apply Q-learning in this compact feature space. A potential issue with this approach is that it is not a data-efficient way of learning policies (weeks of training data are required), i.e., it will be impractical to apply it to a robotic scenario. This data inefficiency is not specific to Q-learning, but a general problem of model-free RL methods [3, 14].

To increase data efficiency, model-based RL methods aim to learn a model of the transition dynamics of the system/robot and subsequently use this model as a surrogate simulator. Recently, the idea of learning predictive models from raw images where only pixel information is available was exploited [15, 16]. The approach taken here follows the idea of Deep Dynamical Models (DDMs) [17]: Instead of learning predictive models for images directly, a detour via a low-dimensional feature space is taken by embedding images into a lower-dimensional feature space, e.g., with a deep auto-encoder. This detour is promising since direct mappings between high-dimensional spaces require large data sets. Whereas Wahlström et al. [15] consider deterministic systems and nonlinear model predictive control (NMPC) techniques for online control, Watter et al. [16] use variational auto-encoders [18], local linearization, and locally linear control methods (iLQR [19] and AICO [20]).

To model the dynamical behavior of the system, the pixels of both the current and previous frame are used. Watter et al. [16] concatenate the input pixels to discover such features, whereas Wahlström et al. [15] concatenate the processed low-dimensional embeddings of the two states. The latter approach requires at least  $\approx 4\times$  fewer parameters, which makes it a promising candidate for more data-efficient learning. Nevertheless, properties such as local linearization [16] can be attractive. However, the complex architecture proposed by Watter et al. [16] is based on very large neural network models with  $\approx 6$  million parameters for learning to swing up a single pendulum. A vast number of training parameters results in higher model complexity, and, thus, decreased statistical efficiency. Hence, an excessive number of training samples, which might not be available, is required to learn the underlying system, taking several days to be trained. These properties make data-efficient learning complicated. Therefore, we propose a relatively lightweight architecture to address the pixels-to-torques problem in a data-efficient manner.

**Contribution** We propose a data-efficient model-based RL algorithm that addresses the pixels-to-torques problem. (1) We devise a data-efficient policy learning framework based on the DDM approach for learning predictive models for images. We use state-of-the-art optimization techniques for training the DDM. (2) Our model profits from a concatenation of low-dimensional features (instead of high-dimensional images) to model dynamical behavior, yielding  $\approx 4$ –20 times fewer model parameters and faster training time than the complex E2C architecture [16]. In practice, our model requires only a few hours of training, while E2C [16] requires days. (3) We introduce a novel training objective that encourages consistency in the latent space paving the way towards more accurate long-term predictions. Overall, we use an efficient model architecture, which can learn tasks of complex non-linear dynamics.

## 2 Problem Set-up and Objective

We consider an  $N$ -step finite-horizon RL setting in which an agent attempts to solve a particular task by trial and error. In particular, our objective is to find a closed-loop policy  $\pi^*$ , that minimizes the long-term cost  $V^\pi = \sum_{t=0}^{N-1} c(\mathbf{s}_t, \mathbf{u}_t)$ , where  $c$  denotes an immediate cost,  $\mathbf{s}_t \in \mathbb{R}^{N_s}$  is the continuous-valued system state, and  $\mathbf{u}_t \in \mathbb{R}^{N_u}$  are continuous control signals.

The learning agent faces the following additional two challenges: (a) The agent has no access to the true state  $\mathbf{s}_t$ , but perceives the environment only through high-dimensional pixel information  $\mathbf{x}_t \in \mathbb{R}^{N_p}$  (images); (b) A good control policy is required in only a few trials. This setting is practically relevant, e.g., when the agent is a robot that is monitored by a video camera based on which the robot has to learn to solve tasks fully autonomously. Therefore, this setting is an instance of the pixels-to-torques problem.

We solve this problem in three key steps, which will be detailed in the following sections: (a) Using a deep auto-encoder architecture we map the high-dimensional pixel information  $\mathbf{x}_t$  to a low-dimensional embedding/feature  $\mathbf{z}_t$ . (b) We combine the  $\mathbf{z}_t$ , and  $\mathbf{z}_{t-1}$  features with the control signal  $\mathbf{u}_t$  to learn a predictive model of the system dynamics for predicting future features  $\mathbf{z}_{t+1}$ . (a) and (b) form a Deep Dynamical Model (DDM) [17]. (c) We apply an adaptive nonlinear model predictive control strategy for optimal closed-loop control and end-to-end learning from pixels to torques.

### 3 Learning a Deep Dynamical Model (DDM)

Our approach to solving the pixels-to-torques problem is based on a deep dynamical model (DDM), see Figure 1, which jointly (a) embeds high-dimensional images in a low-dimensional feature space via deep auto-encoders, and (b) learns a predictive forward model in this feature space, based on the work by Wahlström et al. [17]. In particular, we consider a DDM with control signals  $\mathbf{u}_t$  and high-dimensional observations  $\mathbf{x}_t$  at time-step  $t$ . We assume that the relevant properties of  $\mathbf{x}_t$  can be compactly represented by a feature variable  $\mathbf{z}_t$ . Furthermore,  $\tilde{\mathbf{x}}_t$  is the reconstructed high-dimensional measurement. The two components of the DDM, i.e., the low-dimensional feature and the predictive model, which predicts future features  $\hat{\mathbf{z}}_{t+1}$  and observations  $\hat{\mathbf{x}}_{t+1}$  based on past observations and control signals, are detailed in the following sections.

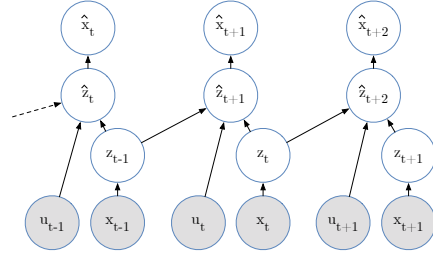


Figure 1: Graphical model of the DDM process for predicting future states.

#### 3.1 Predictive Forward Model

Inspired by the concept of static auto-encoders [21, 22], we turn them into a dynamical model that can predict future features  $\hat{\mathbf{z}}_{t+1}$  and images  $\hat{\mathbf{x}}_{t+1}$ . Our DDM consists of the following elements:

1. An encoder  $f_{\text{enc}}$  mapping high-dimensional observations  $\mathbf{x}_t$  onto low-dimensional features  $\mathbf{z}_t$ ,
2. A decoder  $f_{\text{dec}}$  mapping low-dimensional features  $\mathbf{z}_t$  back to high-dimensional observations  $\tilde{\mathbf{x}}_t$ , and
3. The predictive model  $f_{\text{pred}}$ , which takes  $\mathbf{z}_t, \mathbf{z}_{t-1}, \mathbf{u}_t$  as input and predicts the next latent feature  $\hat{\mathbf{z}}_{t+1}$ .

The  $f_{\text{enc}}, f_{\text{dec}}$  and  $f_{\text{pred}}$  functions of our DDM are neural network models performing the following transformations:

$$\mathbf{z}_t = f_{\text{enc}}(\mathbf{x}_t), \quad (1a)$$

$$\tilde{\mathbf{x}}_t = f_{\text{dec}}(\mathbf{z}_t), \quad (1b)$$

$$\hat{\mathbf{z}}_{t+1} = f_{\text{pred}}(\mathbf{z}_{t-1}, \mathbf{z}_t, \mathbf{u}_t), \quad (1c)$$

$$\hat{\mathbf{x}}_{t+1} = f_{\text{dec}}(\hat{\mathbf{z}}_{t+1}). \quad (1d)$$

We now put these elements together to construct the DDM. The DDM architecture takes the raw images  $\mathbf{x}_{t-1}$  and  $\mathbf{x}_t$  as input and maps them to their low-dimensional features  $\mathbf{z}_{t-1}$  and  $\mathbf{z}_t$  respectively, using  $f_{\text{enc}}$  in (1a). These latent features are then concatenated and, together with the control signal  $\mathbf{u}_t$ , used to predict  $\hat{\mathbf{z}}_{t+1}$  with  $f_{\text{pred}}$  in (1c). Finally, the predicted feature  $\hat{\mathbf{z}}_{t+1}$  is passed through the decoder network  $f_{\text{dec}}$ , to compute the predicted image  $\hat{\mathbf{x}}_{t+1}$ . The overall architecture is depicted in Figure 2.

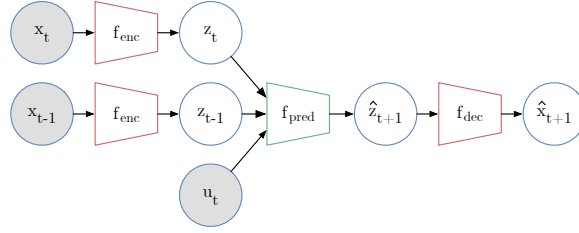


Figure 2: The architecture of the proposed DDM, for extracting the underlying properties of the dynamical system and predicting the next image frame from raw pixel information.

The neural networks  $f_{\text{enc}}$ ,  $f_{\text{dec}}$  and  $f_{\text{pred}}$  that compose the DDM, are parameterized by  $\theta_{\text{enc}}$ ,  $\theta_{\text{dec}}$  and  $\theta_{\text{pred}}$  respectively. These parameters consist of the weights that perform linear transformations of the input data in each neuron.

### 3.2 Training

For training the DDM in (1), we introduce a novel training objective, that encourages consistency in the latent space, paving the way toward accurate long-term predictions. More specifically, for our training objective we define the following cost functions

$$\mathcal{L}_R(\mathbf{x}_t) = \|\tilde{\mathbf{x}}_t - \mathbf{x}_t\|^2, \quad (2a)$$

$$\mathcal{L}_P(\mathbf{x}_{t-1}, \mathbf{x}_t, \mathbf{u}_t, \mathbf{x}_{t+1}) = \|\hat{\mathbf{x}}_{t+1} - \mathbf{x}_{t+1}\|^2, \quad (2b)$$

$$\mathcal{L}_L(\mathbf{z}_{t-1}, \mathbf{z}_t, \mathbf{u}_t, \mathbf{z}_{t+1}) = \|\hat{\mathbf{z}}_{t+1} - \mathbf{z}_{t+1}\|^2, \quad (2c)$$

where  $\mathcal{L}_R$  is the squared deep auto-encoder reconstruction error and  $\mathcal{L}_P$  is the squared prediction error, both operating in image space. Note that  $\hat{\mathbf{x}}_{t+1} = f_{\text{dec}}(f_{\text{pred}}(f_{\text{enc}}(\mathbf{x}_t), f_{\text{enc}}(\mathbf{x}_t), \mathbf{u}_t))$  depends on the parameters of the decoder, the predictive model and the encoder. Additionally, we introduce  $\mathcal{L}_L$  that enforces consistency between the latent spaces of the encoder  $f_{\text{enc}}$  and the prediction model  $f_{\text{pred}}$ . In the big-data regime, this additional penalty in latent space is not necessary, but if not much data is available, this additional term increases the data efficiency as the prediction model is forced to make predictions  $\hat{\mathbf{z}}_{t+1} = f_{\text{pred}}(\mathbf{z}_{t-1}, \mathbf{z}_t, \mathbf{u}_t)$  close to the next embedded feature  $\mathbf{z}_{t+1} = f_{\text{enc}}(\mathbf{x}_{t+1})$ . The overall training objective of the current dataset  $\mathcal{D} = (\mathbf{x}_{0:N}, \mathbf{u}_{0:N})$  is given by

$$\begin{aligned} \mathcal{L}(\mathcal{D}) &= \sum_{t=0}^{N-1} \mathcal{L}_R(\mathbf{x}_t) + \mathcal{L}_P(\mathbf{x}_{t-1}, \mathbf{x}_t, \mathbf{u}_t, \mathbf{x}_{t+1}) + \alpha \mathcal{L}_L(\mathbf{z}_{t-1}, \mathbf{z}_t, \mathbf{u}_t, \mathbf{z}_{t+1}) \\ &= \sum_{t=0}^{N-1} \|\tilde{\mathbf{x}}_t - \mathbf{x}_t\|^2 + \|\hat{\mathbf{x}}_{t+1} - \mathbf{x}_{t+1}\|^2 + \alpha \|\hat{\mathbf{z}}_{t+1} - \mathbf{z}_{t+1}\|^2, \end{aligned} \quad (3)$$

where  $\alpha$  is a parameter that controls the influence of  $\mathcal{L}_L$ . Finally, we train the DDM parameters  $(\theta_{\text{enc}}, \theta_{\text{dec}}, \theta_{\text{pred}})$  by jointly minimizing the overall cost

$$(\hat{\theta}_{\text{enc}}, \hat{\theta}_{\text{dec}}, \hat{\theta}_{\text{pred}}) = \arg \min_{\theta_{\text{enc}}, \theta_{\text{dec}}, \theta_{\text{pred}}} \mathcal{L}(\mathcal{D}). \quad (4)$$

Training jointly leads to good predictions as it facilitates the extraction and separation of the features describing the underlying dynamical system, and not only features for creating good reconstructions [17].

### 3.3 Network Architecture

The neural networks  $f_{\text{enc}}$ ,  $f_{\text{dec}}$  and  $f_{\text{pred}}$  are composed by 3 linear layers, where each of the first 2 are followed by Rectified Linear Unit (ReLU) activation functions [23]. As it has been demonstrated [24], ReLU non-linearities allow the network to train  $\approx 6 \times$  faster than the conventional tanh units, as evaluated on the CIFAR-10 [25] dataset. Furthermore, similar to Watter et al. [16], we use Adam [26] to train the DDM, which is considered the state-of-the-art among the latest methods for stochastic gradient optimization. Finally, after evaluating different weight optimization methods, such as uniform and random Gaussian [24, 27], the weights of the DDM were initialized using orthogonal weight initialization [28], which demonstrated the most efficient training performance, leading to decoupled weights that evolve independently of each other.

## 4 Policy Learning

Our objective is to control the system to a state where a certain target frame  $\mathbf{x}_{\text{ref}}$  without any prior knowledge of the system at hand. To accomplish this we use the DDM for learning a closed-loop policy by means of nonlinear model predictive control (NMPC).

### 4.1 NMPC using the DDM

NMPC finds an optimal sequence of control signals that minimizes a  $K$ -step loss function, where  $K$  is typically smaller than the full horizon. We choose to do the control in the low-dimensional embedded space to reduce the complexity of the control problem.

Our NMPC formulation relies on (a) a target feature  $\mathbf{z}_{\text{ref}}$  and (b) the DDM that allows us to predict future features. The target feature is computed by encoding the target frame  $\mathbf{z}_{\text{ref}} = f_{\text{enc}}(\mathbf{x}_{\text{ref}}, \theta_E)$  provided by the model. Further, with the DDM, future features  $\hat{\mathbf{z}}_1, \dots, \hat{\mathbf{z}}_K$  can be predicted based on a sequence of future (and yet unknown) controls  $\mathbf{u}_0, \dots, \mathbf{u}_{K-1}$  and two initial encoded features  $\mathbf{z}_{-1}, \mathbf{z}_0$  assuming that the current feature is denoted by  $\mathbf{z}_0$ .

Using the dynamical model, NMPC determines an optimal (open-loop) control sequence  $\mathbf{u}_0^*, \dots, \mathbf{u}_{K-1}^*$ , such that the predicted features  $\mathbf{z}_1, \dots, \hat{\mathbf{z}}_K$  gets as close to the target feature  $\mathbf{z}_{\text{ref}}$  as possible, which results in the objective

$$\mathbf{u}_0^*, \dots, \mathbf{u}_{K-1}^* \in \arg \min_{\mathbf{u}_{0:K-1}} \sum_{t=0}^{K-1} \|\hat{\mathbf{z}}_t - \mathbf{z}_{\text{ref}}\|^2 + \lambda \|\mathbf{u}_t\|^2, \quad (5)$$

where  $\|\hat{\mathbf{z}}_t - \mathbf{z}_{\text{ref}}\|^2$  is a cost associated with the deviation of the predicted features  $\hat{\mathbf{z}}_{0:K-1}$  from the reference feature  $\mathbf{z}_{\text{ref}}$ , and  $\|\mathbf{u}_t\|^2$  penalizes the amplitude of the control signals. Here,  $\lambda$  is a tuning parameter adjusting the importance of these two objectives. When the control sequence  $\mathbf{u}_0^*, \dots, \mathbf{u}_{K-1}^*$  is determined, the first control  $\mathbf{u}_0^*$  is applied to the system. After observing the next feature, NMPC repeats the entire optimization and turns the overall policy into a closed-loop (feedback) control strategy.

Overall, we now have an online NMPC algorithm that, given a trained DDM, works indirectly on images by exploiting their feature representation.

### 4.2 Adaptive NMPC for Learning from Scratch

We will now turn over to describe how adaptive NMPC can be used together with our DDM to address the pixels-to-torques problem and to learn from scratch. At the core of our NMPC formulation lies the DDM, which is used to predict future features (and images) from a sequence of control signals. The quality of the NMPC controller is inherently bound to the prediction quality of the dynamical model, which is typical in model-based RL [5, 14, 29].

To learn models and controllers from scratch, we apply a control scheme that allows us to update the DDM as new data arrives. In particular, we use the NMPC controller in an adaptive fashion to gradually improve the model by collected data in the feedback loop without any specific prior knowledge of the system at hand. Data collection is performed in closed-loop (online NMPC), and it is divided into multiple sequential trials. After each trial, we add the data

of the most recent trajectory to the data set, and the model is re-trained using all data that has been collected so far.

Simply applying the NMPC controller based on a randomly initialized model would make the closed-loop system very likely to converge to a point, which is far away from the desired refer-

---

#### Algorithm 1 Adaptive online NMPC in feature space

---

Follow a random control strategy and record data

**loop**

Update DDM with all data collected so far

**for**  $t = 0$  to  $N - 1$  **do**

Get current feature  $\mathbf{z}_t$  via the encoder

$\mathbf{u}_t^* \leftarrow \epsilon$ -greedy NMPC policy using DDM pred.

Apply  $\mathbf{u}_t^*$  and record data

**end for**

**end loop**

---

ence value, due to the poor model that cannot extrapolate well to unseen states. This would in turn imply that no data is collected in unexplored regions, including the region that we are interested in. There are two solutions to this problem: either we use a probabilistic dynamical model [5, 14] to explicitly account for model uncertainty and the implied natural exploration, or we follow an explicit exploration strategy to ensure proper excitation of the system. In this paper, we follow the latter approach. In particular, we choose an  $\epsilon$ -greedy exploration strategy where the optimal feedback  $\mathbf{u}_0^*$  at each time step is selected with a probability  $1 - \epsilon$ , and a random action is selected with probability  $\epsilon$ .

Algorithm 1 summarizes our adaptive online NMPC scheme. We initialize the DDM with a random trial. We use the learned DDM to find an  $\epsilon$ -greedy policy using predicted features within NMPC. This happens online while the collected data is added to the data set, and the DDM is updated after each trial.

## 5 Experimental Evaluation

In this section, we empirically assess the components of our proposed methodology for autonomous learning from high-dimensional synthetic image data, on learning the underlying dynamics of a single and a planar double pendulum. The main lines of the evaluation are: (a) the quality of the learned DDM and (b) the overall learning framework.

In both experiments, we consider the following setting: We take screenshots of a simulated pendulum system at a sampling frequency of 0.2 s. Each pixel  $\mathbf{x}_t^{(i)}$  is a component of the measurement  $\mathbf{x}_t \in \mathbb{R}^{N^*}$  and takes a continuous gray-value in the interval  $[0, 1]$ . The control signals  $\mathbf{u}_t$  are the torques applied to the system. No access to the underlying dynamics nor the state (angles and angular velocities) was available, i.e., we are dealing with a high-dimensional continuous time series. The challenge was to data-efficiently learn (a) a good dynamical model and (b) a good controller from pixel information only.

To speed up the training process, we applied PCA prior to model learning as a pre-processing step to reduce the dimensionality of the original problem. With these inputs, a 3-layer auto-encoder was employed, such that the dimensionality of the features is optimal to model the periodic angle of the pendulums. The features  $\mathbf{z}_{t-1}, \mathbf{z}_t$  and  $\mathbf{u}_t$  were later passed to the 3-layer predictive feedforward neural network generating  $\hat{\mathbf{z}}_{t+1}$ . Furthermore, during training, the  $\alpha$  parameter for encouraging consistent latent space predictions was set to 1 for both experiments. While, in the adaptive NMPC, the  $\lambda$  tuning parameter that penalizes the amplitude of the control signals, was set to 0.01.

### 5.1 Planar Pendulum

The first experiment evaluates the performance of the DDM on a planar pendulum, assembled by 1-link robot arm with length 1 m, weight 1 kg and friction coefficient 1 N s m/rad. The screenshots consist of  $40 \times 40 = 1600$  pixels, and the input dimension has been reduced to  $\dim(\mathbf{x}_t) = 100$  using PCA. These inputs are processed by an encoder  $f_{\text{enc}}$  with architecture:  $100 \times 50 - \text{ReLU} - 50 \times 50 - \text{ReLU} - 50 \times 2$ .

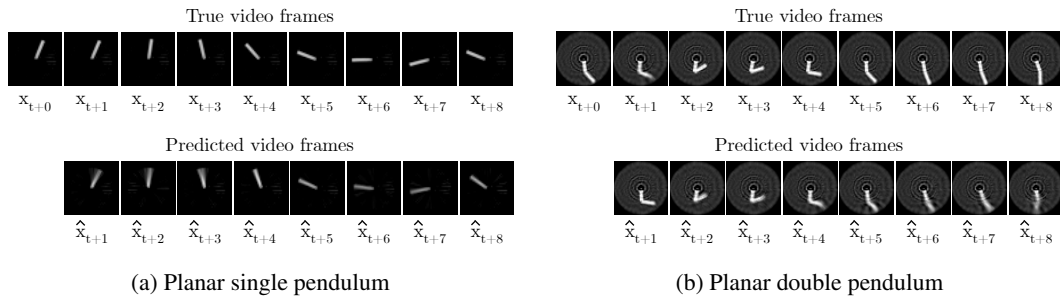


Figure 3: Long-term (up to eight steps) predictive performance of the DDM controlling a planar pendulum (a) and a planar double pendulum (b): true (upper plot) and predicted (lower plot) video frames on test data.

The low-dimensional features are of  $\dim(\mathbf{z}_t) = 2$  in order to model the periodic angle of the pendulum. To capture the *dynamic* properties, such as angular velocity, we concatenate two consecutive features  $\mathbf{z}_{t-1}, \mathbf{z}_t$  with the control signal  $\mathbf{u}_t \in \mathbb{R}^1$  and pass them through the predictive model  $f_{\text{pred}}$ , with architecture:  $5 \times 100 - \text{ReLU} - 100 \times 100 - \text{ReLU} - 100 \times 2$ . Note that the dimensionality of the first layer is given by  $\dim(\mathbf{z}_{t-1}) + \dim(\mathbf{z}_t) + \dim(\mathbf{u}_t) = 5$ . Finally, the predicted feature  $\hat{\mathbf{z}}_{t+1}$ , can be mapped back to  $\hat{\mathbf{x}}_{t+1}$  using our decoder, with architecture:  $2 \times 50 - \text{ReLU} - 50 \times 50 - \text{ReLU} - 50 \times 100$ .

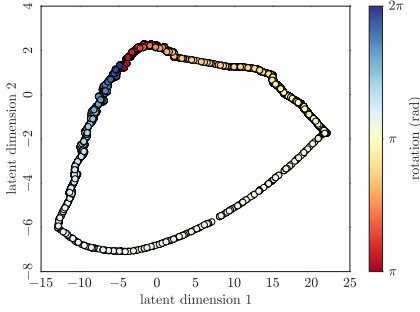


Figure 4: The feature space  $\mathbf{z} \in \mathbb{R}^2$  of the single pendulum experiment for different pendulum angles between  $0^\circ$  and  $360^\circ$ , generated by the DDM.

The performance of the DDM is illustrated in Figure 3 on a test data set. The top row shows the true images and the bottom row shows the DDM’s long-term predictions. The model yields a good predictive performance for both one-step ahead prediction and multiple-step ahead prediction, a consequence of (a) jointly learning predictor and auto-encoder, (b) concatenating features instead of images to model the dynamic behavior.

In Figure 4, we show the learned feature space  $\mathbf{z} \in \mathbb{R}^2$  for different pendulum angles between  $0^\circ$  and  $360^\circ$ . The DDM has learned to generate features that represent the angle of the pendulum, as they are mapped to a circle-like shape accounting for the wrap-around property of an angle.

Finally, in Figure 5, we report results on learning a policy that moves the pendulum from a start position  $\varphi = 0^\circ$  to an upright target position  $\varphi = \pm 180^\circ$ . The reference signal was the screenshot of the pendulum in the target position. For the NMPC controller, we used a planning horizon of  $K = 15$  steps and a control penalty  $\lambda = 0.01$ . For the  $\epsilon$ -greedy exploration strategy we used  $\epsilon = 0.2$ .

Figure 5 shows the learning stages of the system, i.e., the 18 different trials of the NMPC controller. Starting with a randomly initialized model, 500 images were appended to the dataset in each trial. As it can be seen, starting already from the first controlled trial, the system managed to control the pendulum successfully and bring it to a position less than  $10^\circ$  from the target position. This means, the solution is found very data efficiently, especially when we consider that the problem is learned from pixels information without access to the “true” state.

## 5.2 Planar Double Pendulum

In this experiment, a planar double pendulum is considered assembled by 2-link robot arm with length 1m and 1m respectively, weight 1kg and 1kg and friction coefficients 1 N s/m/rad. Torques can be applied at both joints. The screenshots consist of  $48 \times 48 = 2304$  pixels, and the input dimension has been reduced to  $\dim(\mathbf{x}_t) = 512$  prior to model learning using PCA to speed up the

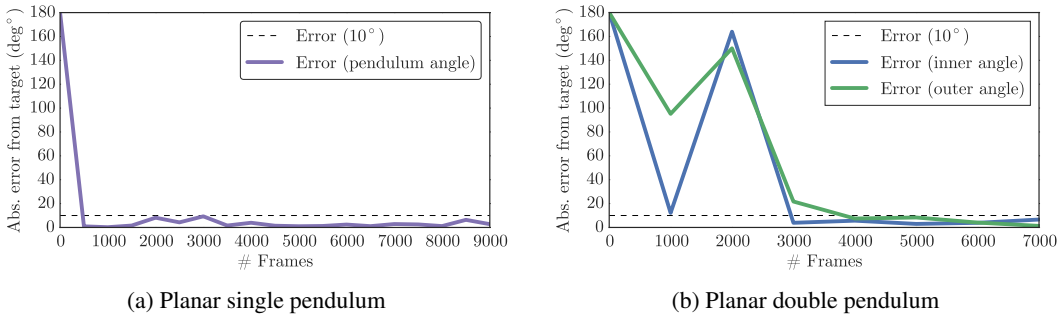


Figure 5: Results on learning a policy that moves the single (a) and the double (b) pendulum systems from  $\varphi = 0^\circ$  to  $\varphi = \pm 180^\circ$ , in 100 time-steps. The horizontal axis shows the learning stages and the corresponding image frames available to the learner. The vertical axis shows the absolute error from the target state, averaged over the last 10 time steps of each test trajectory. The dashed line shows a  $10^\circ$  error, which indicates a “good” solution.

training process. The encoder architecture is:  $512 \times 256 - \text{ReLU} - 256 \times 256 - \text{ReLU} - 256 \times 4$ , and the decoder vice versa. The low-dimensional embeddings  $\dim(\mathbf{z}_t) = 4$  and the architecture of the predictive model was:  $10 \times 200 - \text{ReLU} - 200 \times 200 - \text{ReLU} - 200 \times 4$ .

The predictive performance of the DDM is shown in Figure 3 on a test data set. The performance of the controller is depicted in Figure 5. We used 7 trials with the downward initial position  $\varphi = 0^\circ$  and upward target  $\varphi = \pm 180^\circ$  for the angle of both inner and outer pendulums. The figure shows the error after each trial (1000 frames) and clearly indicates that after three controlled trials a good solution is found, which brings both pendulums within a  $10^\circ$  range to the target angles.

Despite the high complexity of the dynamical system, our learning framework manages to successfully control both pendulums after the third trial in nearly all cases.

### 5.3 Comparison with State-of-the-Art

The same experiments were executed employing PILCO [30], a state of the art policy search method, under the following settings: (a) PILCO has access to the true state, i.e., the angle  $\varphi$  and angular velocity  $\dot{\varphi}$ ; (b) A deep auto-encoder is used to learn two-dimensional features  $\mathbf{z}_t$  from images, which are used by PILCO for policy learning. In the first setting (a) PILCO managed to successfully reach the target after the second and the third trial in the two experiments, respectively. However, in setting (b), PILCO did not manage to learn anything meaningful at all. The reason why PILCO could not learn on auto-encoder features is that these features were only trained to minimize the reconstruction error. However, the auto-encoder did not attempt to map similar images to similar features, which led to zig-zagging around in feature space (instead of following a smooth manifold as in Figure 4), making the model learning part in feature space incredibly hard [17].

We modeled and controlled equally complex models with E2C [16], but at the same time our DDM requires  $\approx 4\times$  fewer neural network parameters if we use the same PCA pre-processing step within E2C. The reason lies in our efficient processing of the dynamics of the model in the feature space instead of the image space. This number increases up to  $\approx 20\times$  fewer parameters than E2C without the PCA pre-processing step.

The number of parameters can be directly translated to reduced training time, and increased data efficiency. Employing the adaptive model predictive control, our proposed DDM model requires significantly less data samples, as it efficiently focuses on learning the latent space towards the reference target state. Furthermore, the control performance of our model is gradually improved in respect to the number of trials. As proved by our experimental evaluation we can successfully control a complex dynamical system, such as the planar double pendulum, with less than 4 000 samples. This adaptive learning approach can be essential in problems with time and hardware constraints.

## 6 Conclusion

We proposed a data-efficient model-based RL algorithm that learns closed-loop policies in continuous state and action spaces directly from pixel information. The key components of our solution are (a) a deep dynamical model (DDM) that is used for long-term predictions via a compact feature space, (b) a novel training objective that encourages consistency in the latent space, paving the way toward more accurate long-term predictions, and (c) an NMPC controller that uses the predictions of the DDM to determine optimal actions on the fly without the need for value function estimation. For the success of this RL algorithm it is crucial that the DDM learns the feature mapping and the predictive model in feature space jointly to capture dynamical behavior for high-quality long-term predictions. Compared to state-of-the-art RL our algorithm learns fairly quickly, scales to high-dimensional state spaces and facilitates learning from pixels to torques.

### Acknowledgements

We thank Roberto Calandra for valuable discussions in the early stages of the project. The Tesla K40 used for this research was donated by the NVIDIA Corporation.



## References

- [1] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*, volume 1. MIT Press, 1998.
- [2] J. Schmidhuber. An on-line algorithm for dynamic reinforcement learning and planning in reactive environments. In *International Joint Conference on Neural Networks (IJCNN)*, pages 253–258. IEEE, 1990.
- [3] C. G. Atkeson and S. Schaal. Learning tasks from a single demonstration. In *Proceedings of the International Conference on Robotics and Automation (ICRA)*. IEEE, 1997.
- [4] J. A. Bagnell and J. G. Schneider. Autonomous helicopter control using reinforcement learning policy search methods. In *Proceedings of the IEEE International Conference on Robotics and Automation*, volume 2, pages 1615–1620, 2001.
- [5] M. P. Deisenroth, D. Fox, and C. E. Rasmussen. Gaussian processes for data-efficient learning in robotics and control. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 37(2):408–423, 2015.
- [6] Y. Pan and E. Theodorou. Probabilistic differential dynamic programming. In *Advances in Neural Information Processing Systems*, pages 1907–1915, 2014.
- [7] S. Levine, C. Finn, T. Darrell, and P. Abbeel. End-to-end training of deep visuomotor policies. *arXiv preprint arXiv:1504.00702*, 2015.
- [8] G. Hinton and R. Salakhutdinov. Reducing the dimensionality of data with neural networks. *Science*, 313(5786):504–507, 2006.
- [9] Y. Bengio, P. Lamblin, D. Popovici, and H. Larochelle. Greedy layer-wise training of deep networks. In *Advances in Neural Information Processing Systems (NIPS)*, pages 153–160. MIT Press, 2007.
- [10] P. Vincent, L. Hugo, Y. Bengio, and P.-A. Manzagol. Extracting and composing robust features with denoising autoencoders. In *International Conference on Machine Learning (ICML)*, pages 1096–1103. ACM, 2008. ISBN 978-1-60558-205-4.
- [11] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- [12] J. Schmidhuber. Deep learning in neural networks: An overview. Technical Report IDSIA-03-14 / arXiv:1404.7828v1 [cs.NE], The Swiss AI Lab IDSIA, 2014.
- [13] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, et al. Human-level control through deep reinforcement learning. *Nature*, 518(7540):529–533, 2015.
- [14] J. G. Schneider. Exploiting model uncertainty estimates for safe dynamic control learning. In *Advances in Neural Information Processing Systems (NIPS)*. 1997.
- [15] N. Wahlström, T. B. Schön, and M. P. Deisenroth. From pixels to torques: Policy learning with deep dynamical models. *arXiv preprint arXiv:1502.02251*, 2015.
- [16] M. Watter, J. T. Springenberg, J. Boedecker, and M. A. Riedmiller. Embed to control: A locally linear latent dynamics model for control from raw images. *arXiv preprint arXiv:1506.07365*, 2015.
- [17] N. Wahlström, T. B. Schön, and M. P. Deisenroth. Learning deep dynamical models from image pixels. In *IFAC Symposium on System Identification (SYSID)*, 2015.
- [18] D. Jimenez Rezende, S. Mohamed, and D. Wierstra. Stochastic backpropagation and variational inference in deep latent Gaussian models. In *International Conference on Machine Learning (ICML)*, June 2014.
- [19] E. Todorov and W. Li. A generalized iterative LQG method for locally-optimal feedback control of constrained nonlinear stochastic systems. In *American Control Conference*, pages 300–306. IEEE, 2005.
- [20] M. Toussaint. Robot trajectory optimization using approximate inference. In *International Conference on Machine Learning (ICML)*, Montreal, QC, Canada, June 2009.

- [21] D. E. Rumelhart, G. E. Hinton, and R. J. Williams. Learning internal representations by error propagation. In D. E. Rumelhart and J. L. McClelland, editors, *Parallel Distributed Processing*, volume 1, pages 318–362. MIT Press, 1986.
- [22] Y. Bengio. Learning deep architectures for AI. *Foundations and trends in Machine Learning*, 2(1):1–127, 2009.
- [23] V. Nair and G. E. Hinton. Rectified linear units improve restricted Boltzmann machines. In *International Conference on Machine Learning (ICML)*, pages 807–814, 2010.
- [24] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1097–1105, 2012.
- [25] A. Krizhevsky. Learning multiple layers of features from tiny images. Master’s thesis, Computer Science Department, University of Toronto, 2009.
- [26] D. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.
- [27] X. Glorot and Y. Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Artificial Intelligence and Statistics (AISTATS)*, pages 249–256, 2010.
- [28] A. M. Saxe, J. L. McClelland, and S. Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *arXiv preprint arXiv:1312.6120*, 2013.
- [29] S. Schaal. Learning from demonstration. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1040–1046. 1997.
- [30] M. P. Deisenroth and C. E. Rasmussen. PILCO: A model-based and data-efficient approach to policy search. In *International Conference on Machine Learning (ICML)*, 2011.