
Generative Temporal Models with Spatial Memory for Partially Observed Environments

Marco Fraccaro^{1*} Danilo Jimenez Rezende² Yori Zwols² Alexander Pritzel² S. M. Ali Eslami² Fabio Viola²

Abstract

In model-based reinforcement learning, generative and temporal models of environments can be leveraged to boost agent performance, either by tuning the agent’s representations during training or via use as part of an explicit planning mechanism. However, their application in practice has been limited to simplistic environments, due to the difficulty of training such models in larger, potentially partially-observed and 3D environments. In this work we introduce a novel action-conditioned generative model of such challenging environments. The model features a non-parametric spatial memory system in which we store learned, disentangled representations of the environment. Low-dimensional spatial updates are computed using a state-space model that makes use of knowledge on the prior dynamics of the moving agent, and high-dimensional visual observations are modelled with a Variational Auto-Encoder. The result is a scalable architecture capable of performing coherent predictions over hundreds of time steps across a range of partially observed 2D and 3D environments.

1. Introduction

Consider a setup in which an agent walks and observes an environment (e.g., a three-dimensional maze) for hundreds of time steps, and is then asked to predict subsequent observations given a sequence of actions. This is a challenging task, as it requires the ability to first remember the visual observations and the position in which they were observed in the environment, and secondly to predict where a possibly long sequence of actions would bring the agent in the environment. Building models that can solve this problem

* Work done during an internship at DeepMind. ¹Technical University of Denmark ²DeepMind, London, UK. Correspondence to: Marco Fraccaro <marfraccaro@gmail.com>.

can be useful for model-based reinforcement learning involving spatial tasks that require long-term memories and other spatial downstream goals (Sutton, 1990; Deisenroth & Rasmussen, 2011; Levine & Abbeel, 2014; Watter et al., 2015; Wahlström et al., 2015; Lenz et al., 2015; Higgins et al., 2017; Finn & Levine, 2017). This requires however agents that are able to remember the past over hundreds of steps, that know both where they are in the environment and how each action changes their position, and that can coherently predict hundreds of steps into the future. Therefore the main focus of this work is to develop an action-conditioned generative model that is able to memorize all the required information while exploring the environment and successfully use it in the prediction phase for long-term generation of high-dimensional visual observations.

Recently, several powerful generative models for sequential data have been proposed in a wide range of applications, such as modelling speech, handwriting, polyphonic music and videos (Chung et al., 2015; Fraccaro et al., 2016; Oh et al., 2015; Chiappa et al., 2017). They build on recurrent neural architectures such as Long Short-Term Memory (LSTM) (Hochreiter & Schmidhuber, 1997) or Gated Recurrent Units (GRU) (Chung et al., 2014), that use an internal state vector to perform computations and store the long-term information needed when making predictions. Since the number of parameters in these models scales quadratically with the dimensionality of the state vector, they are not suitable for applications that require high memory capacity, such as the one considered in this paper. The high dimensional state vector needed to be able to memorize the hundreds of time steps in which the agent has visited the whole environment, make these models in practice both very slow and hard to train. An alternative approach is to use an external memory architecture, for storing of large amount of information while drastically reducing the number of parameters with respect to models with similar memory capacity that build on LSTMs or GRUs. Gemici et al. (2017) present a general architecture for generative temporal models with external memory, and test four different types of memories that are dynamically updated at each time step (Graves et al., 2014; 2016; Santoro et al., 2016). They focus on differentiable addressing mechanisms for memory storage and retrieval (soft-attention), that are based on deep neural

networks that learn to write information to the memory and read from it. While this approach is very general and can be used to model complex long-term temporal dependencies in a wide range of applications, it has not been successful in modeling the data coming from an agent freely moving in a 3d maze, even for a single room [private communications with the authors of (Gemici et al., 2017)].

To define a scalable model capable of exploring larger environments and coherently predicting hundreds of time steps in the future, in this work we build a *spatial memory* architecture that exploits some knowledge of the specific structure of the problem in consideration. In particular, at each time step we split the internal latent representation of the system in to two separate vectors, a low-dimensional one that encodes the position of the agent in the environment and a high dimensional one that encodes what the agent is seeing. We model the low dimensional dynamics of the agent with a state-space model in which we encode prior information on the physical principles that govern the agent’s movements, and learn a higher dimensional latent representation of the visual input (the frames from the environment) with a Variational Auto-Encoder (Kingma & Welling, 2014; Rezende et al., 2014). While exploring the environment, at each time step we store the position of the agent and the corresponding visual information in a Differentiable Neural Dictionary (DND) (Pritzel et al., 2017), a scalable non-parametric memory developed for episodic control. The resulting model is able to coherently generate hundreds of time steps into the future in simulated 3D environments, by retrieving at each time step the observations stored in memory that were collected when passing in nearby positions during the exploration phase. Making predictions with our model is scalable because of the efficient rollouts in a low dimensional space made possible by the state-space assumption and the efficient retrieval of the necessary information from DND. The proposed model can be trained end-to-end on videos with corresponding action sequences of agents walking in an environment. Importantly, unlike the work in (Gemici et al., 2017) we do not need to learn a complex memory addressing mechanisms, as in our model the DND represents a non-parametric component where we store encodings of the positions and visual information that are learned from the data in an unsupervised way.

2. Background

We now provide a brief overview of the building blocks for the model introduced in section 3.

Variational auto-encoders. Variational auto-encoders (VAEs) (Kingma & Welling, 2014; Rezende et al., 2014) define a generative model for high-dimensional data \mathbf{x}_t by introducing a latent state \mathbf{z}_t . The joint probability distribution $p_\theta(\mathbf{x}_t, \mathbf{z}_t)$ is factorized as $p_\theta(\mathbf{x}_t, \mathbf{z}_t) = p_\theta(\mathbf{x}_t|\mathbf{z}_t)p(\mathbf{z}_t)$,

where $p(\mathbf{z}_t)$ is the prior of the latent state and the *decoder* $p_\theta(\mathbf{x}_t|\mathbf{z}_t)$ defines a mapping using deep neural networks parameterized by θ from the states \mathbf{z}_t to the data \mathbf{x}_t . In a VAE, the intractable posterior distribution over the latent states is approximated using the variational distribution $q_\phi(\mathbf{z}_t|\mathbf{x}_t)$, also known as the *encoder* or *inference network*. The parameters θ and ϕ of the decoder and the encoder, respectively, are learned jointly by maximizing the Evidence Lower Bound (ELBO) with stochastic gradient ascent.

DND memory. The Differentiable Neural Dictionary (DND) is a scalable, non-parametric memory module first introduced in Reinforcement Learning (RL) to allow agents to store and retrieve their experiences of an environment (Pritzel et al., 2017). The *write* operation consists of inserting (key, value) pairs into the memory; similarly to a dictionary, this associates a value to each key. Given a query key, we can then *read* from the memory by finding among the keys stored in the DND the nearest neighbours to the query key and returning the corresponding values. The DND can be used in applications that require very large memories, since the nearest-neighbour search can be efficiently approximated using space-partitioning data structures, such as kd-trees (Bentley, 1975).

State-space models. State-space models (SSM) are a class of probabilistic graphical models widely used in the temporal setting to model sequences of vectors $\mathbf{z}_{1:T} = [\mathbf{z}_1, \dots, \mathbf{z}_T]$ conditioned on some actions $\mathbf{a}_{1:T} = [\mathbf{a}_1, \dots, \mathbf{a}_T]$. SSMs introduce at each time step a continuous stochastic variable \mathbf{s}_t , used as a latent representation of the state of the system. The temporal dynamics of the system are described by the *transition density* $p(\mathbf{s}_t|\mathbf{s}_{t-1}, \mathbf{a}_t)$ of the SSM, that defines how to update the state at time t given the previous state \mathbf{s}_{t-1} and the current action \mathbf{a}_t . The output variable \mathbf{z}_t depends on the state \mathbf{s}_t through the *emission density* $p(\mathbf{z}_t|\mathbf{s}_t)$.

3. Model

An important component of model-based reinforcement learning is the ability to plan many steps ahead in time leveraging previous experiences (Sutton, 1990; Racanière et al., 2017). This requires agents that can remember the past and use it to predict what may happen in the future given certain actions. With this purpose in mind, we define an action-conditioned generative model with memory, that can be used within RL agents for model-based planning.

The input of our model consists of T -step videos with corresponding action sequences, generated by an agent acting in an environment. We split each sequence of T time steps into two parts, corresponding to two different model phases:

1. *Memorization phase.* For $t = 1, \dots, \tau$, the model receives at each time step a frame \mathbf{x}_t and action \mathbf{a}_t (e.g. move forwards/backwards, rotate left/right) that led to

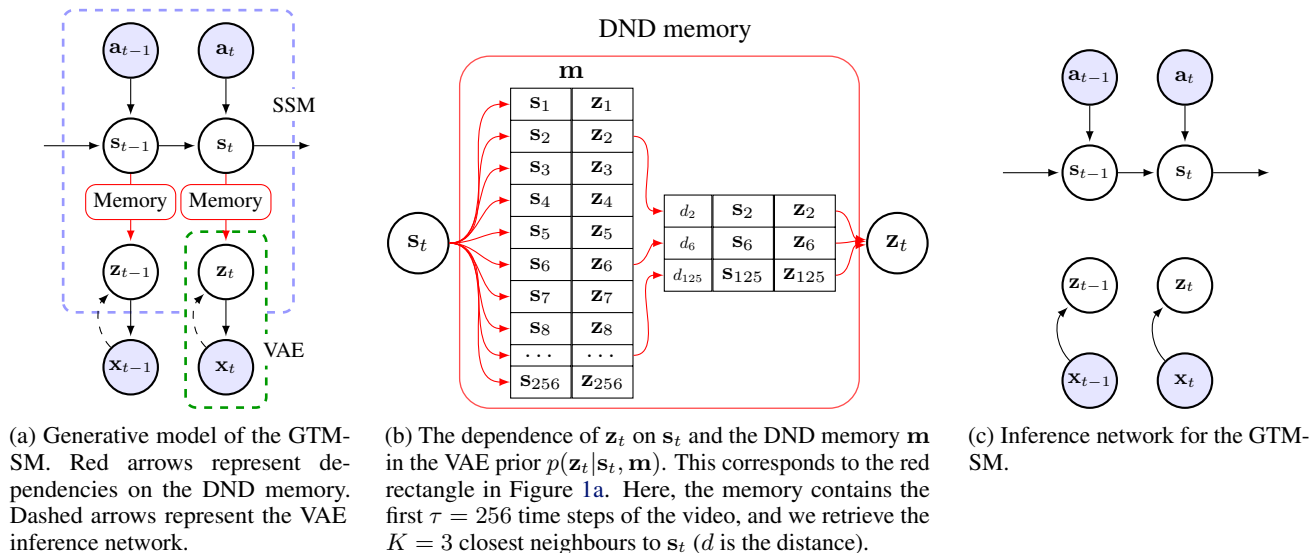


Figure 1: Generative Temporal Model with Spatial Memory

it. In this phase, the model has to store in memory all the information needed in the following prediction phase. During this phase the agent sees most of the environment (but from a restricted set of viewpoints), in order to give the model sufficient information to make accurate predictions in the subsequent phase.

2. *Prediction phase.* For $t = \tau + 1, \dots, T$, the model receives the actions $\mathbf{a}_{\tau+1:T}$ that move the data-generating agent across the previously explored environment (although perhaps viewed from a different angle) and needs to predict the observations $\mathbf{x}_{\tau+1:T}$ using the information about the past that is stored in the memory.

Storing *what* the agent sees at each time step is not sufficient: in order to retrieve the correct information from memory when returning to the same location during the prediction phase, we also need to store *where* the agent is. The location of the agent is a latent variable that can be inferred given the actions as explained in the rest of this section.

As shown in Figure 1a, in a *Generative Temporal Model with Spatial Memory (GTM-SM)* we introduce two sets of latent variables that disentangle visual and dynamics information, similarly to (Fraccaro et al., 2017). At each time step we have a VAE whose latent state \mathbf{z}_t is an encoding of the frame of the video and therefore captures the visual information. The priors of the VAEs are temporally dependent through their dependence on the states \mathbf{s}_t of the SSM, a latent representation of the location of the agent in the environment. The transition density of the SSM is used to include prior knowledge on the environment dynamics, i.e. the underlying physics.

During the initial memorization phase, the GTM-SM infers the states $\mathbf{s}_{1:\tau}$ of the agent (i.e. the position) and the frame encodings $\mathbf{z}_{1:\tau}$, and stores these $(\mathbf{s}_i, \mathbf{z}_i)$ pairs as (key, value) in the DND memory. Probabilistically, we can view this as inserting an approximation of the intractable the posterior $p(\mathbf{s}_{1:\tau}, \mathbf{z}_{1:\tau}|\mathbf{x}_{1:\tau}, \mathbf{a}_{1:\tau})$ into the DND; see Section 3.3 for details. As the latent variables are stochastic, in practice we store in memory the sufficient statistics of the distribution (e.g. the mean and variance in the case of Gaussian variables). To keep the notation simple, we will refer to the information in the DND by the name of the random variable (as done in Figure 1b), rather than introducing a new symbol for the sufficient statistics. An alternative is to insert one or more samples from the distribution into the memory instead, but this would introduce some sampling noise.

In the subsequent prediction phase, we forward-generate from the SSM using the actions $\mathbf{a}_{\tau+1:T}$ to predict $\mathbf{s}_{\tau+1:T}$, and we use the VAE’s generative model to generate the frames $\mathbf{x}_{\tau+1:T}$ given the predicted states and the information from the first τ time steps stored in the DND memory; see Section 3.1 for details. In our experiments, a low-dimensional state vector \mathbf{s}_t (2- or 3-dimensional) suffices. Because of this, we can perform efficient rollouts in latent space without the need to generate high-dimensional frames at each time step as in autoregressive models (Oh et al., 2015; Chiappa et al., 2017; Gemici et al., 2017). Also, thanks to the scalability properties of the DND memory, we can efficiently explore very large environments.

There are three key components that define the GTM-SM and that will be introduced in the following, namely the *generative model*, the *inference network*, and the *past encoder*. As we will see, these components share many parameters.

3.1. Generative model

For brevity, we write the observations and actions in the memorization phase as $\mathbf{v} = \{\mathbf{x}_{1:\tau}, \mathbf{a}_{1:\tau}\}$ and we write the observations and actions in the prediction phase as $\mathbf{x} = \mathbf{x}_{\tau+1:T}$ and $\mathbf{a} = \mathbf{a}_{\tau+1:T}$ respectively. Letting θ be the parameters of the generative model, we model $p_\theta(\mathbf{x}|\mathbf{a}, \mathbf{v})$ as follows. We introduce two sets of latent variables: the frame encodings $\mathbf{z} = \mathbf{z}_{\tau+1:T}$ and the SSM states $\mathbf{s} = \mathbf{s}_{\tau+1:T}$, and define the joint probability density $p_\theta(\mathbf{x}, \mathbf{z}, \mathbf{s}|\mathbf{a}, \mathbf{v})$ following the factorization shown in Figure 1a:

$$p_\theta(\mathbf{x}, \mathbf{z}, \mathbf{s}|\mathbf{a}, \mathbf{v}) = \prod_{t=\tau+1}^T p_\theta(\mathbf{x}_t|\mathbf{z}_t) p_\theta(\mathbf{z}_t|\mathbf{s}_t, \mathbf{m}) \cdot p_\theta(\mathbf{s}_t|\mathbf{s}_{t-1}, \mathbf{a}_t), \quad (1)$$

where $p_\theta(\mathbf{s}_t|\mathbf{s}_{t-1}, \mathbf{a}_t)$ is a Gaussian SSM transition probability density and $p_\theta(\mathbf{x}_t|\mathbf{z}_t)$ is the VAE decoder (Bernoulli or Gaussian distributed, depending on the data). $p_\theta(\mathbf{z}_t|\mathbf{s}_t, \mathbf{m})$ can be seen as the prior of the VAE, that is conditioned at each time step on the current state \mathbf{s}_t and the content of the DND memory $\mathbf{m} = \{\mathbf{s}_{1:\tau}, \mathbf{z}_{1:\tau}\} = \text{PastEncoder}(\mathbf{v})$ as illustrated in Figure 1b (see Section 3.3 for details on the past encoder). Its sufficient statistics are computed as follows. First, we calculate the distances $d_i = d(\mathbf{s}_i, \mathbf{s}_t)$, $i = 1, \dots, \tau$, between \mathbf{s}_t and all the states \mathbf{s}_i in the DND memory. We then retrieve from the memory the K nearest states and the corresponding frame encodings, thus forming a set of triplets $\{(d^{(k)}, \mathbf{s}^{(k)}, \mathbf{z}^{(k)}), k = 1, \dots, K\}$ that will be used as conditioning variables when computing the parameters of the VAE prior $p_\theta(\mathbf{z}_t|\mathbf{s}_t, \mathbf{m})$. Using low-dimensional \mathbf{s}_t and prior knowledge of the environment dynamics when defining $p_\theta(\mathbf{s}_t|\mathbf{s}_{t-1}, \mathbf{a}_t)$, we can make the GTM-SM learn to use \mathbf{s}_t to represent its position in the environment. At each time step the model will then retrieve from the memory what it has seen when it was previously close to the same location, and use this information to generate the current frame \mathbf{x}_t . The exact form of the VAE prior $p_\theta(\mathbf{z}_t|\mathbf{s}_t, \mathbf{m})$ and transition model $p_\theta(\mathbf{s}_t|\mathbf{s}_{t-1}, \mathbf{a}_t)$ is environment-dependent, and will be therefore introduced separately for each experiment in Section 4.

3.2. Inference network

Due to the non-linearities in the VAE and the fact that $p_\theta(\mathbf{z}_t|\mathbf{s}_t, \mathbf{m})$ depends on the DND memory, the posterior distribution $p_\theta(\mathbf{z}, \mathbf{s}|\mathbf{x}, \mathbf{a}, \mathbf{v})$ of the GTM-SM is intractable. We therefore introduce a variational approximation $q_\phi(\mathbf{z}, \mathbf{s}|\mathbf{x}, \mathbf{a})$ that factorizes as

$$q_\phi(\mathbf{z}, \mathbf{s}|\mathbf{x}, \mathbf{a}) = q_\phi(\mathbf{z}|\mathbf{x}) p_\theta(\mathbf{s}|\mathbf{a}) = \prod_{t=\tau+1}^T q_\phi(\mathbf{z}_t|\mathbf{x}_t) p_\theta(\mathbf{s}_t|\mathbf{s}_{t-1}, \mathbf{a}_t). \quad (2)$$

A graphical representation of the inference network of the GTM-SM is shown in Figure 1c. $q_\phi(\mathbf{z}_t|\mathbf{x}_t)$ is an inference network that outputs the mean and variance of a Gaussian distribution, as typically done in VAEs. In (2) we then use the SSM transition probabilities, and we are therefore assuming that the GTM-SM can learn the prior dynamics of the moving agents accurately enough to infer the position of the agent given the sequence of actions. Notice that without this assumption it would be impossible to perform long term generation with the model during the prediction phase. In this phase, we can in fact only rely on the generative model, and not on the inference network as we do not know what the agent is seeing at each time step. To relax this assumption, the inference network could be extended to make use of the information stored in memory, for example by using landmark information when inferring the current position of the agent. This is discussed more in detail in Appendix D in the supplementary material, together with an initial experiment to assess the feasibility of the proposed method.

3.3. Past encoder

The past encoder is used during the memorization phase to extract the information to store in the DND memory. It creates a mapping from \mathbf{s}_t to \mathbf{z}_t that is exploited at times $t = \tau + 1 : T$ in the generative model. During the memorization phase, at each time step we store in the DND the sufficient statistics of the inferred states \mathbf{s}_t and visual information \mathbf{z}_t , $t = 1, \dots, \tau$, obtained from an approximation to the smoothed posterior $p_\theta(\mathbf{z}_t, \mathbf{s}_t|\mathbf{v})$. We first factorize this distribution as $p_\theta(\mathbf{z}_t, \mathbf{s}_t|\mathbf{v}) = p_\theta(\mathbf{s}_t|\mathbf{v}) p_\theta(\mathbf{z}_t|\mathbf{v})$ and only condition on the information up to time t (i.e. we are doing *filtering* instead of smoothing): $p_\theta(\mathbf{z}_t, \mathbf{s}_t|\mathbf{v}) \approx p_\theta(\mathbf{s}_t|\mathbf{v}_{1:t}) p_\theta(\mathbf{z}_t|\mathbf{v}_{1:t})$, where $\mathbf{v}_{1:t} = \{\mathbf{x}_{1:t}, \mathbf{a}_{1:t}\}$. Second, we approximate each of the terms using the inference network, without introducing any additional parameters in the model. Using the VAE encoder, we have $p_\theta(\mathbf{z}_t|\mathbf{v}_{1:t}) \approx q_\phi(\mathbf{z}_t|\mathbf{x}_t)$. We then assume $p_\theta(\mathbf{s}_t|\mathbf{v}_{1:t}) \approx p_\theta(\mathbf{s}_t|\mathbf{v}_{1:t-1}, \mathbf{a}_t)$ with

$$p_\theta(\mathbf{s}_t|\mathbf{v}_{1:t-1}, \mathbf{a}_t) = \int p_\theta(\mathbf{s}_t|\mathbf{s}_{t-1}, \mathbf{a}_t) p_\theta^*(\mathbf{s}_{t-1}) d\mathbf{s}_{t-1},$$

$$p_\theta^*(\mathbf{s}_t) = \int p_\theta(\mathbf{s}_t|\mathbf{s}_{t-1}, \mathbf{a}_t) p_\theta^*(\mathbf{s}_{t-1}) d\mathbf{s}_{t-1}. \quad (3)$$

$p_\theta^*(\mathbf{s}_t)$ is the marginal distribution of \mathbf{s}_t obtained by integrating over the past states (samples from $p_\theta^*(\mathbf{s}_t)$, as needed in the ELBO, are easily obtained with ancestral sampling).

3.4. Training

We learn the parameters θ and ϕ of the GTM-SM by maximizing the ELBO, a lower bound to the log-likelihood $\log p_\theta(\mathbf{x}|\mathbf{a}, \mathbf{v})$ obtained using Jensen's inequality and the

inference network introduced in Section 3.2:

$$\begin{aligned} \log p_\theta(\mathbf{x}|\mathbf{a}, \mathbf{v}) &= \log \int p_\theta(\mathbf{x}, \mathbf{z}, \mathbf{s}|\mathbf{a}, \mathbf{v}) \, d\mathbf{z} \, d\mathbf{s} \\ &\geq \mathbb{E}_{q_\phi(\mathbf{z}, \mathbf{s}|\mathbf{x}, \mathbf{a}, \mathbf{v})} \left[\log \frac{p_\theta(\mathbf{x}, \mathbf{z}, \mathbf{s}|\mathbf{a}, \mathbf{v})}{q_\phi(\mathbf{z}, \mathbf{s}|\mathbf{x}, \mathbf{a}, \mathbf{v})} \right] = \mathcal{F}(\theta, \phi) . \end{aligned}$$

Exploiting the temporal factorization of both the joint distribution $p_\theta(\mathbf{x}, \mathbf{z}, \mathbf{s}|\mathbf{a}, \mathbf{v})$ and the variational approximation $q_\phi(\mathbf{z}, \mathbf{s}|\mathbf{x}, \mathbf{a}, \mathbf{v})$, we obtain after some calculations:

$$\begin{aligned} \mathcal{F}(\theta, \phi) &= \sum_{t=\tau+1}^T \mathbb{E}_{q_\phi(\mathbf{z}_t|\mathbf{x}_t)} [\log p_\theta(\mathbf{x}_t|\mathbf{z}_t)] + \\ &\quad - \mathbb{E}_{p_\theta^*(\mathbf{s}_t)} [KL[q_\phi(\mathbf{z}_t|\mathbf{x}_t)||p_\theta(\mathbf{z}_t|\mathbf{s}_t, \mathbf{m})]] . \end{aligned}$$

The ELBO is then formed by two terms: a reconstruction term and an expected KL divergence for the VAE, which depends on the transition model through the expectation over the marginal distribution $p_\theta^*(\mathbf{s}_t)$ in (3). $\mathcal{F}(\theta, \phi)$ can be maximized jointly with respect to θ and ϕ with stochastic gradient ascent, approximating the intractable expectations with Monte Carlo integration with a single sample and using the reparameterization trick to obtain low-variance gradients (Kingma & Welling, 2014; Rezende et al., 2014). For a given training sequence, we compute the distribution needed in the ELBO as follows: first, we use the past encoder to extract the information to be stored in the DND memory \mathbf{m} . We then obtain samples from the distribution $p_\theta^*(\mathbf{s}_t)$ by iteratively applying the transition model, and use them to read from the DND memory as needed to compute the VAE prior $p_\theta(\mathbf{z}_t|\mathbf{s}_t, \mathbf{m})$. Finally, the variational approximation and the likelihood are computed as usual in VAEs.

4. Experiments

We test the memorization and long-term generation capabilities of the GTM-SM on several 2D and 3D environments of increasing complexity. We use videos with action data from RL agents walking in the environments, that are split so that both the memorization and prediction phase are hundreds of time steps. Experimental details can be found in Appendix A. Videos of long-term generations from the model for all the experiments are available in the supplementary material, see Appendix B for details.

4.1. Image navigation experiment

In this experiment the data-generating agent walks on top of an image and observes a cropped version of the image (centered at the agent’s position). As illustrated in Figure 2a, the 2D environment is a 32x32 image from the CelebA dataset (Liu et al., 2015) and the agent sees an 8x8 crop (the yellow square in the figure). There are five possible actions: move one step up/down/left/right or stay still. At each time step we sample a random action, but to favor exploring the

whole environment the action is repeated in the subsequent time steps. The number of repetitions is sampled from a Poisson distribution. The agent cannot walk outside of the image: the “move right” action, for example, does not change the position of an agent on the right edge of the image. We can interpret this as an environment with walls. The agent walks on the image while adding information in the DND memory for $\tau = 256$ time steps. During training the prediction phase has 32 time steps; during testing we instead generate from the model for 256 time steps, so that $T = 512$. In each of the two dimensions, there are nine possible positions (the crops can overlap). This is illustrated in Figure 2b, which shows the ground truth positions that the agent has visited in the 256 steps of the memorization phase of a test sequence.

We use a 2-dimensional state space that the GTM-SM learns to use to represent the position of the agent. With no walls in the environment all possible transitions are linear, and they can be modelled as $\mathbf{s}_t = \mathbf{s}_{t-1} + M\mathbf{a}_t + \boldsymbol{\varepsilon}_t$ with $\boldsymbol{\varepsilon}_t \sim \mathcal{N}(\mathbf{0}, r^2\mathbf{I})$. In all experiments we use small values for r^2 , that make the transitions close to being deterministic. The transition matrix M can be learned from the data, and describes how to update the state given each of the 5 actions ($M\mathbf{a}_t$ is the *displacement* at time t). The environment in our experiment has walls however, and we need the model to be able to learn not to move with actions that would make it hit a wall. We do this by multiplying the displacement by a neural network σ_d that receives as input the projected position of the agent after taking the action and outputs a value between 0 and 1, and that can therefore learn to cancel out any displacements that would bring the agent out of the environment. These non-linear transitions are therefore modelled as

$$\mathbf{s}_t = \mathbf{s}_{t-1} + M\mathbf{a}_t \cdot \sigma_d(\mathbf{s}_{t-1} + M\mathbf{a}_t) + \boldsymbol{\varepsilon}_t . \quad (4)$$

The VAE prior used in this experiments is obtained by creating a mixture distribution from the sufficient statistics of the frame encodings retrieved from the DND memory, see Appendix A.1 for details.

In Figure 2c we show an example of the states inferred by the model for a test sequence. We see that the model has learned the correct transitions, in a state space that is rotated and stretched with respect to the ground truth one. To test the memorization and prediction capabilities of the GTM-SM, Figure 2d shows a comparison between the ground truth frames of the video and the predicted ones during the prediction phase. The model produces almost perfect predictions, even after more than 200 generation steps ($t = 471$). This shows that it has learned to store all relevant information in the DND, as well as retrieve all relevant information from it. The state-of-the-art generative temporal models with memory introduced in (Gemici et al., 2017), are not able to capture the spatial structure of large environments

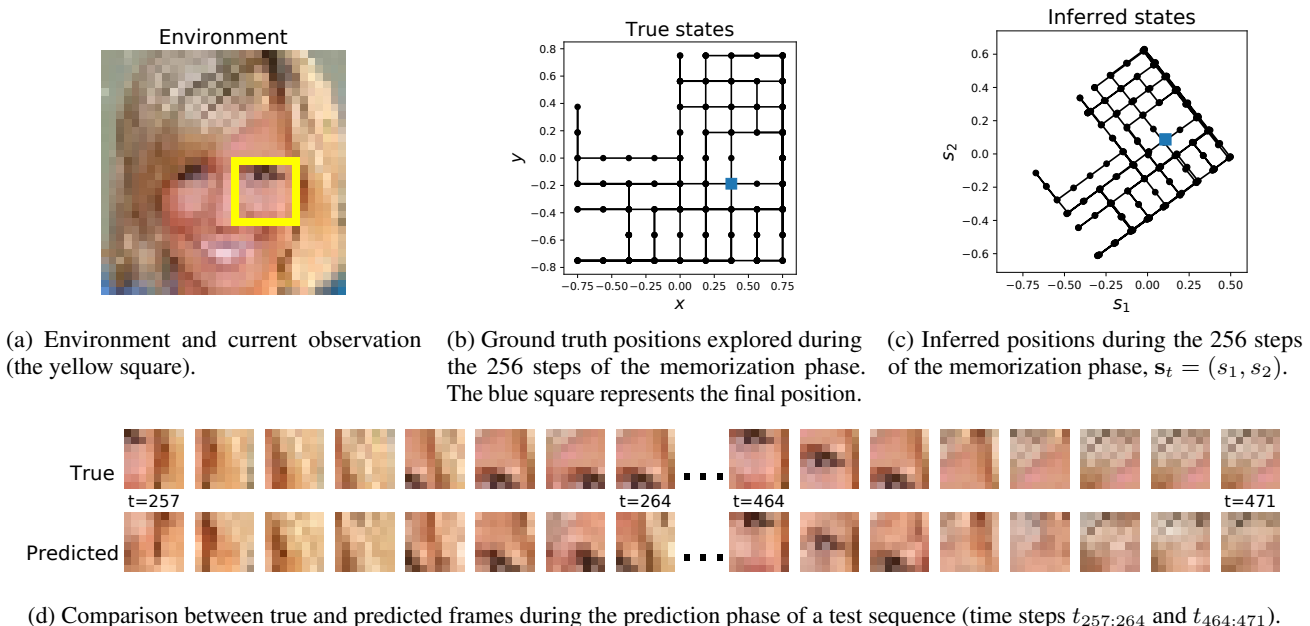


Figure 2: Image navigation experiment

as in the GTM-SM, and would therefore struggle to coherently generate hundreds of time steps into the future. The MNIST maze experiment in (Gemici et al., 2017) can be seen as a simpler version of the image navigation experiment presented above, with agents moving on a 4x4 grid, linear transitions and 25-step sequences.

4.2. Labyrinth experiments

We now show that the GTM-SM is able to remember the past and perform spatio-temporally coherent generations over hundreds of time steps in simulated 3D environments. We use the *Labyrinth* environment (Mnih et al., 2016; Beattie et al., 2016), procedurally-generated 3D mazes with random textures, objects and wall configurations. There are eight possible actions that can both move and rotate the agent in the maze, and we observe images from a first-person point of view. Labyrinth can be seen as a 3D extension of the image navigation experiments in Section 4.1, but in this case the task is much harder for two main reasons that will be tackled below: (1) dealing with rotations and (2) projective transformations in a partially observable environment.

First, the state of the agent is no longer only described by the position, but also from the direction in which the agent is looking and moving. We need to take into account two different coordinate systems, a global one that coincides with the state-space (s -space), and one that is fixed with the agent (agent-space). In the image navigation experiments these coordinate systems coincided. The actions act in agent-space, e.g. a “move right” action will make the agent go right in its reference frame, but depending on its orientation

this could correspond to a move to the left in s -space. To deal with this issues we can introduce a *rotation matrix* R in the state transition equation, that translates a displacement $M\mathbf{a}_t$ in agent-space to a displacements $R M\mathbf{a}_t$ in s -space. More in detail, we consider a 3-dimensional state-space, and define the state transition equations as

$$\mathbf{s}_t = \mathbf{s}_{t-1} + R(\mathbf{s}_{t-1}^{(3)})M\mathbf{a}_t + \varepsilon_t \quad (5)$$

with $\mathbf{s}_{t-1}^{(3)}$ being the 3rd component of the vector \mathbf{s}_{t-1} and

$$R(\mathbf{s}_{t-1}^{(3)}) = \begin{bmatrix} \cos(\mathbf{s}_{t-1}^{(3)}) & -\sin(\mathbf{s}_{t-1}^{(3)}) & 0 \\ \sin(\mathbf{s}_{t-1}^{(3)}) & \cos(\mathbf{s}_{t-1}^{(3)}) & 0 \\ 0 & 0 & 1 \end{bmatrix}.$$

As for the image navigation experiment, we learn the parameters of M . While we do not explicitly tell the GTM-SM to use the first two component of the state vector as a position and the third one as an angle, the model will learn to use them in this way in order to maximize the ELBO. In the nearest neighbor search in the DND memory we need to take into account the periodicity of the angle, e.g. that an agent oriented at 30° or 390° is actually looking in the same direction. When computing distances, instead of using $\mathbf{s}_t^{(3)}$ we then use $\cos(\mathbf{s}_t^{(3)})$ and $\sin(\mathbf{s}_t^{(3)})$, that are possibly passed together with the first two components of \mathbf{s}_t through a linear layer that maps the resulting vector in a learned space where we use the Euclidean distance.

The second challenge arises from the fact that, unlike the image navigation experiment where there were a limited number of possible positions for the agent, in the 3D

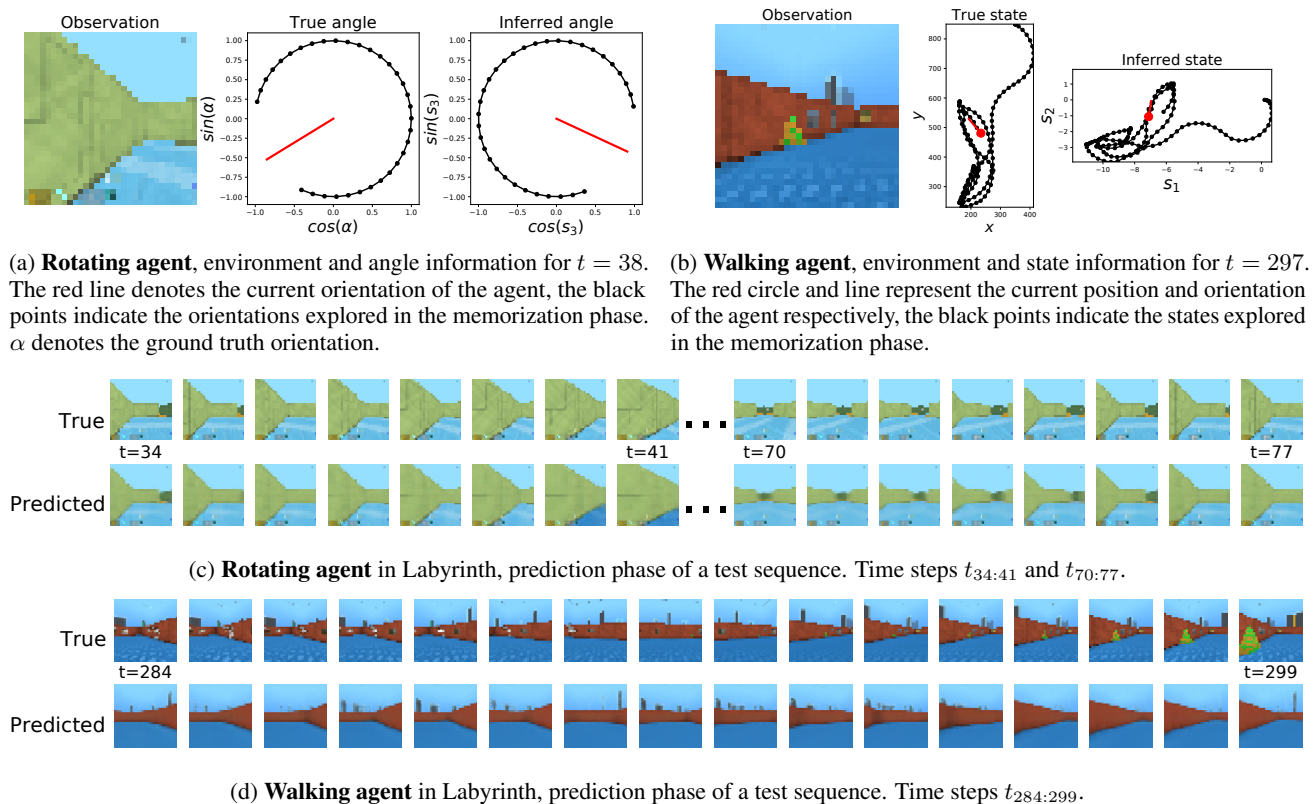


Figure 3: Labyrinth experiment

labyrinth environment it is not reasonable to assume that during the memorization phase the agent will pass in all positions and look from them in all directions. To deal with this issue, we use as VAE prior $p_\theta(\mathbf{z}_t | \mathbf{s}_t, \mathbf{m})$ a *Generative Query Network* (Eslami et al., 2018), a neural architecture that given the frames from the closest positions retrieved from the DND memory learns to combine the different views taking into account projective transformations

4.2.1. ROTATING AGENT IN LABYRINTH

In the first experiment we test the abilities of the GTM-SM to learn to model rotations with the transition model in (5) as well as to combine the information from different views. We use videos with action data of an agent that does two complete rotations while standing still in the same position. The rotational period is around 41 time steps, but we only store in memory the first $\tau = 33$ (approximately 300°). We then ask the model to generate the remaining 60° to finish the first rotation and a whole new rotation. From Figure 3a we see an example of an observation from the test data and that the model has correctly learned to use the third component of the state vector \mathbf{s}_t to represent the orientation of the agent. In the prediction phase in Figure 3c, we notice that the predictions from the model are very close to the ground truth, meaning that the model has learned to

use the memory correctly. In particular, despite the fact that the frames from $t = 33$ to $t = 41$ were never seen during the memorization phase, the GTM-SM has learned to combine the information from other views. Notice that this experiment can be seen as a more challenging version of the Labyrinth rotation experiment of (Gemici et al., 2017), that used a fully observed first rotation with a rotational period of 15 time steps.

4.2.2. WALKING AGENT IN LABYRINTH

We now use videos of a pre-trained RL agent walking in a room and solving a scavenger hunt task. In this case it is fundamental to extend the transition equation in (5) to model more carefully the physics of the walking agent, that make the displacement at a given time step depend not only on the current action, but also on the displacement at the previous time step. The agent is subject to momentum and friction, so that if it is moving in a certain direction at time $t - 1$, it will still continue to move a bit in the same direction even at time t , regardless of the action \mathbf{a}_t . Also, despite the momentum, the displacement of the agent cannot increase indefinitely, i.e. there is saturation. We can model this by extending the way the displacement $\mathbf{d}_t = M\mathbf{a}_t$ is calculated in (5). To take into account momentum and friction, we first add to $\mathbf{d}_t = M\mathbf{a}_t$ a damped version of the displacement at

the previous time step, i.e. $\sigma(\mathbf{c}_f) \odot \mathbf{d}_{t-1}$, where $\sigma(\mathbf{c}_f)$ is a learned vector between 0 and 1 of the same size of \mathbf{d}_t (σ is the sigmoid function and \odot represents the element-wise product). To deal with saturation, we then limit the range of the displacements by squashing them through a \tanh non-linearity that is pre-multiplied by a learned vector \mathbf{c}_s . The resulting transition model then becomes

$$\mathbf{s}_t = \mathbf{s}_{t-1} + R(\mathbf{s}_{t-1}^{(3)}) \underbrace{\mathbf{c}_s \odot \tanh(\sigma(\mathbf{c}_f) \odot \mathbf{d}_{t-1} + M\mathbf{a}_t)}_{\mathbf{d}_t} + \varepsilon_t$$

with M , \mathbf{c}_s and \mathbf{c}_f parameters to be learned. Notice that modelling the displacements in this way essentially corresponds to extending the state space with an additional vector \mathbf{d}_t that models the velocities of the agent. We can then model accelerations/decelerations of the agent by non-linearly changing the velocities over time in this extended state-space. This is a very challenging experiment, as from only one-hot encoded actions and the frames of the video we need to learn a complex transition model. Moreover, due to the low resolution of the images (32x32), small variations in state-space may be impossible to infer using only the images. To solve this, we make the reasonable assumption that at training time the agent feels its movement, i.e. the displacements. We then add a regression loss as an extra term to the objective function, that helps the GTM-SM to learn transition parameters such that the estimated \mathbf{d}_t is close to its true value. We only add the true displacements information as a target in the loss, and never pass it directly into the model.

We let the agent walk around the room while adding information in the DND memory for $\tau = 150$ time steps, and then predict during testing the following 150 time steps ($T = 300$). In Figure 3b, we notice that the GTM-SM is able to learn a very accurate transition model, that provides a sufficiently good approximation of the true state even after $t = 297$ time steps. In Figure 3d we can appreciate the memorization and long-term generation capabilities of the GTM-SM by looking at the comparison between the true and predicted frames of the video in the end of the prediction phase ($t_{284:299}$). We also notice in the predicted frames, that the model correctly draws the walls and the floors but fails to render the objects, probably due to difficulties in modelling with the VAE the very diverse and complex textures that form objects in this environment. We also tested the same trained model on longer videos of larger environments with multiple rooms ($\tau = 150$ and $T = 450$). As explained in detail in Appendix C, the model is able to correctly predict the textures of the environment even after 300 time steps.

5. Related work

A number of recent works have augmented deep generative models with learned external memories, both in the static setting (Li et al., 2016; Bornschein et al., 2017) and in the

temporal one (Gemici et al., 2017). More in general, neural networks have been combined with different memories in a wide range of tasks such as supervised learning (Graves et al., 2014; 2016), reinforcement learning (Oh et al., 2016; Pritzel et al., 2017; Parisotto & Salakhutdinov, 2018), one-shot learning (Santoro et al., 2016), question answering and language modelling (Sukhbaatar et al., 2015; Miller et al., 2016). Each memory architecture uses different addressing mechanisms to write or read information, that are usually chosen depending on the specific application being considered. As discussed in the introduction, our work is closely related to (Gemici et al., 2017), but more suitable for long-term generation for this task and more scalable thanks to the usage of the spatial memory architecture that exploits knowledge on the dynamics of the agent and does not require to learn a parametric memory addressing scheme.

In the deep reinforcement learning community, several works have exploited different memory architectures to store long term information to be used within an agent’s policy, such as in (Zaremba & Sutskever, 2015; Oh et al., 2016). In particular, in (Gupta et al., 2017a;b; Zhang et al., 2017; Parisotto & Salakhutdinov, 2018) the memory architectures have a fixed number of slots that are spatially structured as a 2D grid, and can therefore store information on the moving agent. Similarly to the GTM-SM, these memories are built to exploit the spatial structure of the problem, although for the different task of constructing agents that can learn to navigate and explore the environment, as opposed to the focus on generative modelling of this paper. Simultaneous Localization And Mapping (SLAM) (Smith et al., 1987; Leonard & Durrant-Whyte, 1991) is a popular technique used in robotics to estimate the position of a robot and the map of the environment at the same time using sensor data, recently applied to deep reinforcement learning for example in (Bhatti et al., 2016; Zhang et al., 2017). It is reminiscent to the memorization phase of the GTM-SM, that could be therefore extended using ideas introduced in the visual SLAM community (Taketomi et al., 2017).

6. Conclusion

In this work we introduced an action-conditioned generative model that uses a scalable non-parametric memory to store spatial and visual information. Our experiments on simulated 2D and 3D environments show that the model is able to coherently memorize and perform long-term generation. To our knowledge this is the first published work that builds a generative model for agents walking in an environment that, thanks to the separation of the dynamics and visual information in the DND memory, can coherently generate for hundreds of time steps in a scalable way. Future work will focus on exploiting these capabilities in model-based planning by integrating the GTM-SM within an RL agent.

Acknowledgements

We would like to thank Charles Blundell and Timothy Lillicrap for many helpful discussions, and Ivo Danihelka and Joshua Abramson for generating the datasets used in the Labyrinth experiments.

References

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015.
- Beattie, C., Leibo, J. Z., Teplyaev, D., Ward, T., Wainwright, M., Küttler, H., Lefrancq, A., Green, S., Valdés, V., Sadik, A., Schrittwieser, J., Anderson, K., York, S., Cant, M., Cain, A., Bolton, A., Gaffney, S., King, H., Hassabis, D., Legg, S., and Petersen, S. DeepMind Lab. *CoRR*, abs/1612.03801, 2016.
- Bentley, J. Multidimensional binary search trees used for associative searching. *Commun. ACM*, 1975.
- Bhatti, S., Desmaison, A., Miksik, O., Nardelli, N., Siddharth, N., and Torr, P. H. S. Playing doom with slam-augmented deep reinforcement learning. *arXiv:1612.00380*, 2016.
- Bornschein, J., Mnih, A., Zoran, D., and Jimenez Rezende, D. Variational memory addressing in generative models. In *Advances in Neural Information Processing Systems 30*. 2017.
- Chiappa, S., Racanière, S., Wierstra, D., and Mohamed, S. Recurrent environment simulators. In *ICLR*, 2017.
- Chung, J., Gulcehre, C., Cho, K., and Bengio, Y. Empirical evaluation of gated recurrent neural networks on sequence modeling. *arXiv:1412.3555*, 2014.
- Chung, J., Kastner, K., Dinh, L., Goel, K., Courville, A. C., and Bengio, Y. A recurrent latent variable model for sequential data. In *NIPS*, 2015.
- Deisenroth, M. P. and Rasmussen, C. E. Pilco: A model-based and data-efficient approach to policy search. In *Proceedings of the 28th International Conference on Machine Learning, ICML'11*, 2011.
- Eslami, S. M. A., Jimenez Rezende, D., Besse, F., Viola, F., Morcos, A. S., Garnelo, M., Ruderman, A., Rusu, A. A., Danihelka, I., Gregor, K., Reichert, D. P., Buesing, L., Weber, T., Vinyals, O., Rosenbaum, D., Rabinowitz, N., King, H., Hillier, C., Botvinick, M., Wierstra, D., Kavukcuoglu, K., and Hassabis, D. Neural scene representation and rendering. *Science*, 2018. ISSN 0036-8075.
- Finn, C. and Levine, S. Deep visual foresight for planning robot motion. In *Robotics and Automation (ICRA), 2017 IEEE International Conference on*, pp. 2786–2793. IEEE, 2017.
- Fraccaro, M., Sønderby, S. K., Paquet, U., and Winther, O. Sequential neural models with stochastic layers. In *NIPS*, 2016.
- Fraccaro, M., Kamronn, S., Paquet, U., and Winther, O. A disentangled recognition and nonlinear dynamics model for unsupervised learning. *Advances in Neural Information Processing Systems 30, NIPS*, 2017.
- Gemici, M., Hung, C., Santoro, A., Wayne, G., Mohamed, S., Rezende, D. J., Amos, D., and Lillicrap, T. P. Generative temporal models with memory. *arXiv:1702.04649*, 2017.
- Graves, A., Wayne, G., and Danihelka, I. Neural Turing machines. *arXiv:1410.5401*, 2014.
- Graves, A., Wayne, G., Reynolds, M., Harley, T., Danihelka, I., Grabska-Barwiska, A., Colmenarejo, S. G., Grefenstette, E., Ramalho, T., Agapiou, J., Badia, A. P., Hermann, K. M., Zwols, Y., Ostrovski, G., Cain, A., King, H., Summerfield, C., Blunsom, P., Kavukcuoglu, K., and Hassabis, D. Hybrid computing using a neural network with dynamic external memory. *Nature*, 2016.
- Gupta, S., Davidson, J., Levine, S., Sukthankar, R., and Malik, J. Cognitive mapping and planning for visual navigation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017a.
- Gupta, S., Fouhey, D., Levine, S., and Malik, J. Unifying map and landmark based representations for visual navigation. *arXiv:1712.08125*, 2017b.
- Higgins, I., Pal, A., Rusu, A. A., Matthey, L., Burgess, C., Pritzel, A., Botvinick, M., Blundell, C., and Lerchner, A. DARLA: improving zero-shot transfer in reinforcement learning. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, 2017.
- Hochreiter, S. and Schmidhuber, J. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.

- Kingma, D. and Ba, J. Adam: A method for stochastic optimization. *arXiv:1412.6980*, 2014.
- Kingma, D. and Welling, M. Auto-encoding variational Bayes. In *ICLR*, 2014.
- Lenz, I., Knepper, R. A., and Saxena, A. Deepmpc: Learning deep latent features for model predictive control. In *Robotics: Science and Systems*, 2015.
- Leonard, J. J. and Durrant-Whyte, H. F. Mobile robot localization by tracking geometric beacons. *IEEE Trans. Robotics and Automation*, 1991.
- Levine, S. and Abbeel, P. Learning neural network policies with guided policy search under unknown dynamics. In *Advances in Neural Information Processing Systems 27*. 2014.
- Li, C., Zhu, J., and Zhang, B. Learning to generate with memory. In *Proceedings of the 33rd International Conference on Machine Learning, ICML*, 2016.
- Liu, Z., Luo, P., Wang, X., and Tang, X. Deep learning face attributes in the wild. In *Proceedings of International Conference on Computer Vision (ICCV)*, 2015.
- Miller, A. H., Fisch, A., Dodge, J., Karimi, A.-H., Bordes, A., and Weston, J. Key-value memory networks for directly reading documents. In *EMNLP*, 2016.
- Mnih, V., Badia, A. P., Mirza, M., Graves, A., Lillicrap, T., Harley, T., Silver, D., and Kavukcuoglu, K. Asynchronous methods for deep reinforcement learning. In *ICML*, 2016.
- Oh, J., Guo, X., Lee, H., Lewis, R. L., and Singh, S. Action-conditional video prediction using deep networks in atari games. In *NIPS*, 2015.
- Oh, J., Chockalingam, V., Singh, S. P., and Lee, H. Control of memory, active perception, and action in minecraft. In *Proceedings of the 33rd International Conference on Machine Learning, ICML*, 2016.
- Parisotto, E. and Salakhutdinov, R. Neural map: Structured memory for deep reinforcement learning. In *ICLR*, 2018.
- Pritzel, A., Uria, B., Srinivasan, S., Badia, A. P., Vinyals, O., Hassabis, D., Wierstra, D., and Blundell, C. Neural episodic control. In *Proceedings of the 34th International Conference on Machine Learning*, Proceedings of Machine Learning Research, 2017.
- Racanière, S., Weber, T., Reichert, D., Buesing, L., Guez, A., Jimenez Rezende, D., Puigdomènech Badia, A., Vinyals, O., Heess, N., Li, Y., Pascanu, R., Battaglia, P., Hassabis, D., Silver, D., and Wierstra, D. Imagination-augmented agents for deep reinforcement learning. In *Advances in Neural Information Processing Systems 30*. 2017.
- Rezende, D. J., Mohamed, S., and Wierstra, D. Stochastic backpropagation and approximate inference in deep generative models. In *ICML*, 2014. ISBN 9781634393973.
- Santoro, A., Bartunov, S., Botvinick, M., Wierstra, D., and Lillicrap, T. P. One-shot learning with memory-augmented neural networks. *arXiv:1605.06065*, 2016.
- Smith, R., Self, M., and Cheeseman, P. Estimating uncertain spatial relationships in robotics. In *Proceedings. 1987 IEEE International Conference on Robotics and Automation*, 1987.
- Sukhbaatar, S., Szlam, A., Weston, J., and Fergus, R. End-to-end memory networks. In *Neural Information Processing Systems, NIPS'15*, 2015.
- Sutton, R. S. Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In *In Proceedings of the Seventh International Conference on Machine Learning*, 1990.
- Taketomi, T., Uchiyama, H., and Ikeda, S. Visual slam algorithms: a survey from 2010 to 2016. *IPSJ Transactions on Computer Vision and Applications*, 9:1–11, 2017.
- Wahlström, N., Schön, T. B., and Deisenroth, M. P. From pixels to torques: Policy learning with deep dynamical models. *arXiv:1502.02251*, 2015.
- Watter, M., Springenberg, J., Boedecker, J., and Riedmiller, M. Embed to control: A locally linear latent dynamics model for control from raw images. In *NIPS*, 2015.
- Zaremba, W. and Sutskever, I. Reinforcement learning neural turing machines. *arXiv:1505.00521*, 2015.
- Zhang, J., Tai, L., Boedecker, J., Burgard, W., and Liu, M. Neural slam: Learning to explore with external memory. *arXiv:1706.09520*, 2017.