

# Model-based Reinforcement Learning in Computer Systems

Sean J. Parker  
Clare Hall



UNIVERSITY OF  
CAMBRIDGE

*A dissertation submitted to the University of Cambridge  
in partial fulfilment of the requirements for the degree of  
Master of Philosophy in Advanced Computer Science*

University of Cambridge  
Computer Laboratory  
William Gates Building  
15 JJ Thomson Avenue  
Cambridge CB3 0FD  
UNITED KINGDOM

Email: [sjp240@cam.ac.uk](mailto:sjp240@cam.ac.uk)

April 10, 2021



# Declaration

I, Sean J. Parker of Clare Hall, being a candidate for the M.Phil in Advanced Computer Science, hereby declare that this report and the work described in it are my own work, unaided except as may be specified below, and that the report does not contain material that has already been used to any substantial extent for a comparable purpose.

Total word count: 0

**Signed:**

**Date:**

This dissertation is copyright ©2021 Sean J. Parker.

All trademarks used in this dissertation are hereby acknowledged.



# Acknowledgements

# Abstract

Write a summary of the whole thing. Make sure it fits in one page.

# Contents

<b>1</b>	<b>Introduction</b>	<b>vii</b>
<b>2</b>	<b>Background and Related Work</b>	<b>1</b>
2.1	Introduction to Deep Learning Models . . . . .	1
2.1.1	Current approaches to optimising deep learning models	2
2.2	Reinforcement Learning . . . . .	3
2.2.1	Model-Free and Model-Based RL . . . . .	5
2.2.2	World Models . . . . .	6
2.3	Graph Neural Networks . . . . .	10
<b>3</b>	<b>Title</b>	<b>11</b>
<b>4</b>	<b>Title</b>	<b>13</b>
<b>5</b>	<b>Title</b>	<b>15</b>





# List of Figures

2.1	Single perceptron as a dataflow (computation) graph . . . . .	2
2.2	Model-based Reinforcement Learning End-To-End System . .	7
2.3	Structure of an unrolled LSTM . . . . .	9



# List of Tables



# Chapter 1

## Introduction

# Chapter 2

## Background and Related Work

### 2.1 Introduction to Deep Learning Models

This section discusses the way in which machine learning models are represented for efficient execution on physical hardware devices. First, we discuss how the mapping of tensor operations to computation graphs is performed followed by an overview of recent approaches that optimise computation graphs to minimise execution time.

Over the past decade, there has been a rapid development of various deep learning architectures that aim to solve a specific task. Common examples include convolutional networks (popularised by AlexNet then ResNets, etc), transformer networks that have seen use in the modelling and generation of language. Recurrent networks that have shown to excel at learning long and short trends in data.

Importantly, the fundamental building blocks of the networks have largely remained unchanged. As the networks become more complex, it becomes untenable to manually optimise the networks to reduce the execution time on hardware. Therefore, there is extensive work in ways to both automatically optimise the models, or, alternatively apply a set of hand-crafted optimisations.

Computation graphs are a way to graphically represent both the individual tensor operations in a model, and the connections (or data-flow) along the edges between nodes in the graph. Figure 2.1 shows how the expression,  $y = \text{ReLU}(\mathbf{w} \cdot \mathbf{x} + b)$ , can be represented graphically in a computation graph.

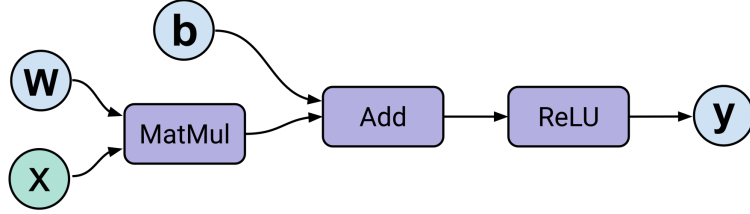


Figure 2.1: The operations shown in purple are the nodes of the computation graph which take an arbitrary number of inputs, performs a computation at the node and produces an output. The blue nodes represent the input nodes for tensors. The directed edges show the flow of tensors through the graph.

Similarly, the whole model can be converted into a stateful dataflow graph in this manner. By using a stateful dataflow (or computation) graph, we can use any optimisation technique for backpropagation of the model loss through the graph [TODO rewrite last sentence]. We consider two key benefits of this representation. First, we can execute the model on any hardware device as the models have a single, uniform representation. Secondly, it allows for pre-execution optimisations based on the host device, for example, we may perform different optimisations for executing on a GPU compared to a TPU.

### 2.1.1 Current approaches to optimising deep learning models

Due to the prevalence and importance of machine learning, especially deep networks, there is a focus on finding ways decrease the inference runtime and by extension, increasing the model throughput. All major frameworks such as TensorFlow [1], PyTorch [2], MXNet [3], and Caffe [4] have some level of support for performing pre-execution optimisations. However, the process of performing such optimisations is often time-consuming and cannot

be completed in real-time. Rather, it is common to use a deep learning optimisation library such as cuDNN [5] or cuBLAS [6] that instead directly optimise individual tensor operations.

Alternatively, TVM and TensorRT can be used to optimise deep learning models and offer greater performance gains compared to the more commonly used frameworks such as TensorFlow and PyTorch. They also use greedy rule-based optimisation approaches. TODO - either expand or remove

Rather than using a rule-based optimisation approach, it is possible to use more sophisticated algorithms to optimise deep learning models at the expense of computation time. Jia et al. used a cost-based backtracking search to iteratively search through the state space of possible graphs that are provably equivalent [7]. As opposed to using rule-based optimisation that applied hand-crafted optimisations, TASO generates the candidate subgraphs automatically and formally proves the transformations are equivalent using an automated theorem prover. Furthermore, Jia et al. showed that by jointly optimising both the data layout of the subgraph transformation, and the transformation itself, TASO achieves a speedup compared to performing the operations sequentially.

A key benefit of using a cost-based approach is that the search can take into account far more complex interactions between the transformed kernels. For example, if we apply a series of transformations  $T_1, \dots, T_i$ , the runtime may increase, and, due to the first set of transformations, we can now apply  $T_{i+1}, \dots, T_{i+j}$ , after all transformations have been applied, it is possible that we see a net decrease in runtime. By increasing the search space of transformations in this way, it is possible to increase the runtime of deep learning models by up to 3x [7, 8].

## 2.2 Reinforcement Learning

Reinforcement learning (RL) is a sub-field in machine learning, broadly, it aims to compute a control policy such that an agent can maximise its cu-



mulative reward from the environment. It has powerful applications in environments where a model that describes the semantics of the system are not available and the agent must itself discover the optimal strategy via a reward signal.

- TODO Can also mention POMDPs

Formally, RL is a class of learning problem that can be framed as a Markov decision processes (MDP) when the MDP that describes the system is not known [9]; they are represented as a 5-tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}_a, \mathcal{R}_a, \rho_0 \rangle$  where:

- $\mathcal{S}$ , is a finite set of valid states
- $\mathcal{A}$ , is a finite set of valid actions
- $\mathcal{P}_a$ , is the transition probability function that an action  $a$  in state  $s_t$  leads to a state  $s'_{t+1}$
- $\mathcal{R}_a$ , is the reward function, it returns the reward from the environment after taking an action  $a$  between state  $s_t$  and  $s'_{t+1}$
- $\rho_0$ , is the starting state distribution

We aim to compute a policy, denoted by  $\pi$ , that when given a state  $s \in \mathcal{S}$ , returns an action  $a \in \mathcal{A}$  with the optimisation objective being to find a control policy  $\pi^*$  that maximises the *expected reward* from the environment defined by 2.1. Importantly, we can control the ‘far-sightedness’ of the policy by tuning the discount factor  $\gamma \in [0, 1)$ . As  $\gamma$  tends to 1, the policy will consider the rewards further in the future but with a lower weight as the distant expected reward may be an imperfect prediction.

$$\pi^* = \arg \max_{\pi} \mathbb{E} \left[ \sum_{t=0}^{\infty} \gamma^t \mathcal{R}_t \right] \quad (2.1)$$

Classic RL problems are formulated as MDPs in which we have a finite state space, however, such methods quickly become inefficient with large state spaces that we consider with applications such as Atari and Go. Therefore, we take advantage of modern deep learning function approximators, such as

neural networks, that makes learning the solutions far more efficient in practise. We have seen many successful applications in a wide range of fields, for example, robotic control tasks [10], datacenter power management, device placement, and, playing both perfect and imperfect information games to a super-human level. Reinforcement learning excels when applied to environments in which actions may have long-term, inter-connected dependencies that are difficult to learn or model with traditional machine learning techniques.

In the following sections we discuss the two key paradigms that exist in reinforcement learning and the current research in both areas and the application to systems tasks.

### 2.2.1 Model-Free and Model-Based RL

Model-free and model-based are the two main approaches to reinforcement learning, however, with recent work such as [11, 12, 13], the distinction between the two is becoming somewhat nebulous; it is possible to use a hybrid approach that aims to improve the sample efficiency of the agent by training model-free agents directly in the imagined environment.

The major branching point that distinguishes between model-free and model-based approaches is in what the agent learns during training. A model-free agent, in general, could learn a governing policy, action-value function, or, environment model. On the other hand, model-based agents commonly either learn an explicit representation of the parameterised policy  $\pi_\theta$  using planning, such as AlphaZero [14] or ExIt [15]. Alternatively, we can use data augmentation methods to learn a representation of the underlying environment behaviour, and either only use fictitious model, or augment real experiences to train an agent in the domain [12, 16, 17].

Understandably, a relevant question is why one would prefer a model-free over model-based approach and what are the benefits of the respective methods. The primary benefit of model-based RL is that it has far greater sample

efficiency, meaning, the agent requires in total, less interactions with the real environment than the model-free counterparts. If we can either provide, or learn, a model of the environment it allows the agent to plan ahead, choosing from a range of possible trajectories its actions to maximise its reward. The agent that acts in this “*imagined*” or “*hallucinogenic*” environment can be a simple MLP [18] to a model-free agent trained using modern algorithms such as PPO, A2C or Q-learning. Further, training an agent in the world model is comparatively cheap, especially in the case of complex systems environments where a single episode can be on the order of hundreds of milliseconds.

Unfortunately, learning a model of the environment is not trivial. The most challenging problem that must be overcome is that if the model is imperfect, the agent may learn to exploit the model’s deficiencies, thus making it effectively useless in the real environment.

Model-based approaches have been successfully applied in various domains such as board games, video games, systems optimisation and robotics. Despite the apparent advantages of model-based RL with regards to reduced computation time, model-free reinforcement learning is by far the most popular approach and massive amounts of compute, typically by distributed training on clusters of GPUs/TPUs, is required to overcome the sample inefficiency of model-free algorithms.

### 2.2.2 World Models

World models, first introduced by Ha and Schmidhuber, motivated and described an approach to model-based reinforcement learning in which we learn a model of the real environment using function approximators and train an agent using only predictions from the world model. Figure 2.2 shows the design to utilise a world model as substitute for the real environment. In practice, a world model can be broken down into three main components. A visual model,  $V$ , that encodes the input into a latent vector  $z$ , a memory model,  $M$  that integrates the historical state to produce a representation that can be used as planning for future actions and rewards. Finally, a controller,

$C$  that uses both  $V$  and  $M$  to predict an action from the action set,  $a \in \mathcal{A}$ .

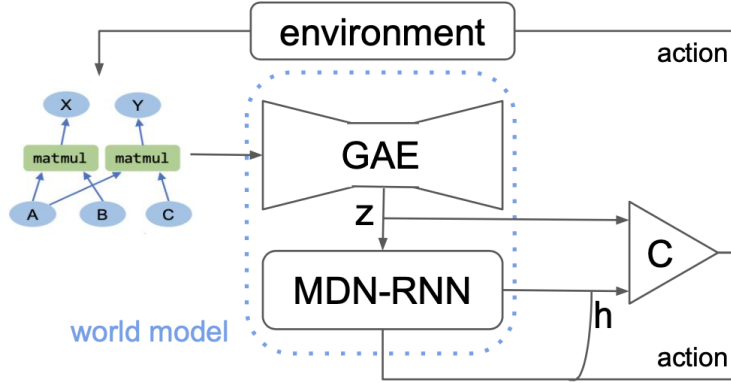


Figure 2.2: Diagrammatic representation of a world-model made up from an encoder ( $V$ ) that transforms the input into latent space, a ‘memory’ module,  $M$ , that learns the behaviour of the environment from the latent vector  $z$  and a controller,  $C$ , that is trained using the latent vector of the encoder and the output features from the memory to choose an action which is either applied to the real or imagined environment. Figure adapted from [18].

Typically, a world model is trained using rollouts of the real environment that have been sampled using a random agent acting in the environment. The aim is to learn to accurately predict, given a state  $s_t$ , the next state  $s_{t+1}$  and the associated reward  $r_{t+1}$ . After training, the controller,  $C$ , can either learn using only observations from the world model, so called “training in a dream”. Alternatively, the world model can be used to augment the observations from the real environment samples or used only for planning. To construct the world model, if the environment is simple and deterministic, it is possible to use a deep neural network to act as the world model, however, for environments that are only partially observable, a more complex model is required such as Recurrent Neural Networks (RNNs) or Long-short term memory (LSTMs). The following two sub-sections describe the fundamental concepts required to construct a world model.

## Mixture Density Networks

Mixture Density Networks (MDNs) are a class of neural networks first described by Christopher Bishop [19] that were designed to deal with problems where there is an inherent uncertainty in the predictions. Given an input to the network, we wish to output a range of possible outputs conditioned on the input where we can assess the probability of each outcome. MDNs are commonly parameterised by a neural network that is trained using supervised learning and outputs the parameters for multiple mixture of Gaussians.

## Recurrent Neural Networks

Recurrent Neural Networks are class of neural networks that allows for previous outputs to be re-used as inputs to sequential nodes while maintaining and updating their own hidden state. Primarily, RNNs are commonly used in the field of speech recognition and natural language processing as they can process inputs of an arbitrary length with a constant model size. In practise however, RNNs suffer from being unable to utilise long chains of information due to the vanishing/exploding gradient problem; the gradient can change exponentially changing in proportion to the number of layers in the network [20].

Motivated by the desire to overcome the limitations of RNNs, Hochreiter et al. [21] developed long-short term memory by describing Constant Error Carousel (CECs). The idea was further improved by Gers et al. [22] with the modern LSTM that is made up of four gates, each with a specific purpose that influences the behaviour of each cell and in combination, the properties of the network as a whole. Figure 2.3 shows the internal structure of an LSTM module.

An LSTM can be described used four “gates”, where a gate influences a property of the behaviour of the LSTM cell. The *forget* gate dictates if the information stored in the cell should be erased by observing the inputs  $[h_{t-1}, x_t]$  it outputs a value in the range  $[0, 1]$ , using the sigmoid function



Figure 2.3: LSTM

$\sigma$ , where 1 means to completely forget the state. Secondly, the *input* gate calculates the new information to be stored in the cell state, generating a vector of candidates  $\tilde{C}_t$ . The *update* gate is to determine how much of the past state sequence should be considered using the outputs from the *forget* gate, the prior state  $C_{t-1}$  and the input gate  $\tilde{C}_t$ . Finally, the *output* gate determines the LSTM cell output based on the current, filtered state of the cell.

$$\begin{aligned}
 f_t &= \sigma(W_f \cdot [h_{t-1}, x_t] + b_f) & (1) \text{ Forget gate} \\
 i_t &= \sigma(W_i \cdot [h_{t-1}, x_t] + b_i) & (2) \text{ Input gate} \\
 \tilde{C}_t &= \tanh(W_C) \cdot [h_{t-1}, x_t] + b_C & (3) \text{ Candidate value} \\
 C_t &= f_t C_{t-1} + i_t \tilde{C}_t & (4) \text{ Update previous cell state} \\
 o_t &= \sigma(W_o [h_{t-1}, x_t] + b_o) & (5) \text{ Output gate} \\
 h_t &= o_t \cdot \tanh(C_t) & (6) \text{ Hidden state}
 \end{aligned}$$

There are a number of popular variants of LSTM cells such as peephole LSTMs, GRUs, Grid LSTMs and ConvLSTM. There are many areas which have been revolutionised by the usage of LSTM cells in the network architecture. As we will describe in [TODO: ref chapter], LSTM cells are a key to allow world models to learn to simulate the behaviour of the environment state-action transitions.

## 2.3 Graph Neural Networks

---

**Algorithm 1:** Steps of computation in a full GN block

---

```

for  $k \in \{1 \dots N^e\}$  do
     $\mathbf{e}'_k \leftarrow \phi^e(\mathbf{e}_k, \mathbf{v}_{r_k}, \mathbf{v}_{s_k}, \mathbf{u})$ 
end
for  $i \in \{1 \dots N^n\}$  do
    let  $E'_i = \{(\mathbf{e}'_k, r_k, s_k)\}_{r_k=i, k=1:N^e}$ 
     $\bar{\mathbf{e}}'_i \leftarrow \rho^{e \rightarrow v}(E'_i)$ 
     $\bar{\mathbf{v}}'_i \leftarrow \phi^v(\bar{\mathbf{e}}'_i, \mathbf{v}_i, \mathbf{u})$ 
end
let  $V' = \mathbf{v}'_{i=1:N^v}$ 
let  $E' = (\mathbf{e}'_k, r_k, s_k)_{k=1:n^e}$ 
 $\bar{\mathbf{e}}' \leftarrow \rho^{e \rightarrow u}(E')$ 
 $\bar{\mathbf{v}}' \leftarrow \rho^{v \rightarrow u}(V')$ 
 $\bar{\mathbf{u}}' \leftarrow \phi^u(\bar{\mathbf{e}}', \bar{\mathbf{v}}', \mathbf{u})$ 
return  $(E', V', \bar{\mathbf{u}}')$ 

```

---

- Form of neural network that has seen lots of active research recently
- Can be viewed in similar way to CNNs
- Graph network, Edge block -  $\bar{\cdot}$  Node block -  $\cdot$  global block
- We use it to process the graphs and produce a latent tensor (from the global block) which we can learn using

## Chapter 3

### Title





# Chapter 4

## Title



# Chapter 5

## Title



# Bibliography

- [1] Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen, Craig Citro, Greg S. Corrado, Andy Davis, Jeffrey Dean, Matthieu Devin, Sanjay Ghemawat, Ian Goodfellow, Andrew Harp, Geoffrey Irving, Michael Isard, Yangqing Jia, Rafal Jozefowicz, Lukasz Kaiser, Manjunath Kudlur, Josh Levenberg, Dandelion Mané, Rajat Monga, Sherry Moore, Derek Murray, Chris Olah, Mike Schuster, Jonathon Shlens, Benoit Steiner, Ilya Sutskever, Kunal Talwar, Paul Tucker, Vincent Vanhoucke, Vijay Vasudevan, Fernanda Viégas, Oriol Vinyals, Pete Warden, Martin Wattenberg, Martin Wicke, Yuan Yu, and Xiaoqiang Zheng. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.
- [2] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. Pytorch: An imperative style, high-performance deep learning library. In *Advances in Neural Information Processing Systems 32*, pages 8024–8035. Curran Associates, Inc., 2019.
- [3] Tianqi Chen, Mu Li, Yutian Li, Min Lin, Naiyan Wang, Minjie Wang, Tianjun Xiao, Bing Xu, Chiyuan Zhang, and Zheng Zhang. Mxnet: A flexible and efficient machine learning library for heterogeneous distributed systems, 2015.
- [4] Yangqing Jia, Evan Shelhamer, Jeff Donahue, Sergey Karayev, Jonathan Long, Ross Girshick, Sergio Guadarrama, and Trevor Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22nd ACM International Conference on Multimedia*, MM '14, page 675678, New York, NY, USA, 2014. Association for Computing Machinery.

- [5] Sharan Chetlur, Cliff Woolley, Philippe Vandermersch, Jonathan Cohen, John Tran, Bryan Catanzaro, and Evan Shelhamer. cudnn: Efficient primitives for deep learning, 2014.
- [6] NVIDIA. cuBLAS Library. <https://developer.nvidia.com/cublas>, 2008.
- [7] Zhihao Jia, Oded Padon, James Thomas, Todd Warszawski, Matei Zaharia, and Alex Aiken. Taso: Optimizing deep learning computation with automatic generation of graph substitutions. In *Proceedings of the 27th ACM Symposium on Operating Systems Principles, SOSP '19*, page 4762, New York, NY, USA, 2019. Association for Computing Machinery.
- [8] Zhihao Jia, James Thomas, Tod Warszawski, Mingyu Gao, Matei Zaharia, and Alex Aiken. Optimizing dnn computation with relaxed graph substitutions. *SysML 2019*, 2019.
- [9] Richard Bellman. A markovian decision process. *Journal of Mathematics and Mechanics*, 6(5):679–684, 1957.
- [10] OpenAI, Ilge Akkaya, Marcin Andrychowicz, Maciek Chociej, Mateusz Litwin, Bob McGrew, Arthur Petron, Alex Paino, Matthias Plappert, Glenn Powell, Raphael Ribas, Jonas Schneider, Nikolas Tezak, Jerry Tworek, Peter Welinder, Lilian Weng, Qiming Yuan, Wojciech Zaremba, and Lei Zhang. Solving rubik’s cube with a robot hand, 2019.
- [11] Harrison Brown, Kai Fricke, and Eiko Yoneki. World-models for bitrate streaming. *Applied Sciences*, 10(19), 2020.
- [12] Lukasz Kaiser, Mohammad Babaeizadeh, Piotr Milos, Blazej Osinski, Roy H Campbell, Konrad Czechowski, Dumitru Erhan, Chelsea Finn, Piotr Kozakowski, Sergey Levine, Afroz Mohiuddin, Ryan Sepassi, George Tucker, and Henryk Michalewski. Model-based reinforcement learning for atari, 2020.
- [13] Jan Robine, Tobias Uelwer, and Stefan Harmeling. Smaller world models for reinforcement learning, 2021.
- [14] David Silver, Thomas Hubert, Julian Schrittwieser, Ioannis Antonoglou, Matthew Lai, Arthur Guez, Marc Lanctot, Laurent Sifre, Dhharshan Kumaran, Thore Graepel, Timothy Lillicrap, Karen Simonyan, and Demis Hassabis. Mastering chess and shogi by self-play with a general reinforcement learning algorithm, 2017.
- [15] Thomas Anthony, Zheng Tian, and David Barber. Thinking fast and slow with deep learning and tree search, 2017.

- [16] Vladimir Feinberg, Alvin Wan, Ion Stoica, Michael I. Jordan, Joseph E. Gonzalez, and Sergey Levine. Model-based value estimation for efficient model-free reinforcement learning, 2018.
- [17] C. Daniel Freeman, Luke Metz, and David Ha. Learning to predict without looking ahead: World models without forward prediction, 2019.
- [18] David Ha and Jürgen Schmidhuber. Recurrent World Models Facilitate Policy Evolution, 2018. <https://worldmodels.github.io>.
- [19] Christopher M Bishop. Mixture density networks, 1994.
- [20] S. Hochreiter, Y. Bengio, P. Frasconi, and J. Schmidhuber. Gradient flow in recurrent nets: the difficulty of learning long-term dependencies. In S. C. Kremer and J. F. Kolen, editors, *A Field Guide to Dynamical Recurrent Neural Networks*. IEEE Press, 2001.
- [21] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [22] Felix A Gers, Jürgen Schmidhuber, and Fred Cummins. Learning to forget: Continual prediction with lstm. 1999.