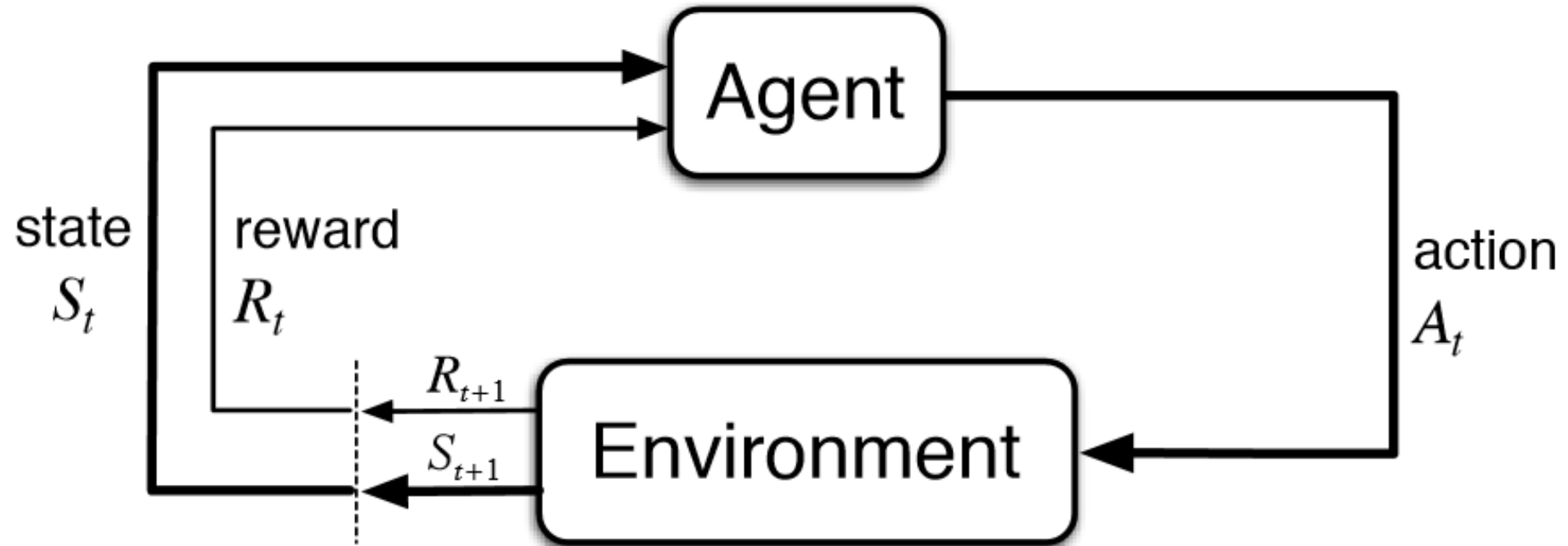


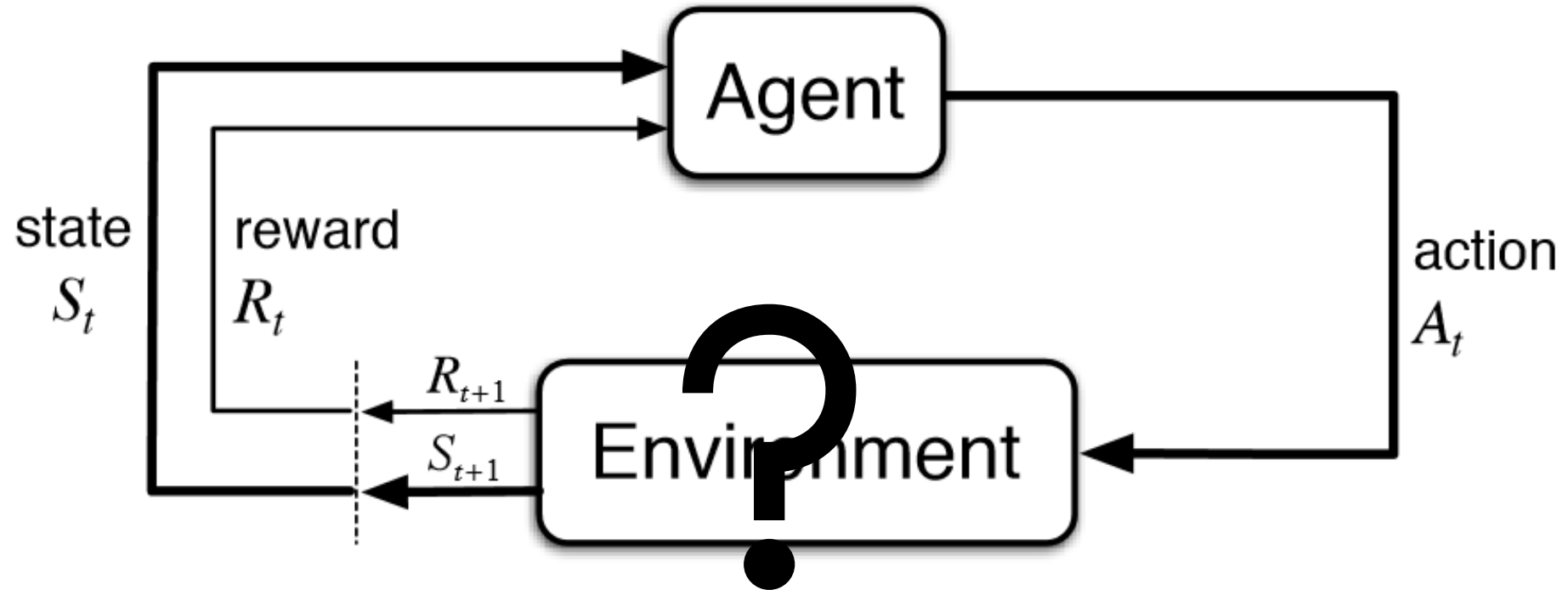
# Model Based Reinforcement Learning

Presenter: Adrian Hoffmann

# Basic Set-Up



# Basic Set-Up



# Overview

- Definitions and their problems
- Why I want a model
- Paper “The Effect of Planning Shape on Dyna-style Planning in High-dimensional State Spaces”
- Paper “World Models”







Definition Model-Based

# Definition Model-Based

---

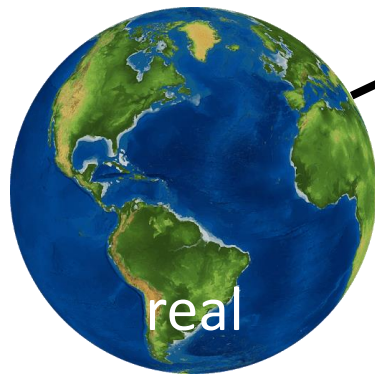
**Algorithm 1** Model-based reinforcement learning

---

```
1: Input: state sample procedure  $d$ 
2: Input: model  $m$ 
3: Input: policy  $\pi$ 
4: Input: predictions  $v$ 
5: Input: environment  $\mathcal{E}$ 
6: Get initial state  $s \leftarrow \mathcal{E}$ 
7: for iteration  $\in \{1, 2, \dots, K\}$  do
8:   for interaction  $\in \{1, 2, \dots, M\}$  do
9:     
10:    
11:    Interactions with the real Environment
12:    
13:    
14:   end for
15:   for planning step  $\in \{1, 2, \dots, P\}$  do
16:     
17:     Take advantage of the model
18:     
19:   end for
20: end for
```

---

# Definition Model-Based



---

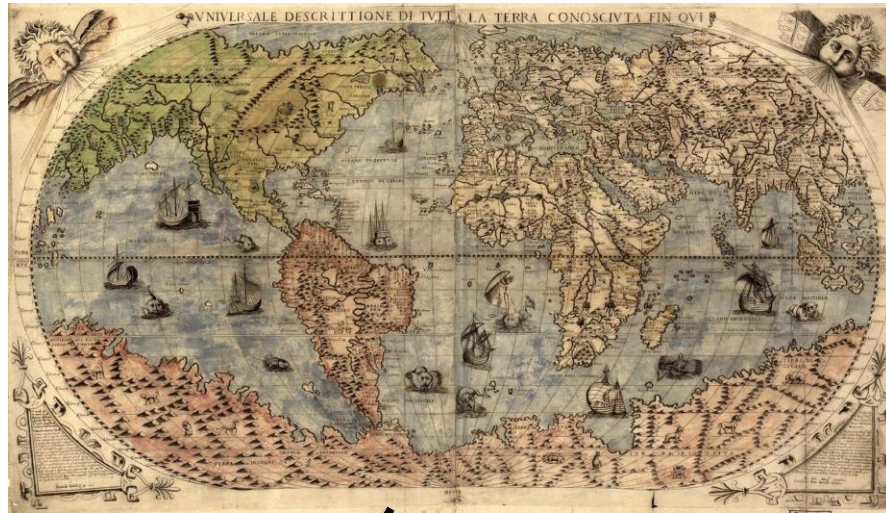
## Algorithm 1 Model-based reinforcement learning

---

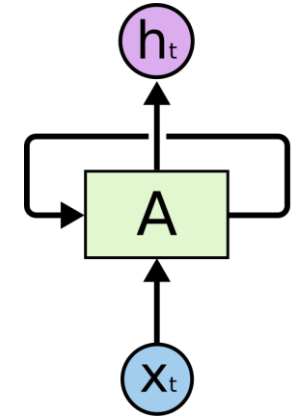
```
1: Input: state sample procedure  $d$ 
2: Input: model  $m$ 
3: Input: policy  $\pi$ 
4: Input: predictions  $v$ 
5: Input: environment  $\mathcal{E}$ 
6: Get initial state  $s \leftarrow \mathcal{E}$ 
7: for iteration  $\in \{1, 2, \dots, K\}$  do
8:   for interaction  $\in \{1, 2, \dots, M\}$  do
9:     Interactions with the real Environment
10:
11:
12:
13:
14:   end for
15:   for planning step  $\in \{1, 2, \dots, P\}$  do
16:     Take advantage of the model
17:
18:
19:   end for
20: end for
```

---

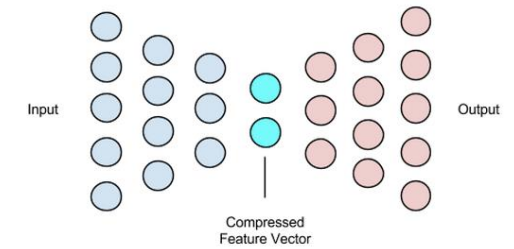
# Interlude – Model



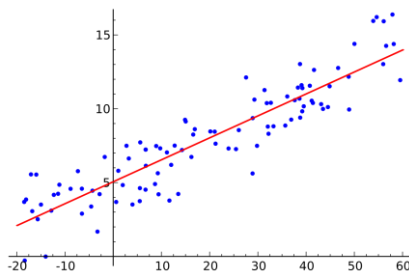
RNN, LSTM



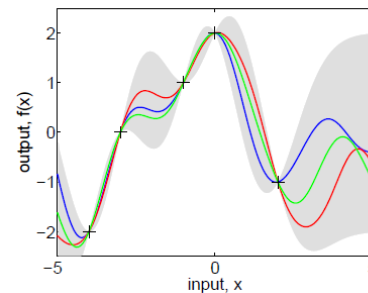
(Variational Auto) Encoders



Statistical Models

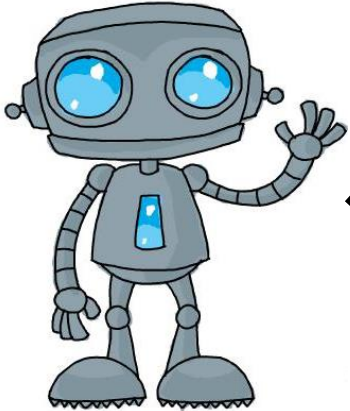


Gaussian Process





# Definition Model-Based



agent

## Algorithm 1 Model-based reinforcement learning


```
1: Input: state sample procedure  $d$ 
2: Input: model  $m$ 
3: Input: policy  $\pi$ 
4: Input: predictions  $v$ 
5: Input: environment  $\mathcal{E}$ 
6: Get initial state  $s \leftarrow \mathcal{E}$ 
7: for iteration  $\in \{1, 2, \dots, K\}$  do
8:   for interaction  $\in \{1, 2, \dots, M\}$  do
9:     [Interactions with the real Environment]
10:
11:
12:
13:
14:   end for
15:   for planning step  $\in \{1, 2, \dots, P\}$  do
16:     [Take advantage of the model]
17:
18:
19:   end for
20: end for
```

# Definition Model-Based

---

**Algorithm 1** Model-based reinforcement learning

---

```
1: Input: state sample procedure  $d$ 
2: Input: model  $m$ 
3: Input: policy  $\pi$ 
4: Input: predictions  $v$ 
5: Input: environment  $\mathcal{E}$ 
6: Get initial state  $s \leftarrow \mathcal{E}$ 
7: for iteration  $\in \{1, 2, \dots, K\}$  do
8:   for interaction  $\in \{1, 2, \dots, M\}$  do
9:     Generate action:  $a \leftarrow \pi(s)$ 
10:    Generate reward, next state:  $r, s' \leftarrow \mathcal{E}(a)$ 
11:     $m, d \leftarrow \text{UPDATEMODEL}(s, a, r, s')$ 
12:     $\pi, v \leftarrow \text{UPDATEAGENT}(s, a, r, s')$ 
13:    Update current state:  $s \leftarrow s'$ 
14:   end for
15:   for planning step  $\in \{1, 2, \dots, P\}$  do
16:     
17:     Take advantage of the model
18:   end for
19: end for
20: end for
```

---

# Definition Model-Based

---

**Algorithm 1** Model-based reinforcement learning

---

```
1: Input: state sample procedure  $d$ 
2: Input: model  $m$ 
3: Input: policy  $\pi$ 
4: Input: predictions  $v$ 
5: Input: environment  $\mathcal{E}$ 
6: Get initial state  $s \leftarrow \mathcal{E}$ 
7: for iteration  $\in \{1, 2, \dots, K\}$  do
8:   for interaction  $\in \{1, 2, \dots, M\}$  do
9:     Generate action:  $a \leftarrow \pi(s)$ 
10:    Generate reward, next state:  $r, s' \leftarrow \mathcal{E}(a)$ 
11:     $m, d \leftarrow \text{UPDATEMODEL}(s, a, r, s')$ 
12:     $\pi, v \leftarrow \text{UPDATEAGENT}(s, a, r, s')$ 
13:    Update current state:  $s \leftarrow s'$ 
14:   end for
15:   for planning step  $\in \{1, 2, \dots, P\}$  do
16:     Generate state, action  $\tilde{s}, \tilde{a} \leftarrow d$ 
17:     Generate reward, next state:  $\tilde{r}, \tilde{s}' \leftarrow m(\tilde{s}, \tilde{a})$ 
18:      $\pi, v \leftarrow \text{UPDATEAGENT}(\tilde{s}, \tilde{a}, \tilde{r}, \tilde{s}')$ 
19:   end for
20: end for
```

---

# Definition (purely) Model-Based

---

**Algorithm 1** Model-based reinforcement learning

---

```
1: Input: state sample procedure  $d$ 
2: Input: model  $m$ 
3: Input: policy  $\pi$ 
4: Input: predictions  $v$ 
5: Input: environment  $\mathcal{E}$ 
6: Get initial state  $s \leftarrow \mathcal{E}$ 
7: for iteration  $\in \{1, 2, \dots, K\}$  do
8:   for interaction  $\in \{1, 2, \dots, M\}$  do
9:     Generate action:  $a \leftarrow \pi(s)$ 
10:    Generate reward, next state:  $r, s' \leftarrow \mathcal{E}(a)$ 
11:     $m, d \leftarrow \underline{\text{UPDATEMODEL}(s, a, r, s')}$ 
12:
13:    Update current state:  $s \leftarrow s'$ 
14:   end for
15:   for planning step  $\in \{1, 2, \dots, P\}$  do
16:     Generate state, action  $\tilde{s}, \tilde{a} \leftarrow d$ 
17:     Generate reward, next state:  $\tilde{r}, \tilde{s}' \leftarrow m(\tilde{s}, \tilde{a})$ 
18:      $\pi, v \leftarrow \underline{\text{UPDATEAGENT}(\tilde{s}, \tilde{a}, \tilde{r}, \tilde{s}'})$ 
19:   end for
20: end for
```

---

# Definition Model-Free

---

**Algorithm 1** Model-based reinforcement learning

---

```
1:
2:
3: Input: policy  $\pi$ 
4: Input: predictions  $v$ 
5: Input: environment  $\mathcal{E}$ 
6: Get initial state  $s \leftarrow \mathcal{E}$ 
7: for iteration  $\in \{1, 2, \dots, K\}$  do
8:   for interaction  $\in \{1, 2, \dots, M\}$  do
9:     Generate action:  $a \leftarrow \pi(s)$ 
10:    Generate reward, next state:  $r, s' \leftarrow \mathcal{E}(a)$ 
11:
12:     $\pi, v \leftarrow \text{UPDATEAGENT}(s, a, r, s')$ 
13:    Update current state:  $s \leftarrow s'$ 
14:   end for
15:
16:
17:
18:
19:
20: end for
```

---

# Definition Model-Based

---

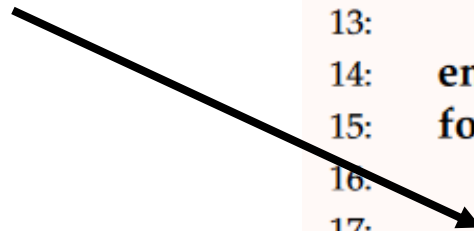
**Algorithm 1** Model-based reinforcement learning

---

```
1: Input: state sample procedure  $d$ 
2: Input: model  $m$ 
3: Input: policy  $\pi$ 
4: Input: predictions  $v$ 
5: Input: environment  $\mathcal{E}$ 
6: Get initial state  $s \leftarrow \mathcal{E}$ 
7: for iteration  $\in \{1, 2, \dots, K\}$  do
8:   for interaction  $\in \{1, 2, \dots, M\}$  do
9:
10:    Interactions with the real
11:    Environment
12:
13:
14:   end for
15:   for planning step  $\in \{1, 2, \dots, P\}$  do
16:
17:    Take advantage of the model
18:
19:   end for
20: end for
```

---

DQN is model-based  
due to its replay buffer



# Definition Model-Based

The fourth and final element of some reinforcement learning systems is a *model* of the environment. This is something that mimics the behavior of the environment, or more generally, that allows inferences to be made about how the environment will behave.

# Definition Model-Based

The fourth and final element of some reinforcement learning systems is a *model* of the environment. This is something that mimics the behavior of the environment, or more generally, that allows inferences to be made about how the environment will behave.

Again, DQN's replay  
buffer fits this description



# Distinction to normal Control Problems

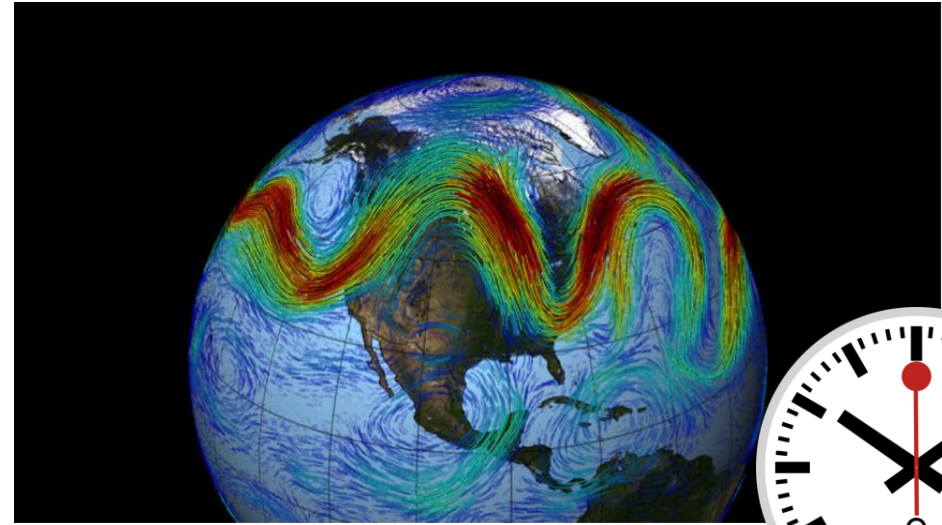


Why I want a model

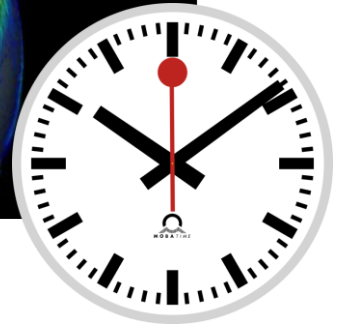
# Why I want a model – Sample efficiency



Save exploration



High computational cost



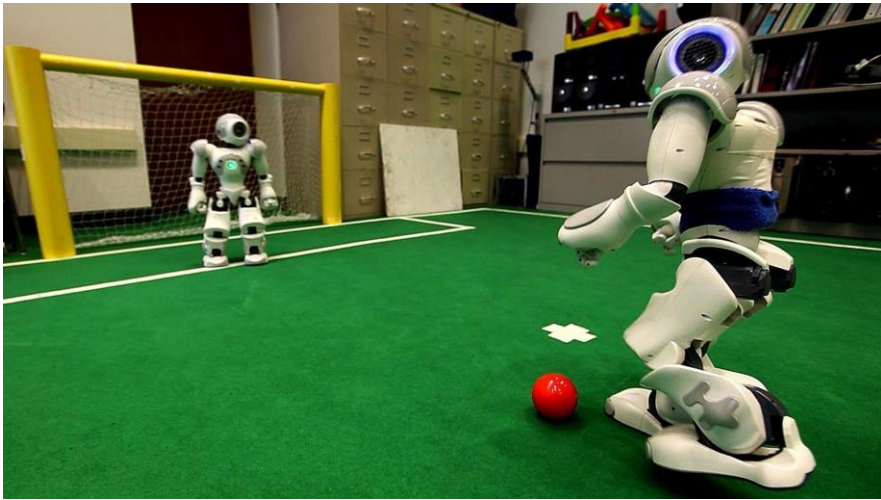
Create Samples with Model

# Why I want a model – Planning

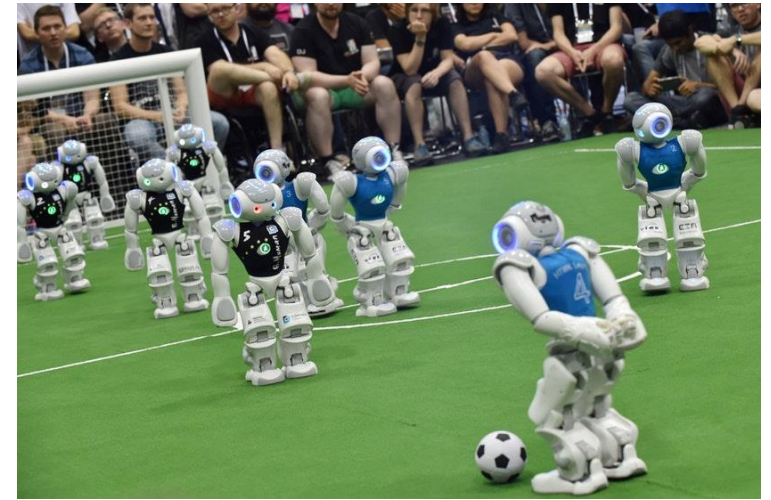


Learn a model and then use a planning algorithm

# Why I want a model – Transfer of Knowledge



Learn how to kick the ball



Concentrate on teamplay

# The effect of Planning Shape On Dyna-style planning in High- dimensional State Spaces

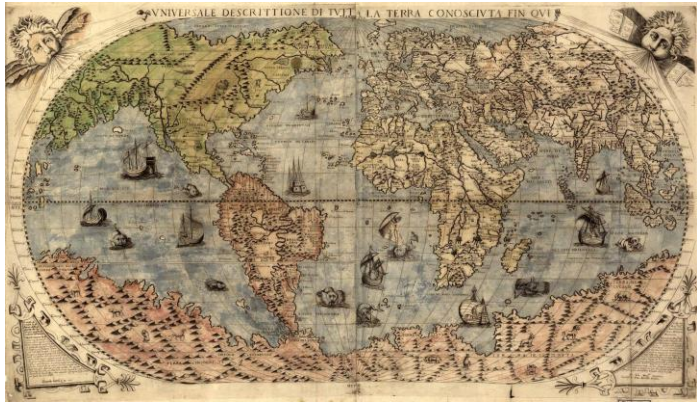
(Zacharias Holland et al. 2018, arXiv:1806.01825)



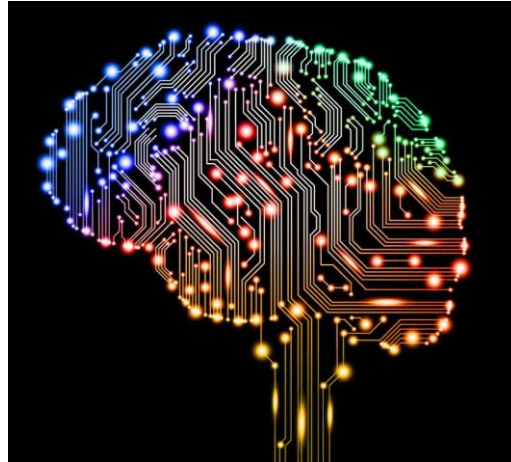
How does Dyna perform at Arcade games in a Deep Learning Setting?



# Important parts



Model

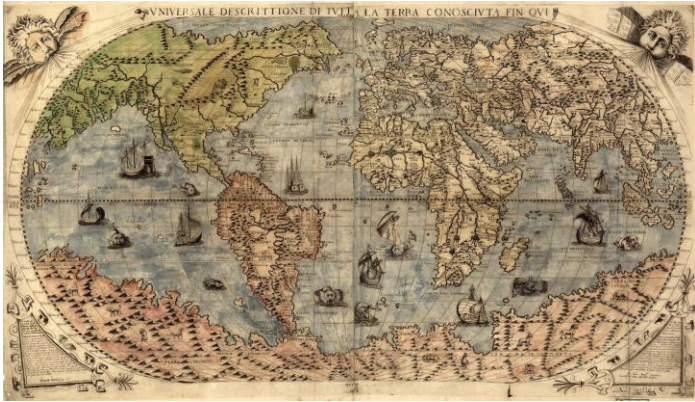


Learner

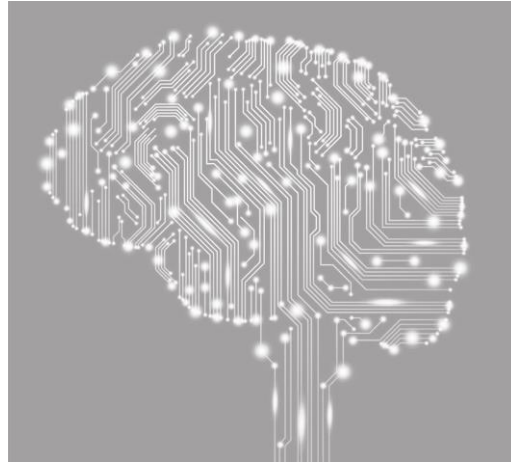


Roll-out Shapes

# Important parts



Model

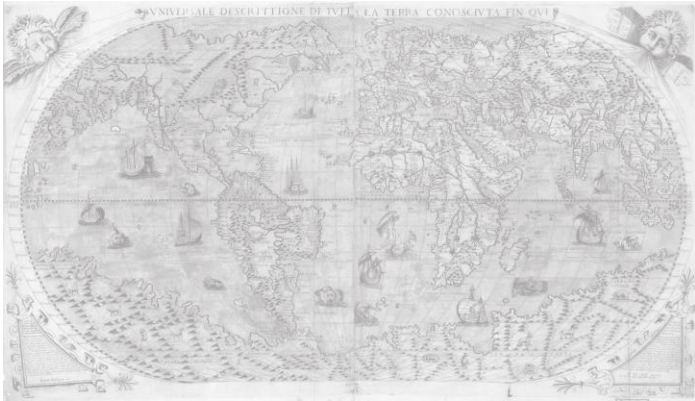


Learner

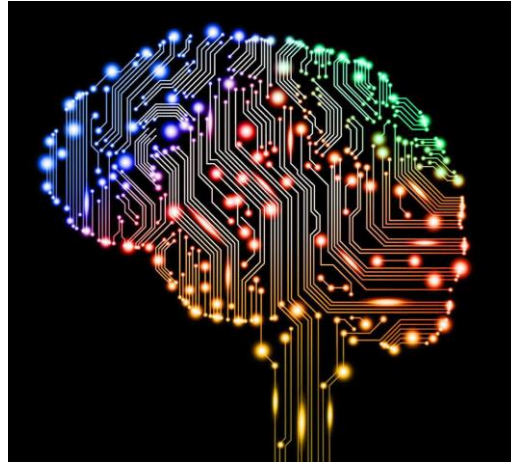


Roll-out Shapes

# Important parts



Model



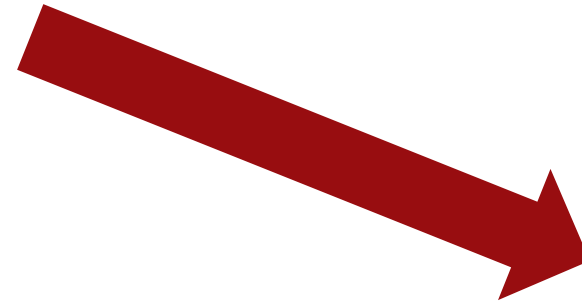
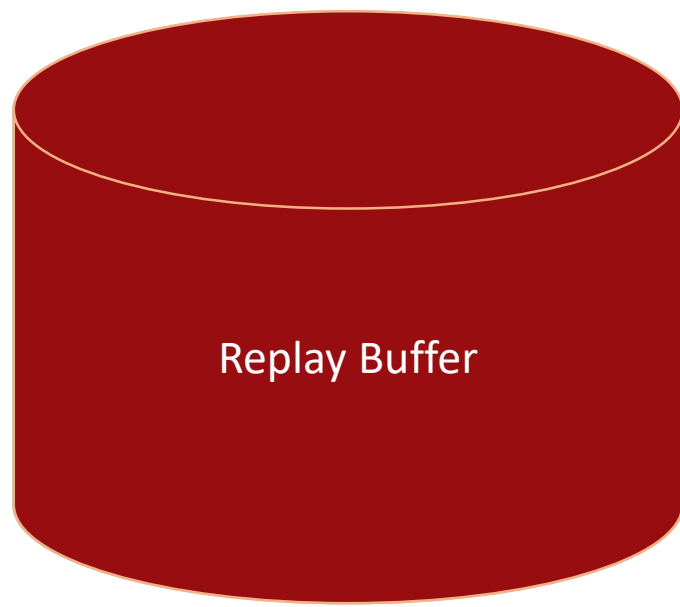
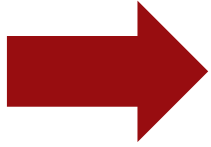
Learner



Roll-out Shapes



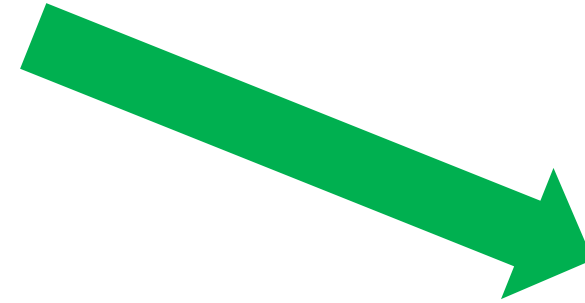
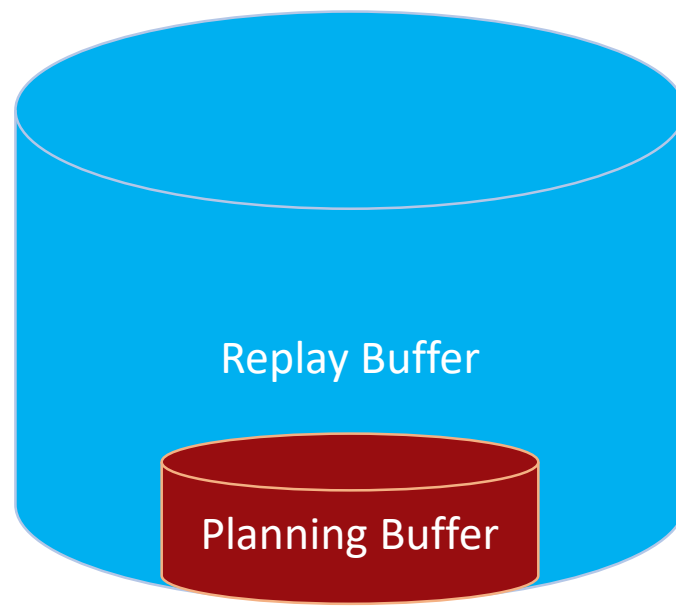
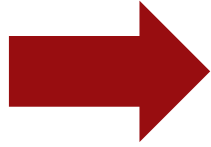
$\pi$



Train your Q-network



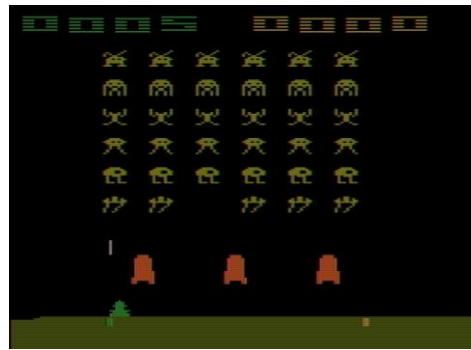
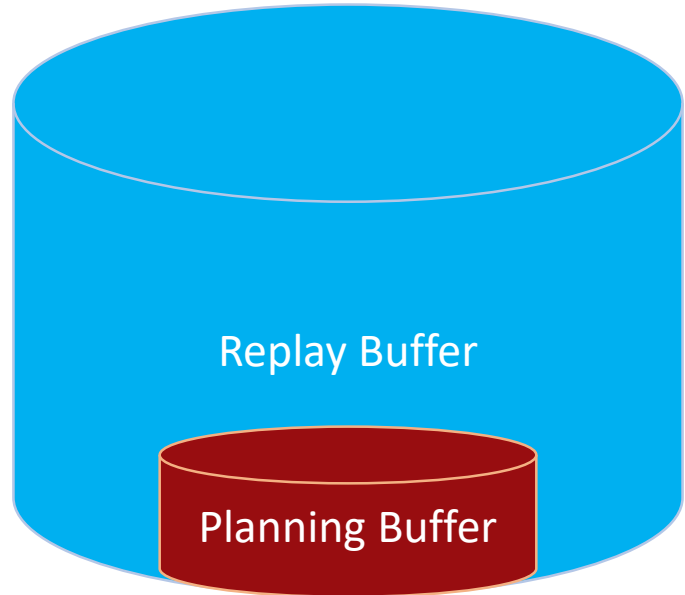
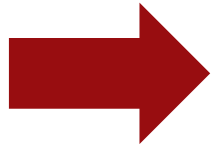
$\pi$



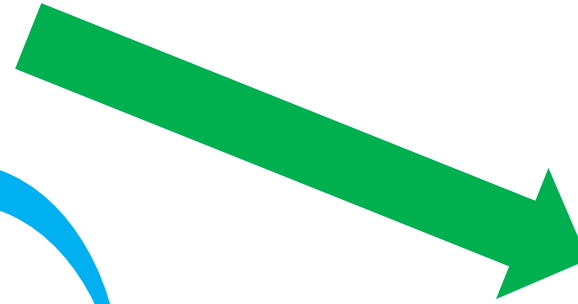
Train your Q-network



$\pi$

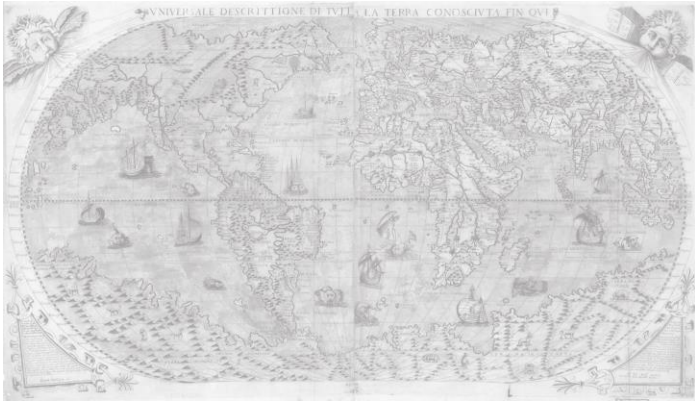


$\pi$   $M$

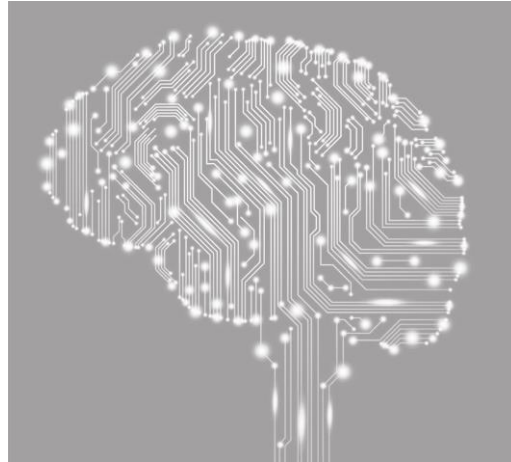


Train your Q-network

# Important parts



Model

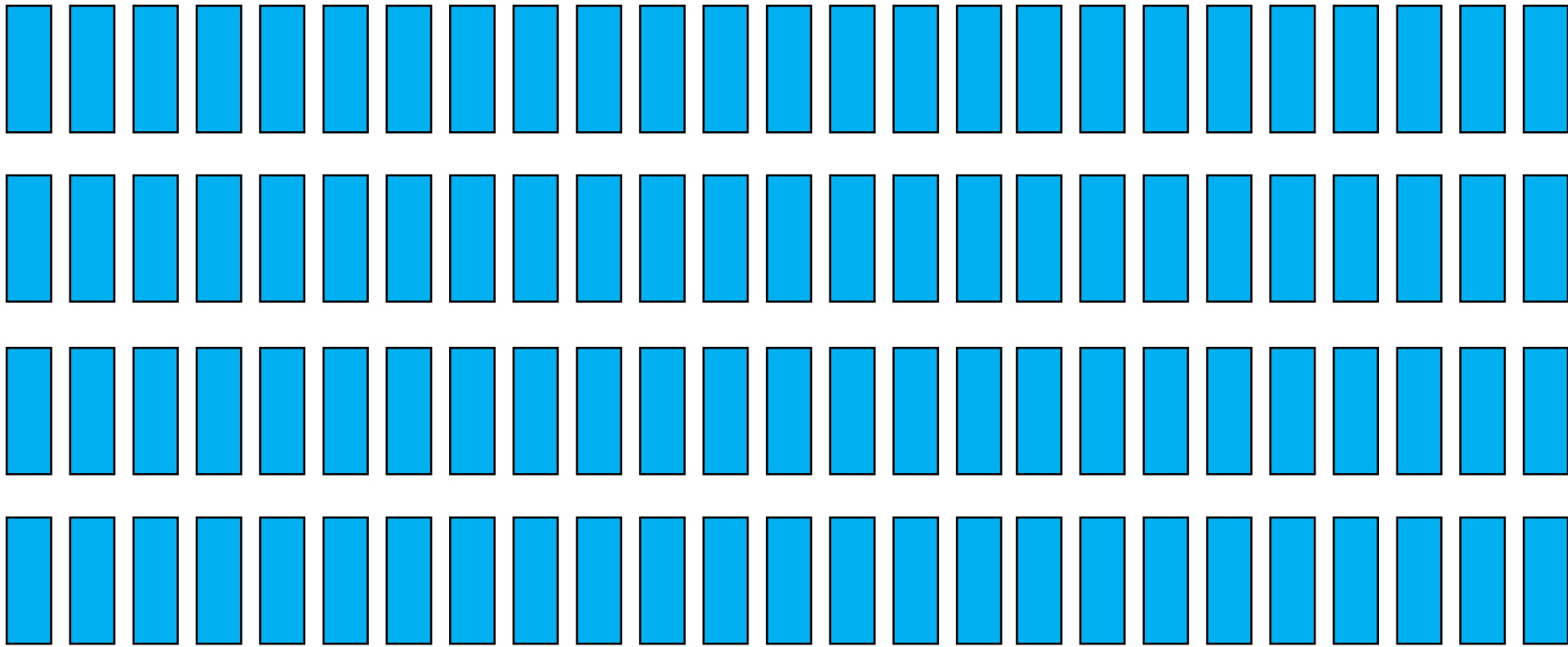


Learner



Roll-out Shapes

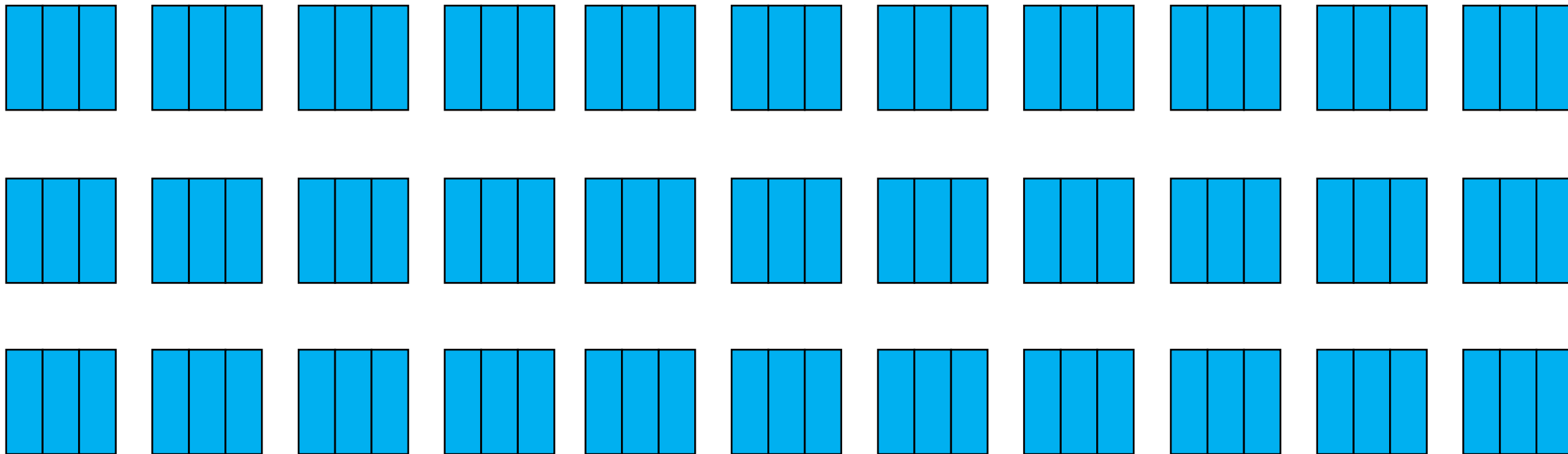
# Roll-out Shapes



100 x 1

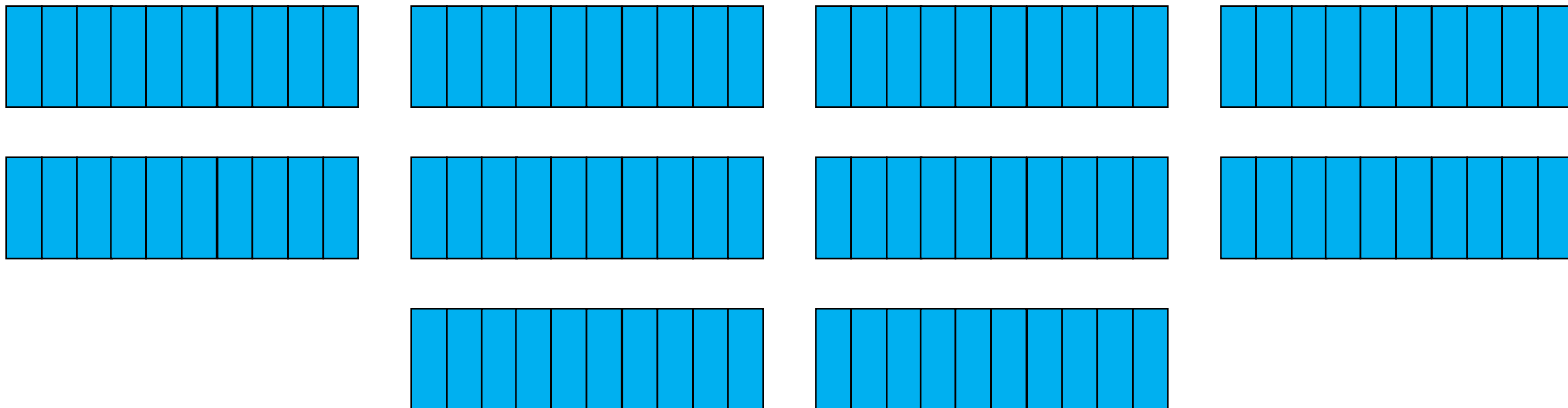


# Roll-out Shapes



33 x 3

# Roll-out Shapes

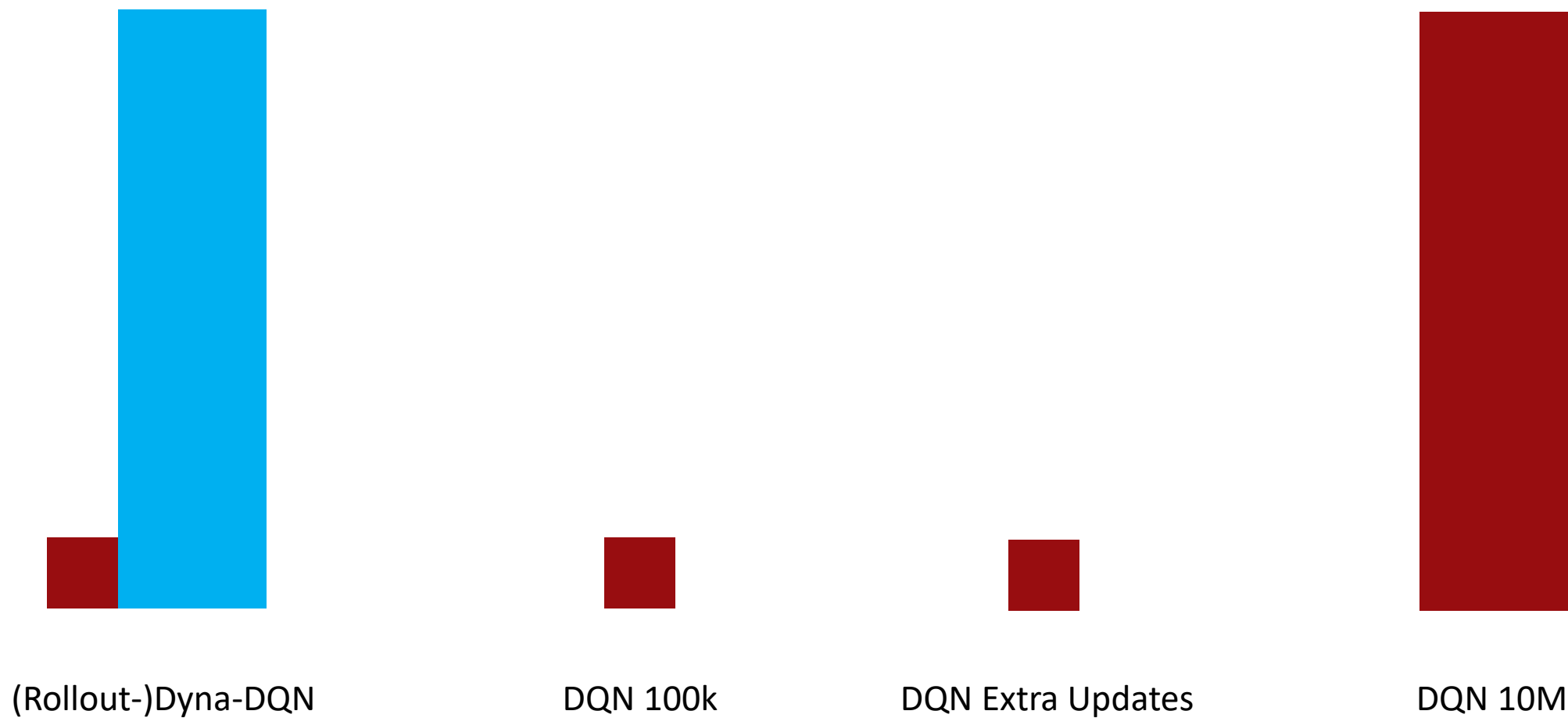


10 x 10

# Experiments

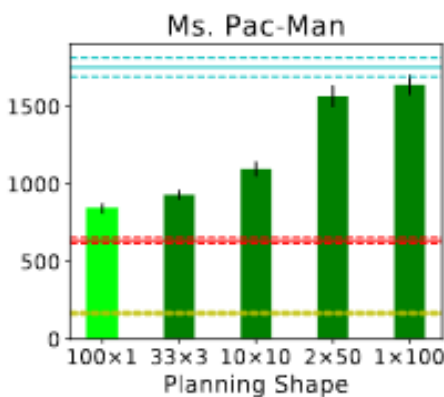
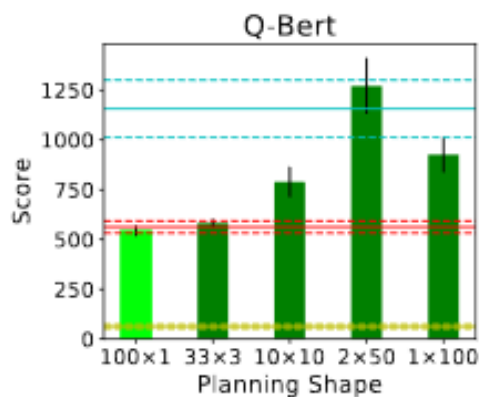
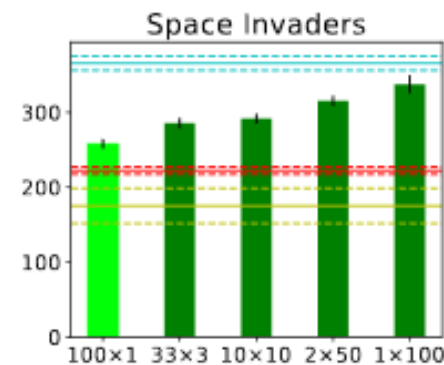
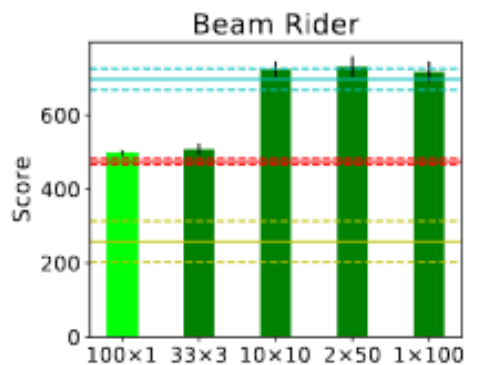
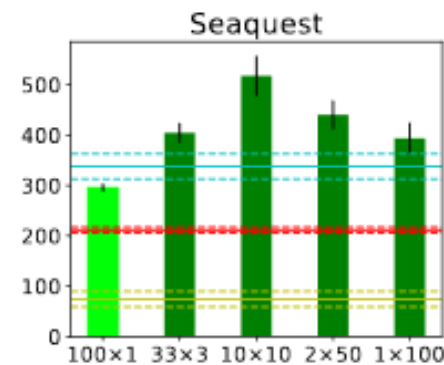
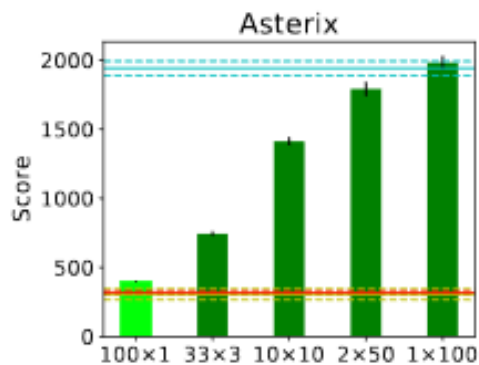


# Number of Samples & Benchmarks



■ DQN 100k   
 ■ DQN Extra Updates   
 ■ DQN 10M

■ Dyna-DQN   
 ■ Rollout-Dyna-DQN



# Relaxing the model

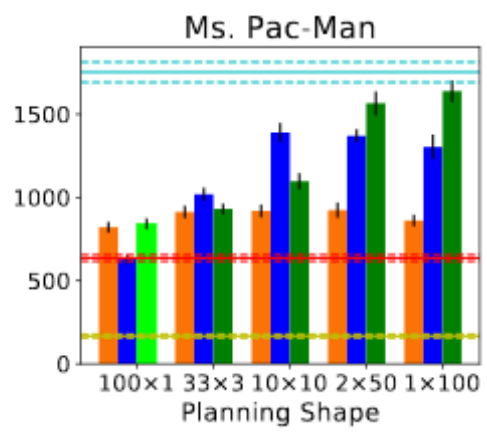
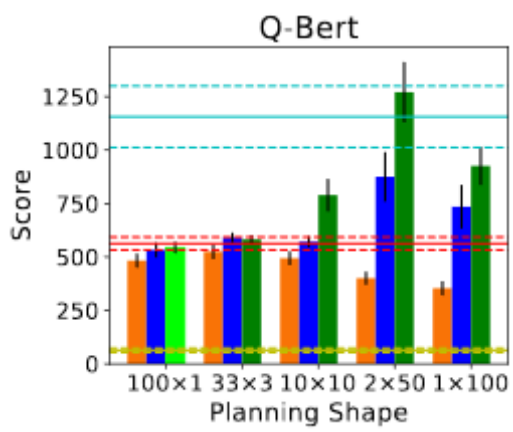
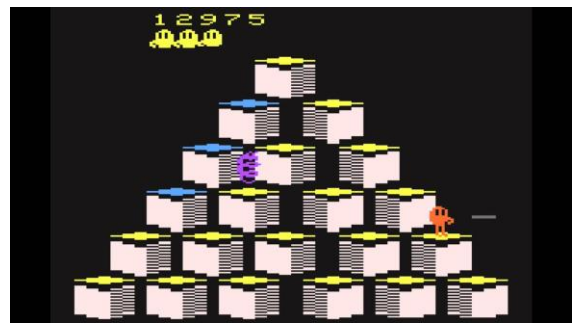
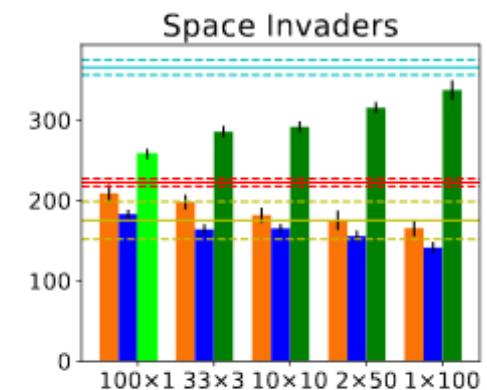
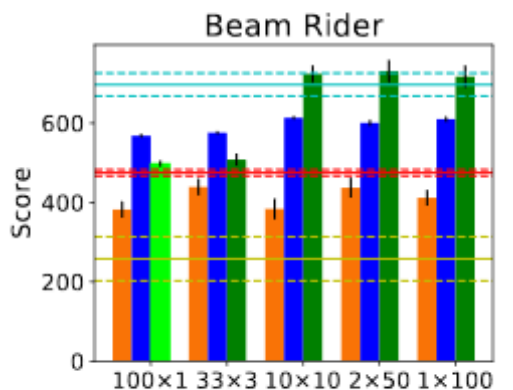
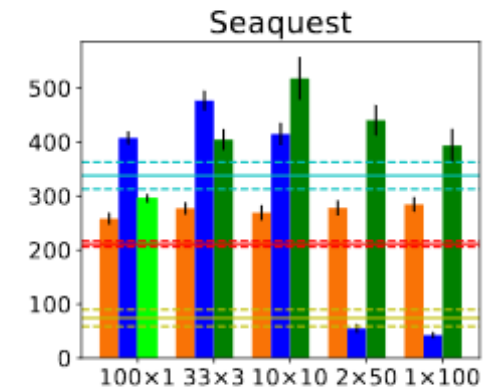
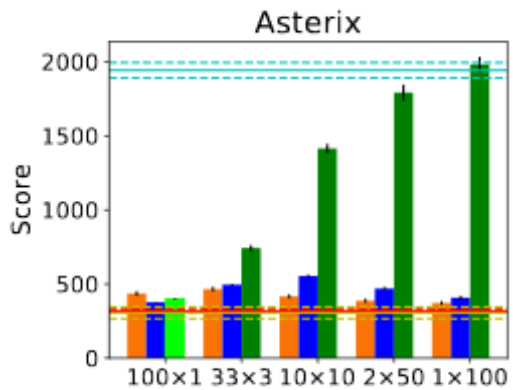


One model is pretrained  
on expert data



One model learns in an online  
fashion

Online learned (orange) Pre-trained (blue) Perfect (green)



# World Models

(Ha & Schmidhuber 2018, arXiv:1803.10122)



This thought bubble is our goal

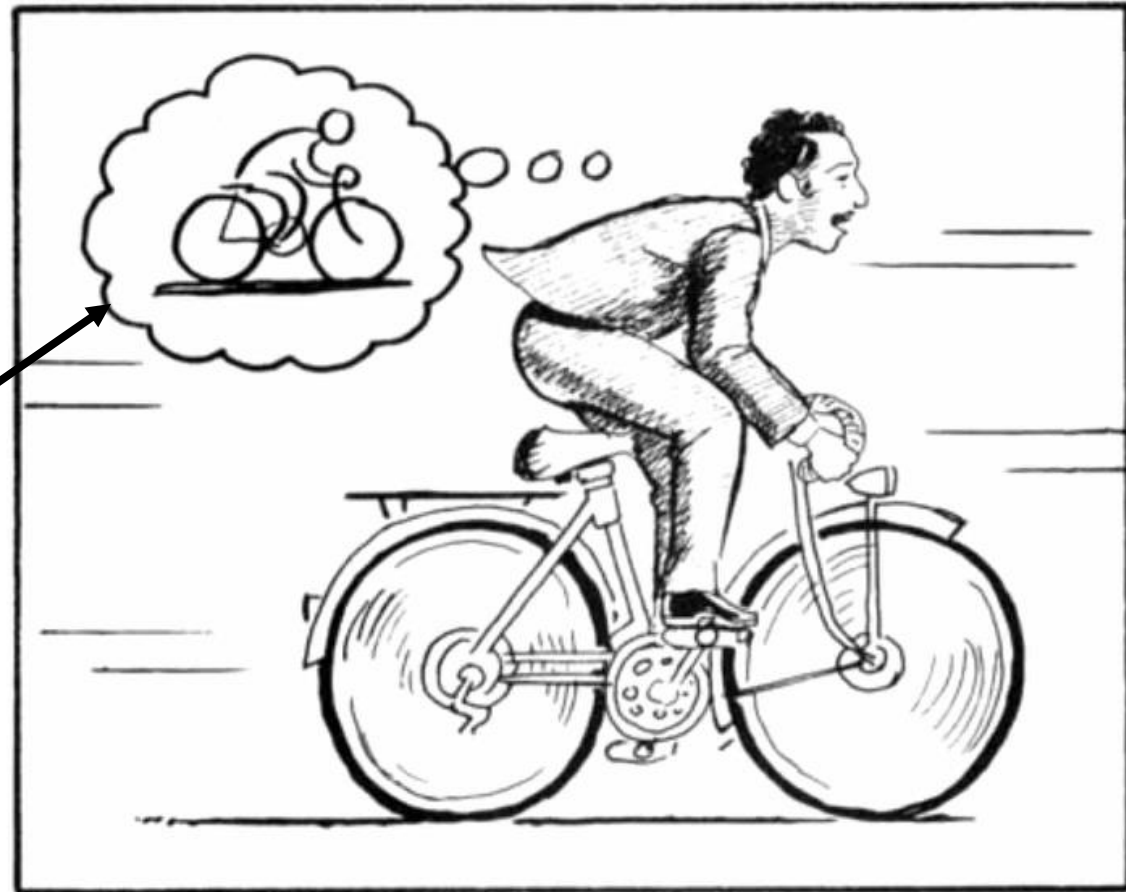
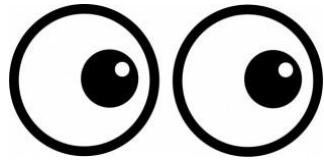


Figure 1. A World Model, from Scott McCloud's *Understanding Comics*. (McCloud, 1993; E, 2012)

# Learner over all

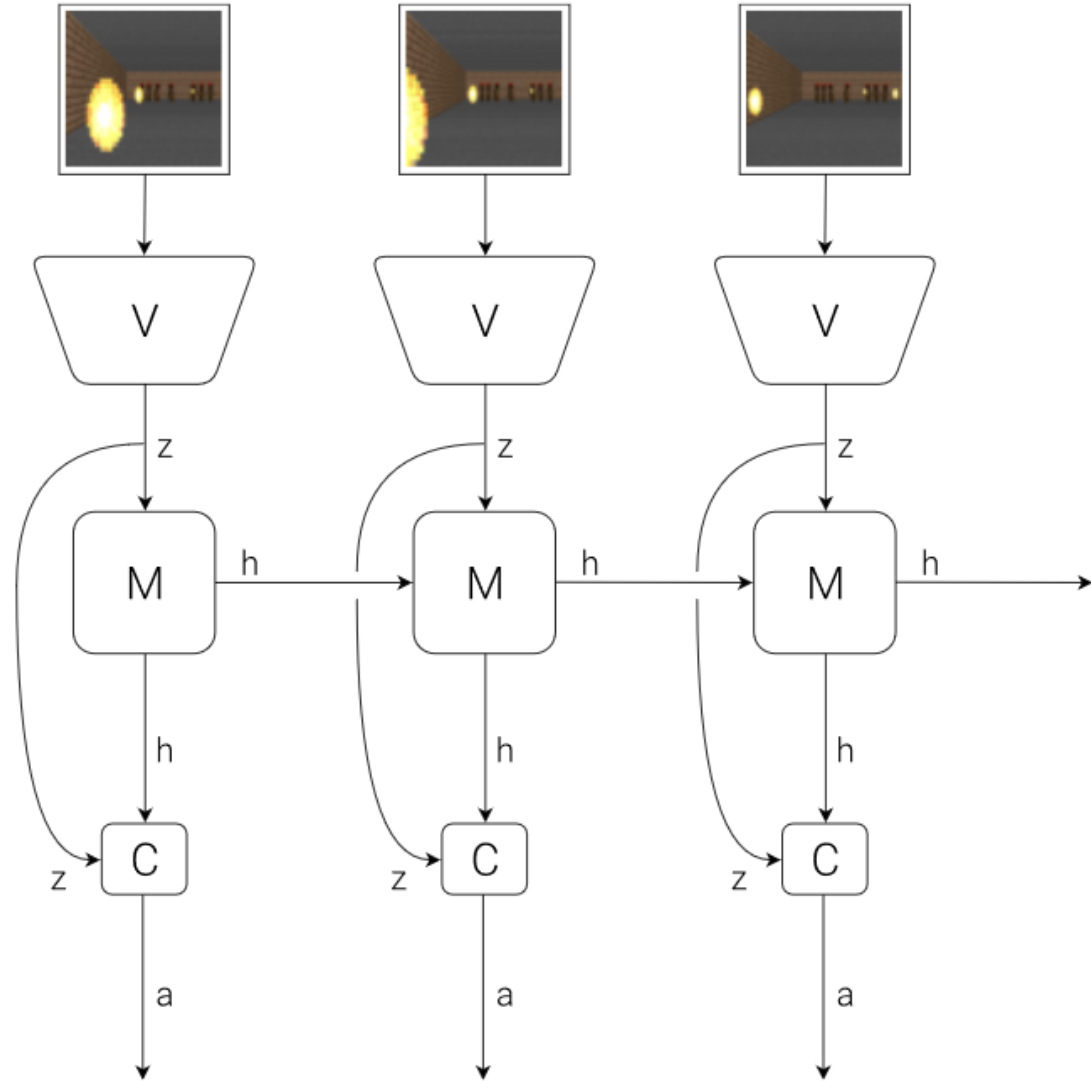
Vision



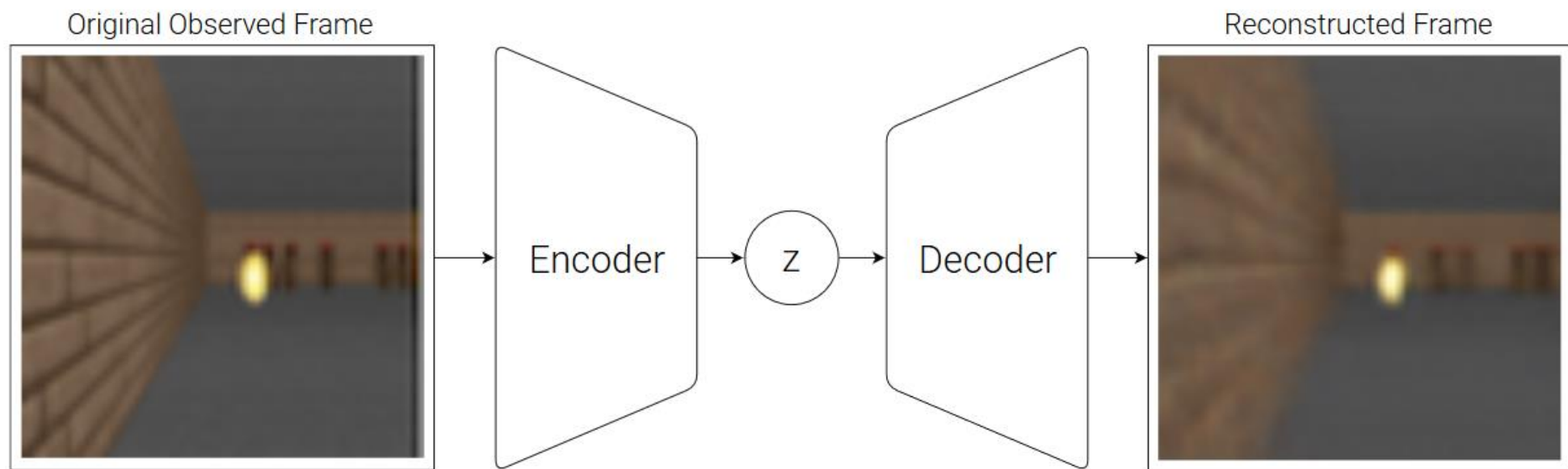
Memory



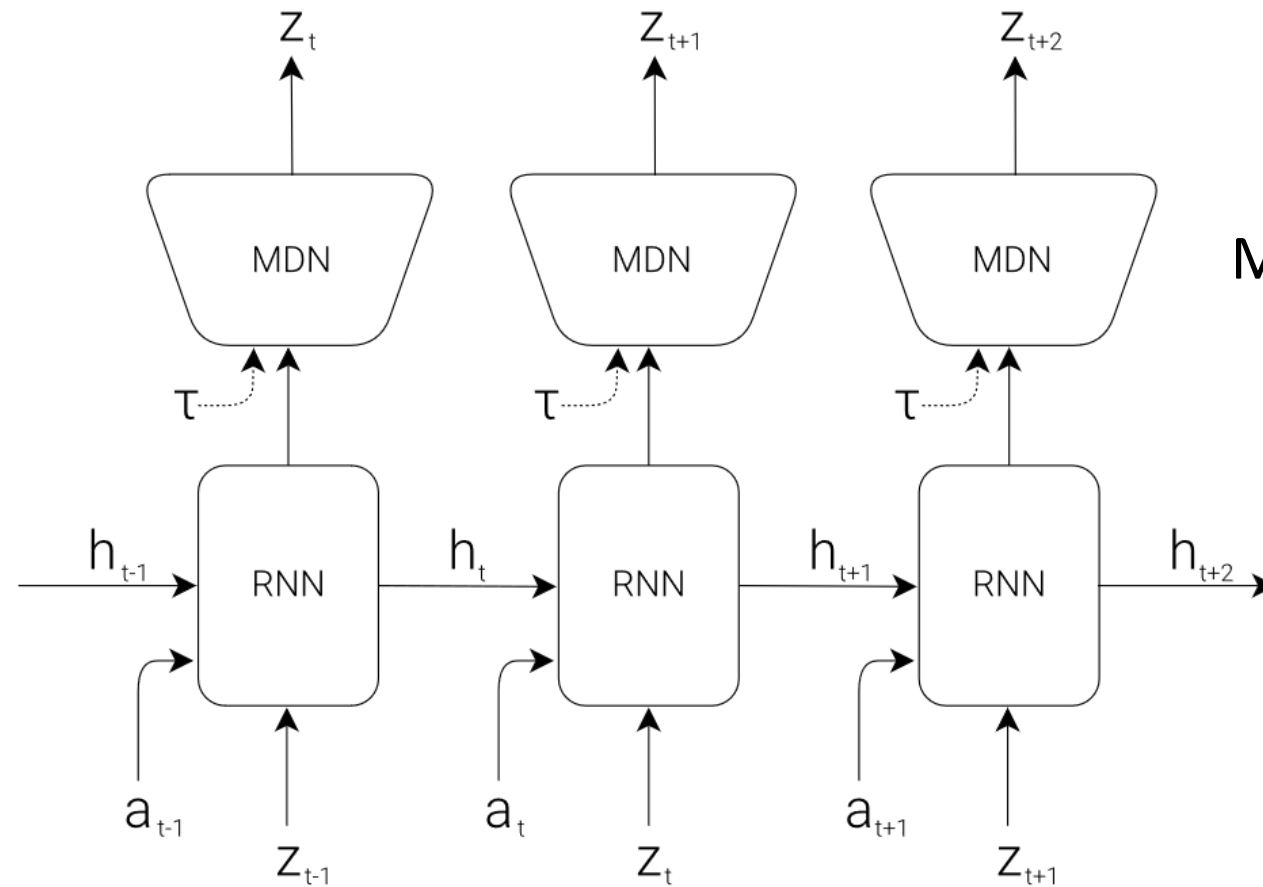
Controller



# Vision (V) Model



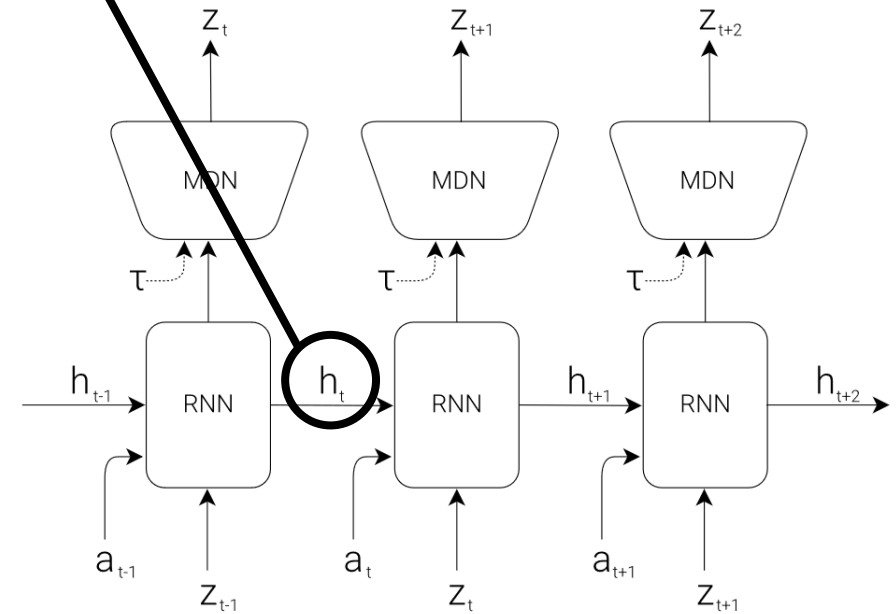
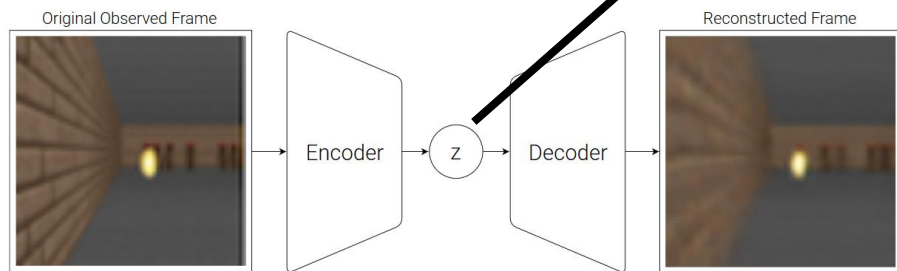
# Memory (M) Model



Mixture Density Network

# Controller (C) Model

$$a_t = W_c [z_t \ h_t] + b_c$$



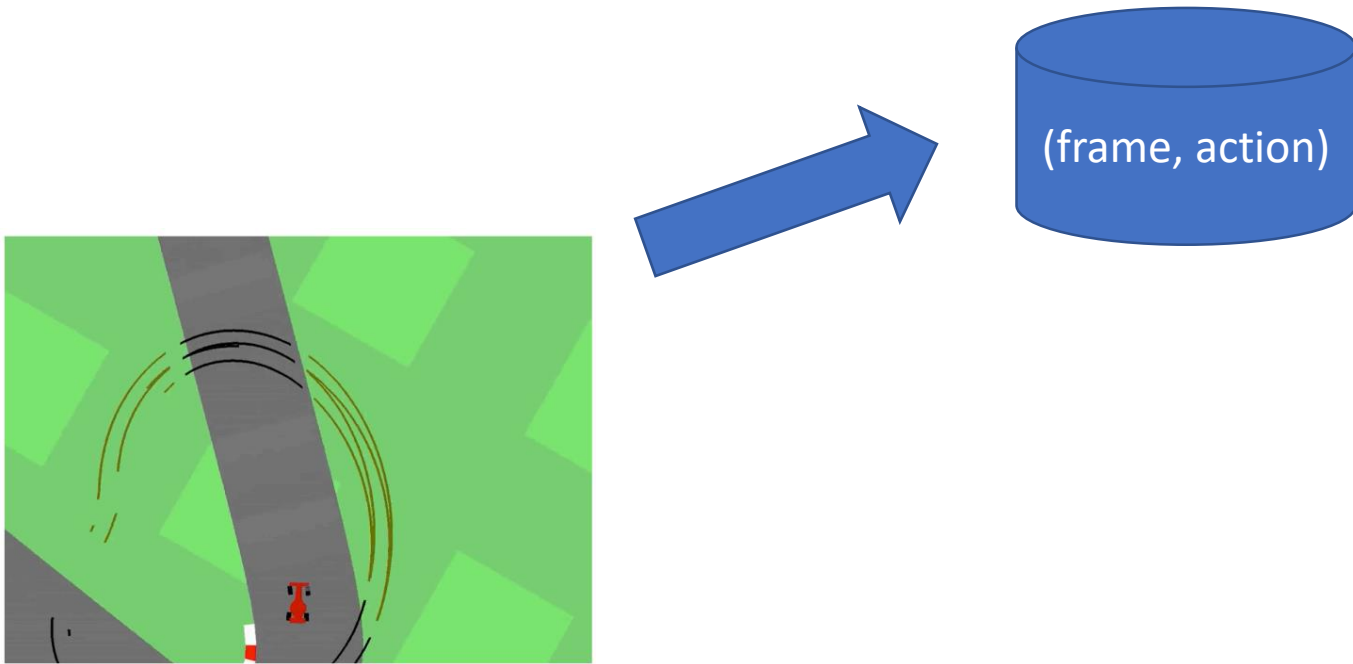
# Experiments



# Experiment 1 – Training the model

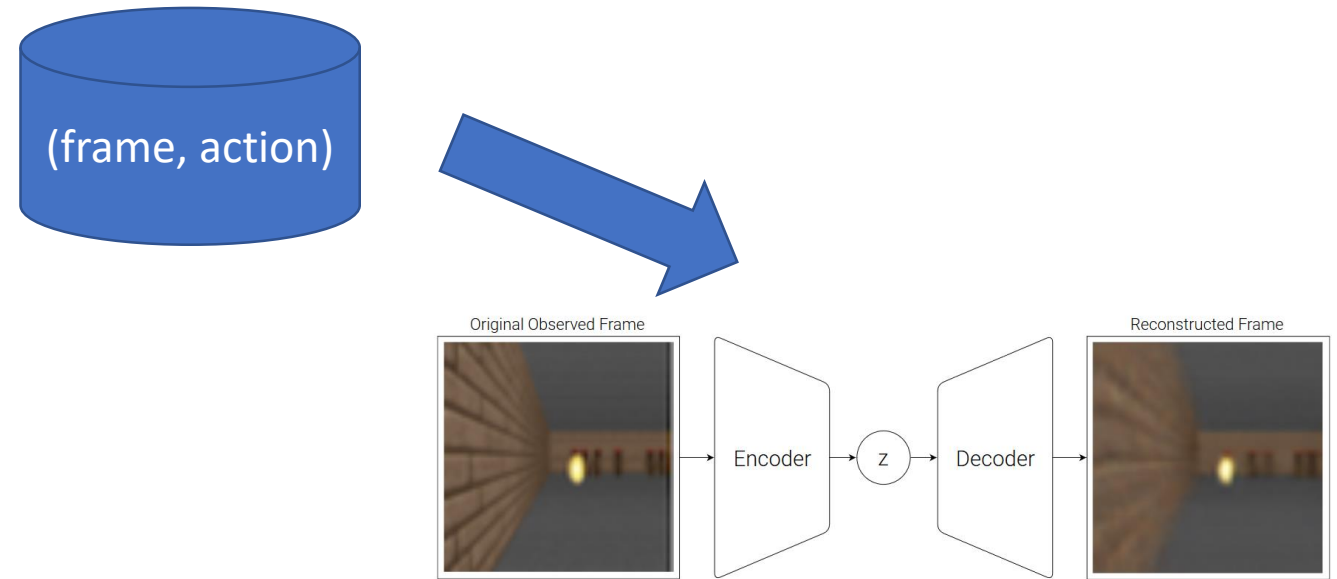


# Experiment 1 – Training the model

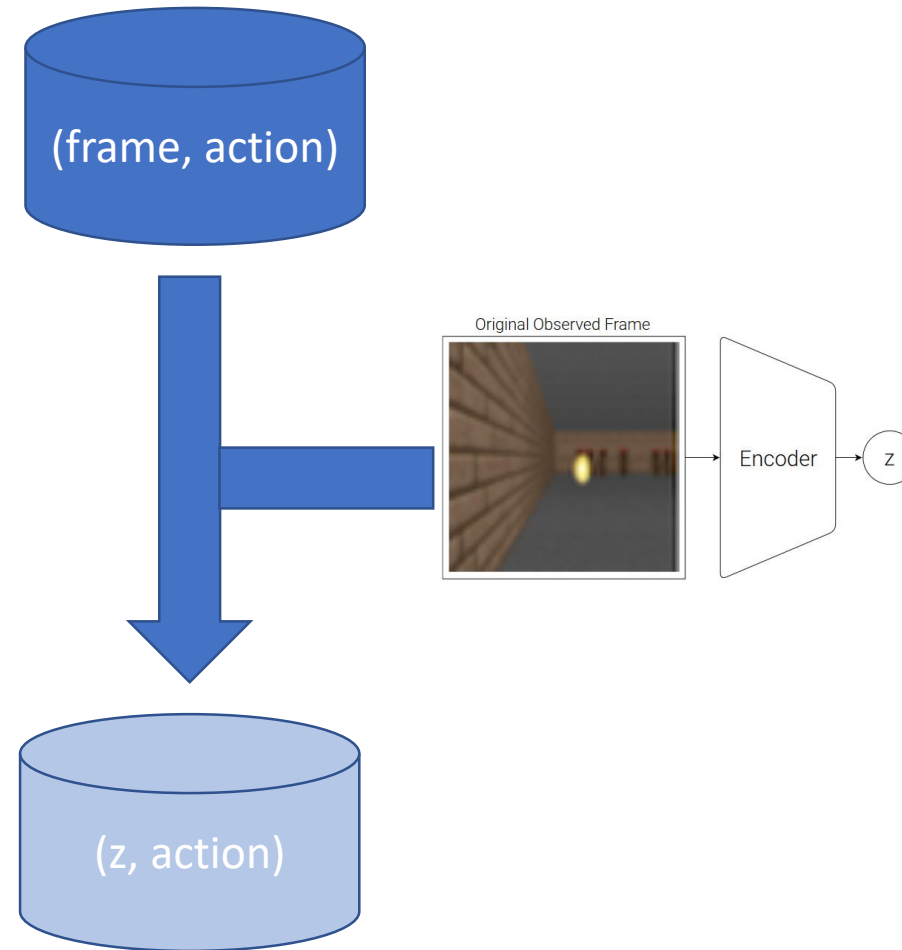




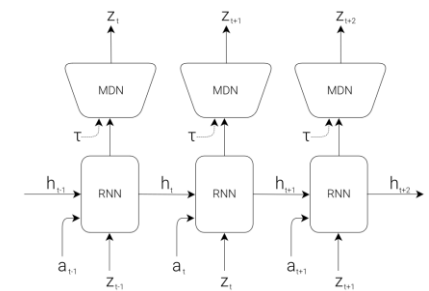
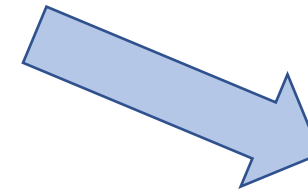
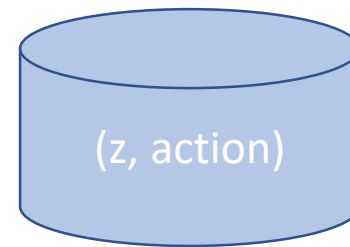
# Experiment 1 – Training the model



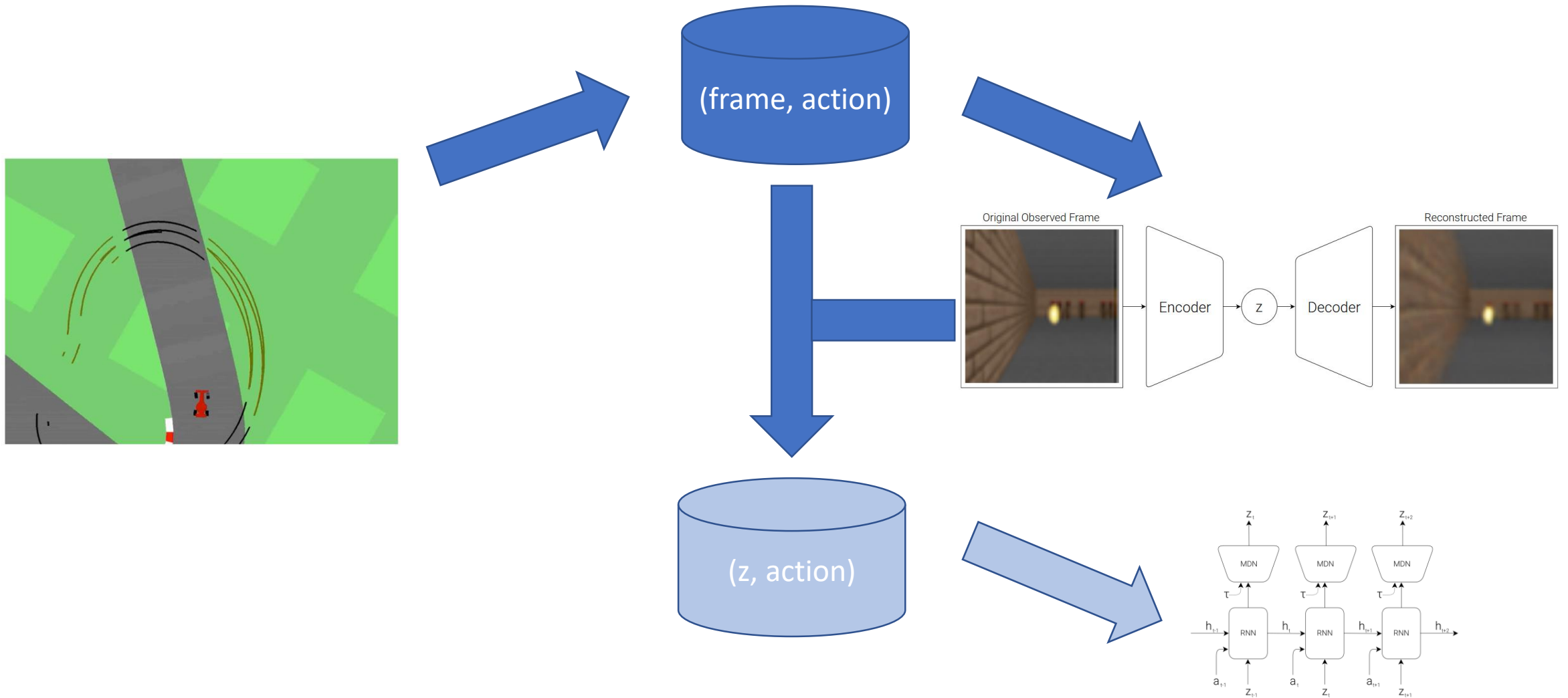
# Experiment 1 – Training the model



# Experiment 1 – Training the model



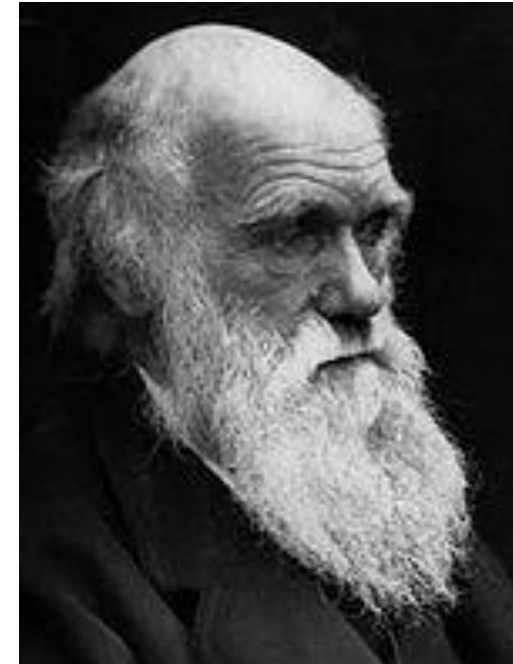
# Experiment 1 – Training the model



# Experiment 1 – Training the Controller

$$a_t = W_c [z_t \ h_t] + b_c$$

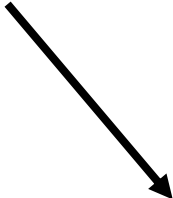
Evolution!




# Experiment 1

METHOD	AVG. SCORE
DQN (PRIEUR, 2017)	343 ± 18
A3C (CONTINUOUS) (JANG ET AL., 2017)	591 ± 45
A3C (DISCRETE) (KHAN & ELIBOL, 2016)	652 ± 10
CEOBILLIONAIRE (GYM LEADERBOARD)	838 ± 11
V MODEL	632 ± 251
V MODEL WITH HIDDEN LAYER	788 ± 141
<b>FULL WORLD MODEL</b>	<b>906 ± 21</b>

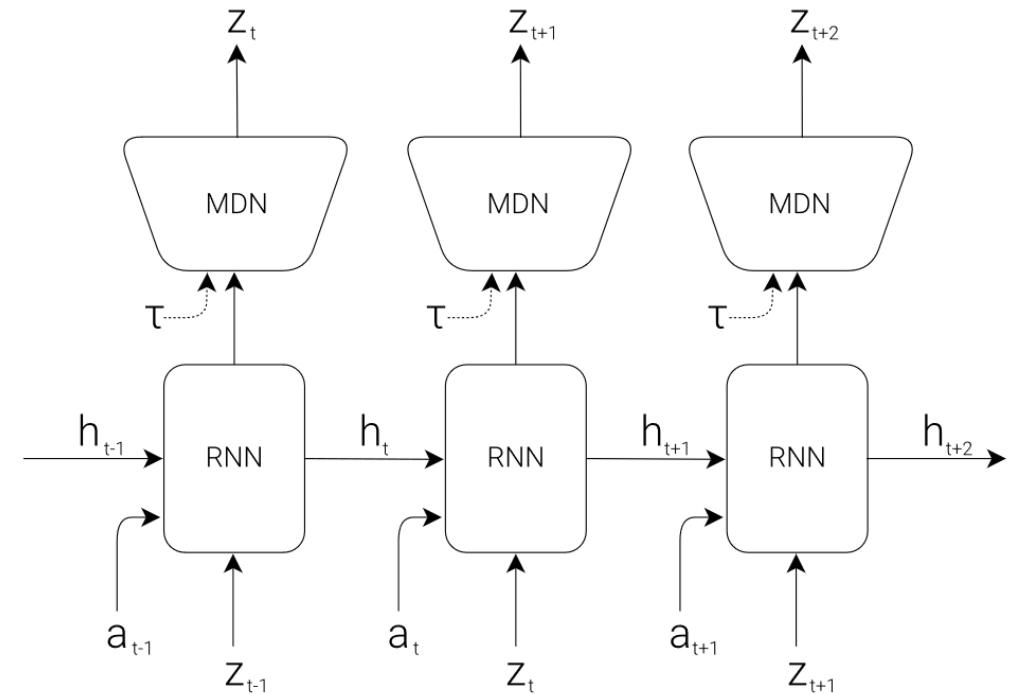
Removed memory  
Model



Removed memory  
Model, but with  
hidden layer in  
Controller



# Experiment 2 – train in dream world



# Experiment 2 – what about rewards?



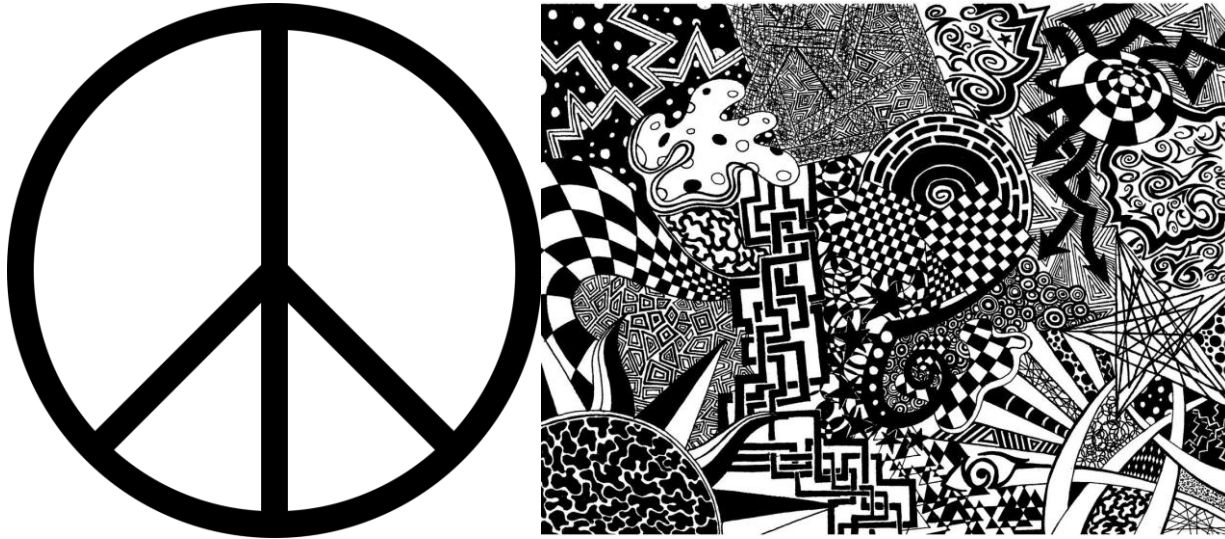
Maximize survival time

Train MDN-RNN (M) to model  $P(z_{t+1}, \underline{d_{t+1}} \mid a_t, z_t, h_t)$ .

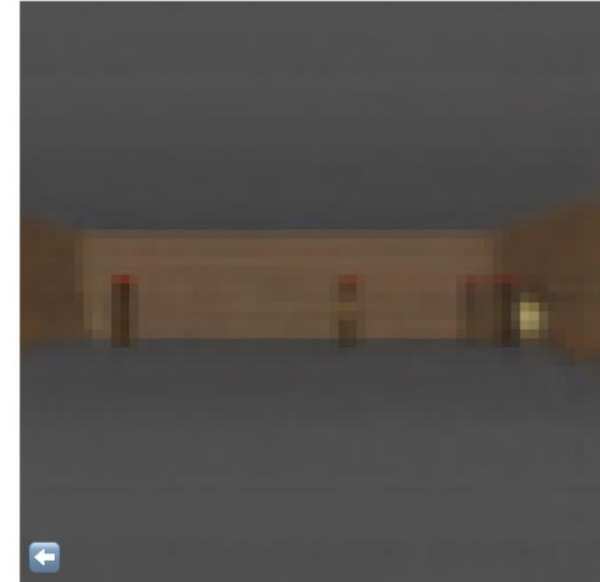
Memory model predicts death



# Experiment 2 – Problems



Too low temperature -> no shooting  
Too high temperature -> chaos



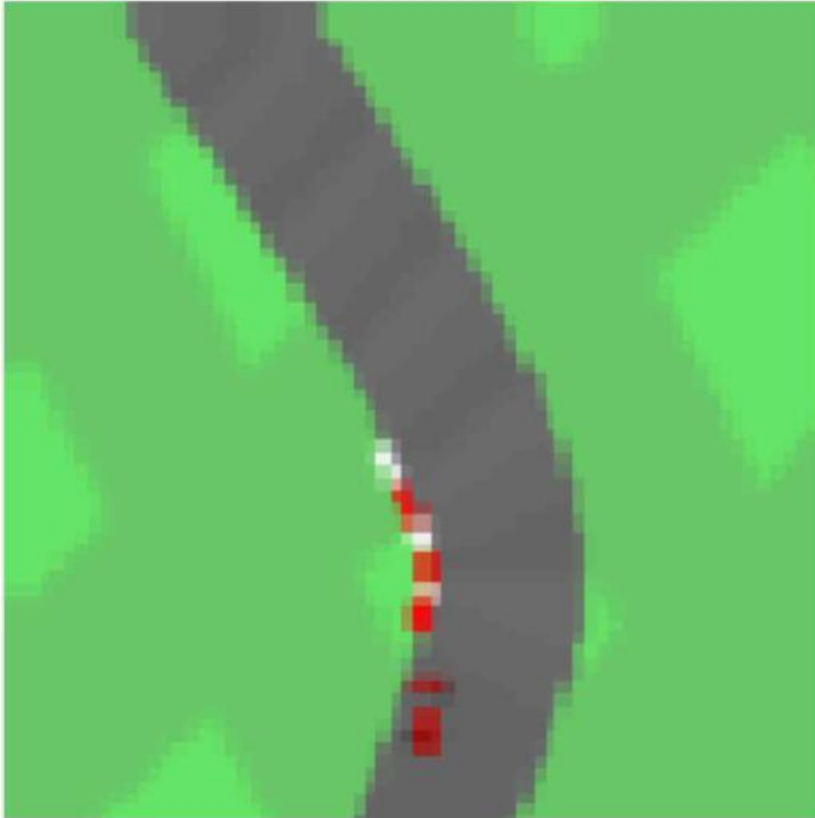
Controller exploits inaccuracies  
of the model

# Experiment 2

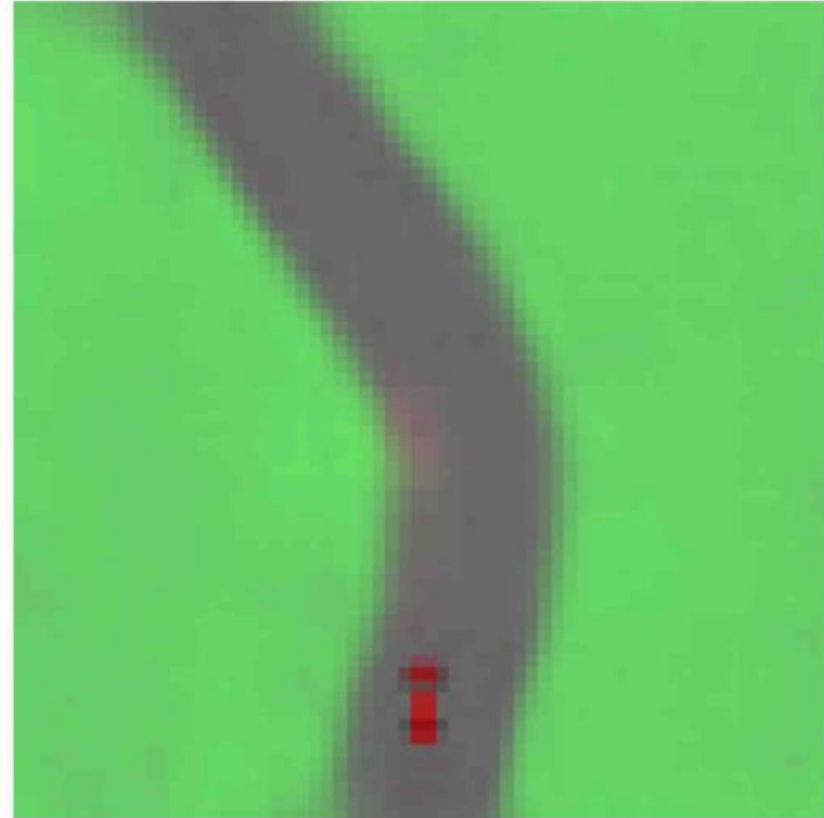
TEMPERATURE $\tau$	VIRTUAL SCORE	ACTUAL SCORE
0.10	2086 $\pm$ 140	193 $\pm$ 58
0.50	2060 $\pm$ 277	196 $\pm$ 50
1.00	1145 $\pm$ 690	868 $\pm$ 511
1.15	918 $\pm$ 546	1092 $\pm$ 556
1.30	732 $\pm$ 269	753 $\pm$ 139
RANDOM POLICY	N/A	210 $\pm$ 108
GYM LEADER	N/A	820 $\pm$ 58

# World Models – Shortcoming

Screenshot Image



Reconstruction



Latent representation  $z$  optimized for reconstruction and not for task solving

See you on Piazza