# When to use parametric models in reinforcement learning?

Hado van Hasselt, DeepMind

Tea-time talk, Montreal, July 11, 2019
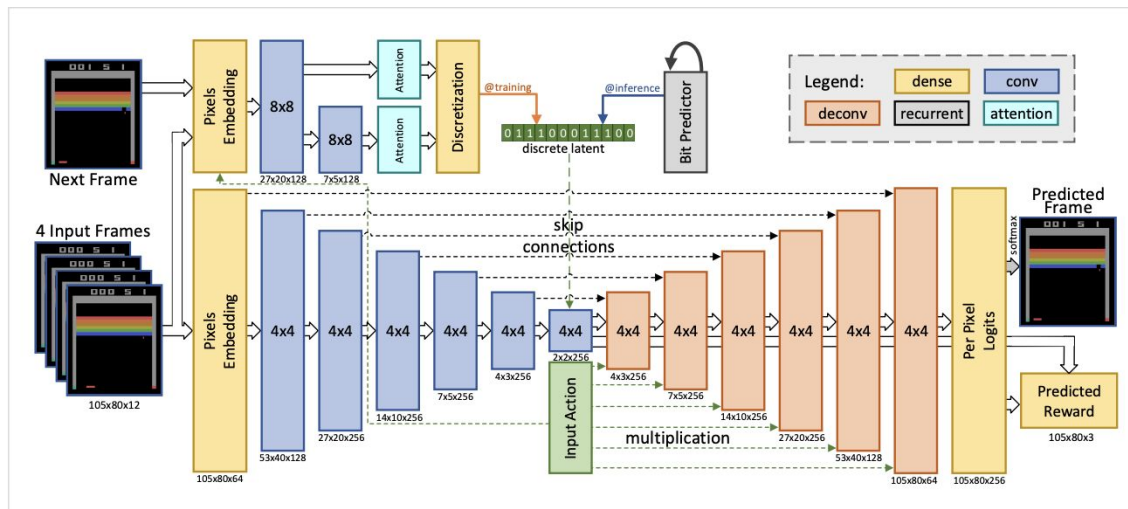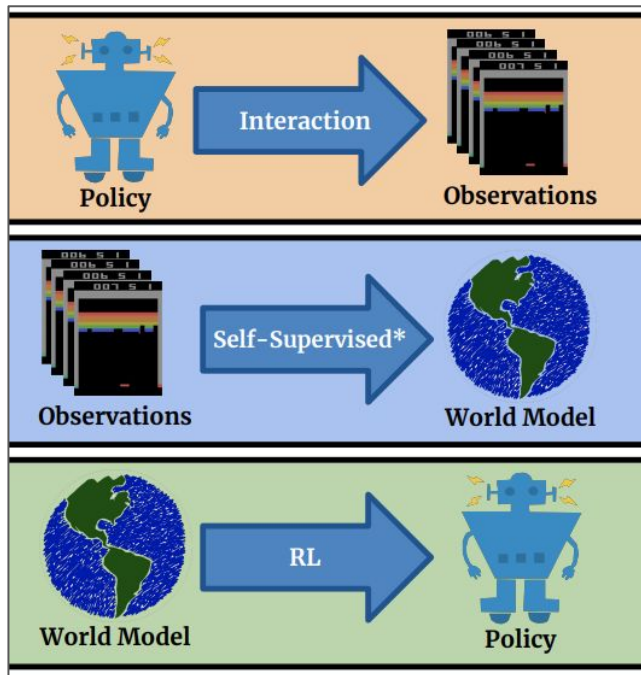
Related paper: https://arxiv.org/abs/1906.05243

# Motivation

## Model Based Reinforcement Learning for Atari

Łukasz Kaiser [*1]   Mohammad Babaeizadeh [*23]   Piotr Miłoś [*45]   Błażej Osiński [*453]
Roy H Campbell [2]   Konrad Czechowski [4]   Dumitru Erhan [1]   Chelsea Finn [1]   Piotr Kozakowski [4]   Sergey Levine [1]
Afroz Mohiuddin [1]   Ryan Sepassi [1]   George Tucker [1]   Henryk Michalewski [45]
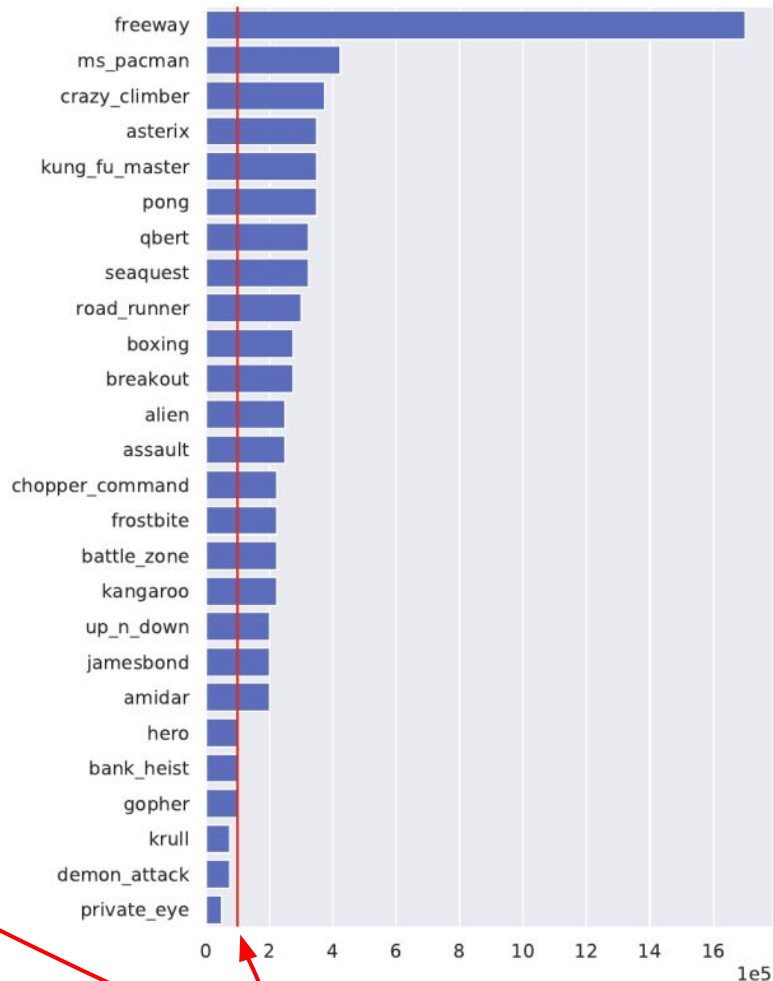
# Motivation

# Motivation

Performance (=score)
after 100,000 steps (=400,000 frames)

Baseline: Rainbow DQN

# Question

Why does the parametric model perform better than replay?

# Models and planning

A **model** is a function:

$$r, \ s' \ = \ m(s, a)$$

We can use models to **plan**: spend more compute to improve prediction & policies.

We can also plan with **experience replay**:

$$r\_\{n+1\}, \ s\_\{n+1\} \ = \ replay(s\_n, a\_n)$$

- Experience replay is similar to a non-parametric model
- But we can only query it at observed state action pairs $(s\_n, a\_n)$, $n < t$.

# Replay and models, properties

Typically models use **less memory** and **more compute** than replay

But what about **data efficiency & performance?**

# Algorithms

$$\textbf{for } \text{iteration} \in \{1, 2, \ldots, K\} \textbf{ do}$$
$$\quad \textbf{for } \text{interaction} \in \{1, 2, \ldots, M\} \textbf{ do}$$
$$\quad\quad \text{Generate action: } a \leftarrow \pi(s)$$
$$\quad\quad \text{Generate reward, next state: } r, s' \leftarrow \mathcal{E}(a)$$
$$\quad\quad m, d \leftarrow \textsc{UpdateModel}(s, a, r, s')$$
$$\quad\quad \pi, v \leftarrow \textsc{UpdateAgent}(s, a, r, s')$$
$$\quad\quad \text{Update current state: } s \leftarrow s'$$
$$\quad \textbf{end for}$$
$$\quad \textbf{for } \text{planning step} \in \{1, 2, \ldots, P\} \textbf{ do}$$
$$\quad\quad \text{Generate state, action } \tilde{s}, \tilde{a} \leftarrow d$$
$$\quad\quad \text{Generate reward, next state: } \tilde{r}, \tilde{s}' \leftarrow m(\tilde{s}, \tilde{a})$$
$$\quad\quad \pi, v \leftarrow \textsc{UpdateAgent}(\tilde{s}, \tilde{a}, \tilde{r}, \tilde{s}')$$
$$\quad \textbf{end for}$$

**K**: iterations
**M**: real steps / iteration
**P**: planned steps / iteration

State-sampling distribution

# Algorithms

| | Iterations (**K**) | Real steps per iteration (**M**) | Planned steps per iteration (**P**) |
|---|---|---|---|
| SimPLe | 16 | 6400 | 800,000 |
| Rainbow DQN | 12,500,000 | 4 | 32 |
| Data-efficient Rainbow DQN | 100,000 | 1 | 32 |

# Algorithms

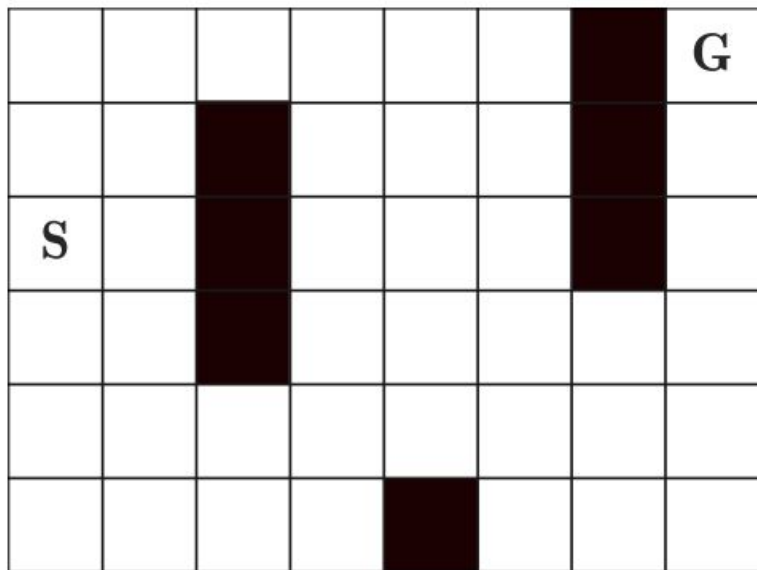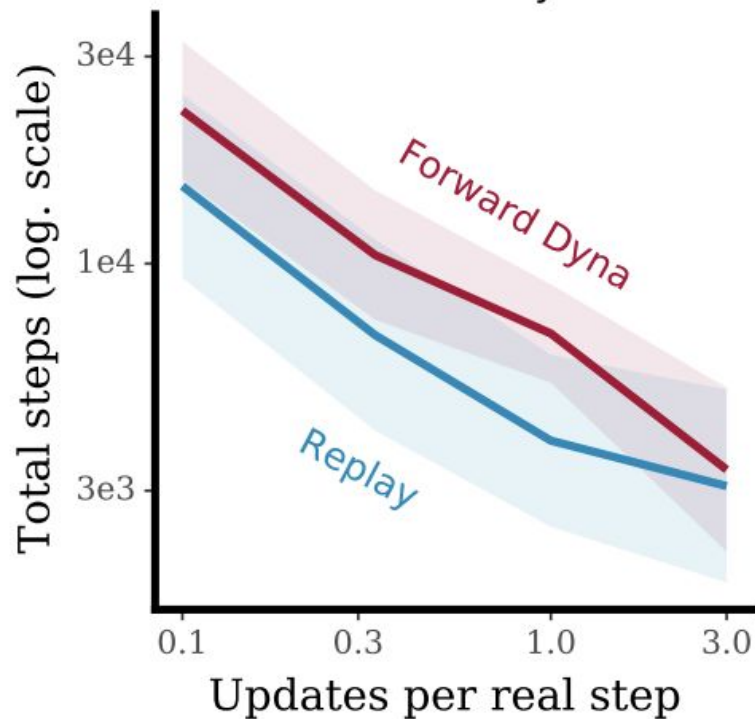| | Total real experience (**K** x **M**) | Total planned experience (**K** x **P**) |
|---|---|---|
| SimPLe | **100,000**<br>(400K frames) | **15,200,000** |
| Rainbow DQN | **50,000,000**<br>(200M frames) | **400,000,000** |
| Data-efficient Rainbow DQN | **100,000**<br>(400K frames) | **3,200,000** |

# Algorithms

# When do models help performance?

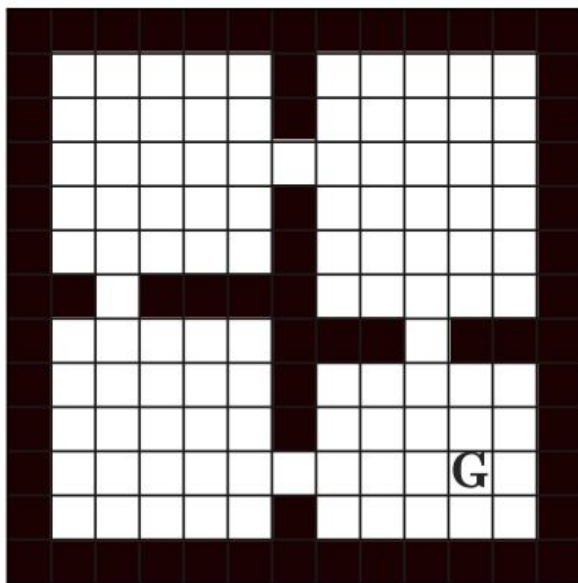# Forward planning for credit assignment
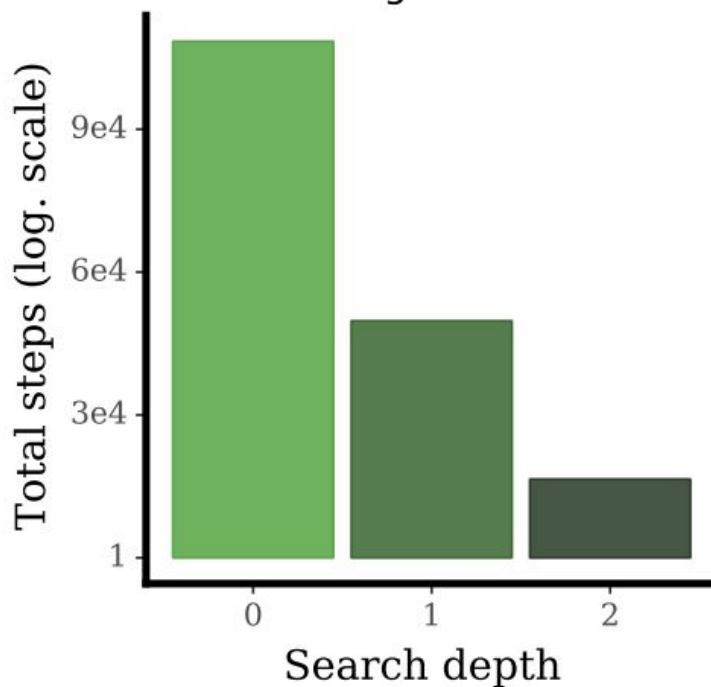


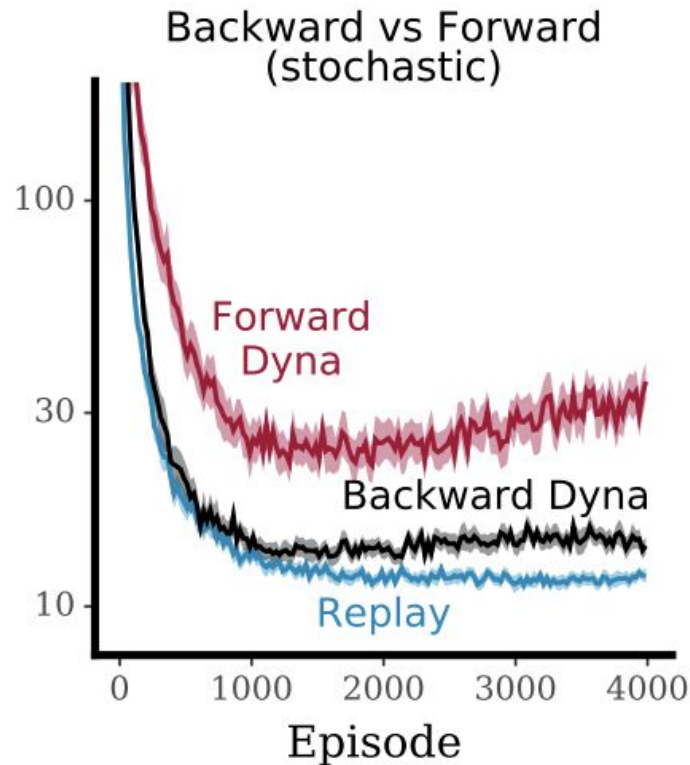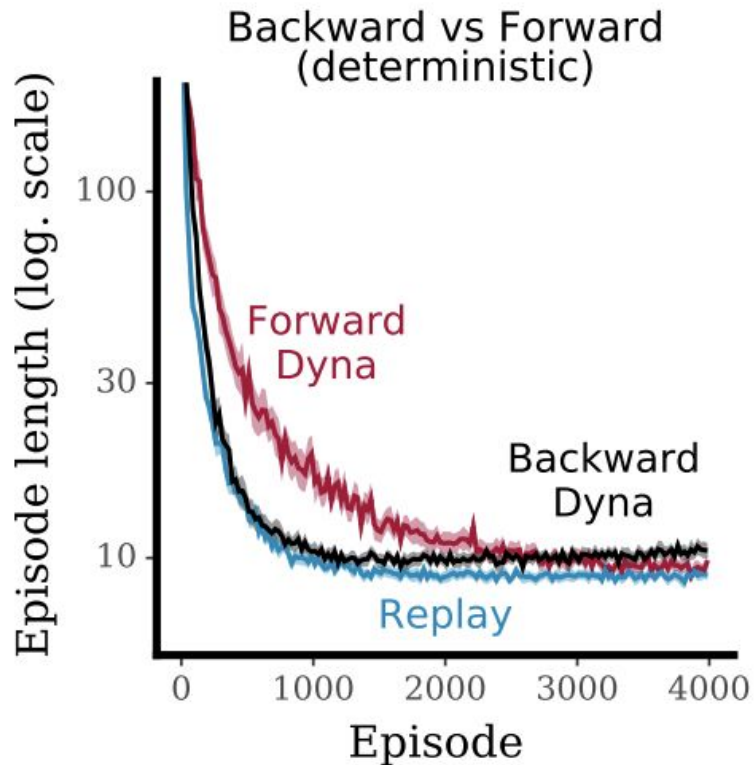The maze



Scalability

# Forward planning for behaviour



Four rooms

Planning in the now

# Backward planning for credit assignment



Backward vs Forward (deterministic)
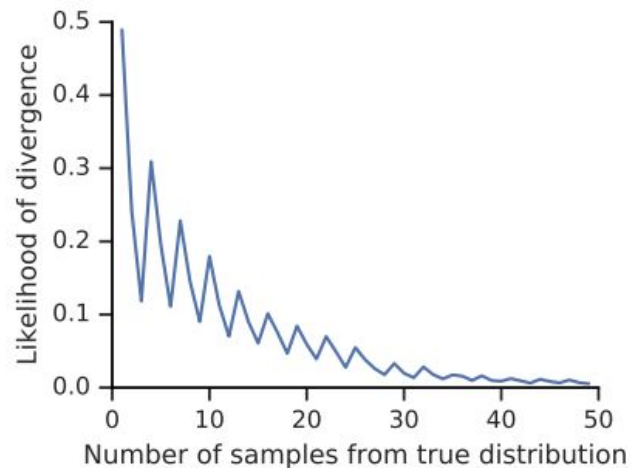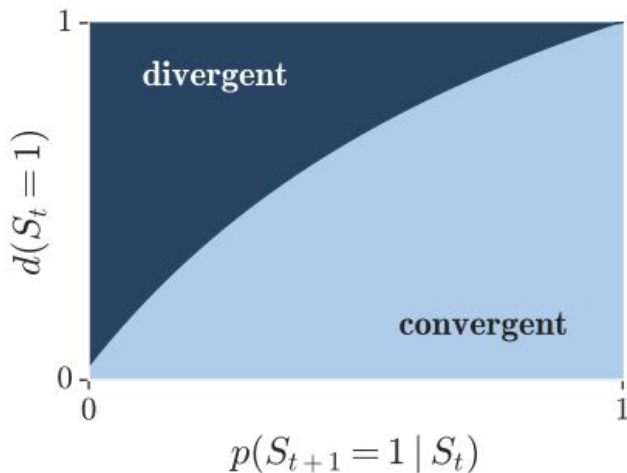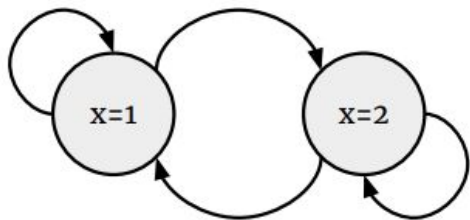
Backward vs Forward (stochastic)

# Conclusions

1. **Replay** can be used for **planning**
2. There are different ways to use models:
   a. **Forward** planning     for **credit assignment**
   b. **Forward** planning     for **immediate behaviour**
   c. **Backward** planning    for **credit assignment**

# Thank you

More details: https://arxiv.org/abs/1906.05243

# Bonus slide: Surprising instabilities



- Even with **perfect models** learning can be unstable
- This happens surprisingly easily!
- Related to the *deadly triad*:
  - the state sampling distribution $d$ and the model $m$ may mismatch, even if $m = p$ is perfect.

# Bonus slide: Rainbow DQN hyperparameters changes

| Hyper-parameter | canonical | data-efficient |
| --- | --- | --- |
| Training frames | 200,000,000 | 400,000 |
| Min replay size for sampling | 20,000 | 1600 |
| Memory size | 1,000,000 steps | unbounded |
| Replay period every | 4 steps | 1 steps |
| Multi-step return length | 3 | 20 |
| Q network: channels | 32, 64, 64 | 32, 64 |
| Q network: filter size | $8 \times 8, 4 \times 4, 3 \times 3$ | $5 \times 5, 5 \times 5$ |
| Q network: stride | 4, 2, 1 | 5, 5 |
| Q network: hidden units | 512 | 256 |
| Optimizer: learning rate | 0.0000625 | 0.0001 |