

Todo list

This paragraph looks redundant. Should I trim this into two sentences instead?	11
This formula looks redundant. Should I cut it?	12
my notation is messed up but I will fix them later for sure.	19
(20 - 25 pages)	21
reference training	27
This terminology is aligned with MuZero's pseudo-code and MCTX's real code. Should I mention this?	30
make this image smaller	33
add figure here	36
add figure here	37
can't think of a better title.	39
Do I need to show what the environments look like? Those will be static images so the readers will be clueless still...	42
(3 pages)	47
(1 page)	49

UNIVERSITY OF ALBERTA

DOCTORAL THESIS

High Performance Game-playing System that Plans with a Learned Model

Author:
Zeyi Wang

Supervisor:
Martin Müller

*A thesis submitted in fulfillment of the requirements
for the degree of Master of Science (Thes)*

Department of Computing Science

August 13, 2022

UNIVERSITY OF ALBERTA

Abstract

Department of Computing Science

Master of Science (Thesis)

High Performance Game-playing System that Plans with a Learned Model

by Zeyi Wang

The Thesis Abstract is written here (and usually kept to just this page). The page is kept centered vertically so can expand into the blank space above the title too...

For/Dedicated to/To my...

Declaration of Authorship

I, Zeyi Wang, declare that this thesis titled, “High Performance Game-playing System that Plans with a Learned Model” and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

Date:

Contents

List of Tables

List of Figures

List of Symbols

a	distance	m
P	power	W (J s ⁻¹)
ω	angular frequency	rad

Symbol	Description	Reference
s	state	??
a	action	??
r	reward	??
t	timestep	??
T	terminal timestep	??
γ	discount	??
o	partially observable environment frame	
G	N-step return	??
G^N	N-step return	??
V	value function	??
Q	state-action value function	??
δ	TD-error or value diff	??
\mathcal{A}^e	environment action space	??
\mathcal{A}^a	agent action space	??
\mathcal{S}	state space	
\mathcal{O}	observation space	
\mathcal{T}_t	step sample	??
\mathcal{T}	trajectory sample	
\mathcal{L}	loss function	
\mathbf{x}	hidden state	
h	representation function	
g	dynamics function	
f	prediction function	
ϱ	projection function	
v^i	value prediction	
B	batch size	
H	height	??
W	width	??
C_e	environment channels	??
C_h	history channels	??
C_x	hidden space channels	
K	number of unrolled steps	??
L	history length	??
N	bootstrap length for N-step return	??
A	dimension of action	
Z	support of the scalar transformation	??, ??

List of Symbols

Symbol	Description	Reference
s	state	??
a	action	??
r	reward	??
t	timestep	??
T	terminal timestep	??
γ	discount	??
o	partially observable environment frame	
G	N-step return	??
G^N	N-step return	??
V	value function	??
Q	state-action value function	??
δ	TD-error or value diff	??
\mathcal{A}^e	environment action space	??
\mathcal{A}^a	agent action space	??
\mathcal{S}	state space	
\mathcal{O}	observation space	
\mathcal{T}_t	step sample	??
\mathcal{T}	trajectory sample	
\mathcal{L}	loss function	
\mathbf{x}	hidden state	
h	representation function	
g	dynamics function	
f	prediction function	
ϱ	projection function	
v^i	value prediction	
B	batch size	
H	height	??
W	width	??
C_e	environment channels	??
C_h	history channels	??
C_x	hidden space channels	
K	number of unrolled steps	??
L	history length	??
N	bootstrap length for N-step return	??
A	dimension of action	
Z	support of the scalar transformation	??, ??

Chapter 1

Introduction

Deep Learning (DL) is a branch of **Artificial Intelligence (AI)** that emphasizes the use of neural networks to fit any arbitrary function represented by a dataset. The training of a neural network is done by computing a loss function from a batch of data, back-propagate gradients with respect to the loss, and updates weights and biases based on the gradients. Deep learning techniques have been widely adopted in many domains, including computer vision, natural language processing, and robotics.

Reinforcement Learning (RL) is a branch of AI that emphasizes on solving decision making problems through trials and errors with delayed rewards. RL had most success in the domain of **game playing**, in which the algorithm is represented as an **agent** and interacts with the game environments, such as boardgames and Atari games. An extension to game playing is **general game playing (GGP)**, whose goal is to design a single agent that can play many different games without having much prior knowledge of the games.

Deep Reinforcement Learning (DRL) is a rising branch that combines DL and RL to solve decision making problems. In a DRL system, the RL techniques lay out the structure of the algorithm such as the use the **agent-environment interface**, a value function, a reward signal, e.t.c., while the DL techniques are used approximate specific functions and learn representations.

Planning refers to any computational process that analyzes generated actions and their consequences in an environment. In the RL terms, planning specifically means the use of a model to improve a policy. In boardgames where perfect models are accessible, planning with these models yield great performance. The most significant achievement of planning with a perfect model is AlphaGo beating human champion in Computer Go. However, how to plan in games where no perfect models available remains a challenging problem to researchers.

A **distributed system** is a computer system that uses multiple processes with various purposes to complete tasks. DRL systems for solving large problems are both data and compute intensive. Utilizing concurrency to increase efficiency and throughput for these DRL systems are sometimes necessary. Building a distributed system to achieve such concurrency is a common practice in the industry but requires significant engineering effort.

1.1 Motivation

MasteringAtariGo'Schrittwieser.Antonoglou.ea'2020 developed MuZero, an algorithm that plans with a learned model (??). This algorithm achieved the state-of-the-art in playing both Atari games and boardgames. However, the source code of the algorithm is not publicly available, and the pseudo-code provided with the paper isn't sufficient to reproduce the full algorithm. Moreover, MuZero

requires much more computations than other model-free RL algorithms, and an inefficient implementation will drastically slows down experimentation. We need an efficient and publicly available implementation of an algorithm that plans with a learned model. This helps researchers understand how does the algorithm plan with its learned model, and facilitates future research.

1.2 Contribution

In this thesis we present the project **MooZi**, a system that play games by planning with a learned model. This project includes:

- A collection of environments bridges that connect the system to various common RL environments.
- Neural networks that learns representation and can be used for planning.
- A MCTS based planner that uses the learned model to perform planning.
- A distributed training system that efficiently trains the algorithm.
- A thesis with empirical studies and analysis.

Chapter 2

Literature Review

2.1 Planning and Search

Many AI problems can be reduced to a search problem [ArtificialIntelligenceGames'Yannakakis.Togelius'2018]. Such search problems can be solved by determining the best plan, path, model, function, and so on, based on some metrics of interest. Therefore, search has played a vital role in AI research since its dawn. The terms **planning** and **search** are widely used across different domains. Here we adopt the definition by ReinforcementLearningIntroduction'Sutton.Barto'2018 [ReinforcementLearningIntroduction'Sutton.Barto'2018].

Planning refers to any process by which the agent updates the action selection policy $\pi(a | s)$ or the value function $V_\pi(s)$. We will focus on the case of improving the policy in our discussion. We view the planning process as an operator \mathcal{I} that takes the policy as input and outputs an improved policy $\mathcal{I}\pi$.

Planning methods can be categorized based on the target state s they aim to improve. If the method improves the policy for arbitrary states, we call it **background planning**. That is, for any timestep t and a set of states $S' \subset \mathcal{S}$:

$$\pi(a | s) \leftarrow \mathcal{I}\pi(a | s), \forall s \in S'$$

Typical background planning methods include **dynamic programming** and **Dyna-Q** [ReinforcementLearningIntroduction'Sutton.Barto'2018]. In the case of dynamic programming, a full sweep of the state space is performed and all states are updated. In the case of Dyna, a subset of the state space is selected for update.

An other type of planning focuses on improving the policy of the current state s_t . We call this **decision-time planning**. That is, for any timestep t :

$$\pi(a | s) \leftarrow \mathcal{I}\pi(a | s), s = s_t$$

Algorithms such as AlphaGo use both types of planning when they use self-play for training. For decision-time planning, a tree search is performed at the root node and updates the policy of the current state. For background planning, a neural network uses past experience to train and updates policy for all states.

An early example of the use of search as a planning method is the **A*** algorithm. In 1968, FormalBasisHeuristic'Hart.Nilsson.ea'1968 designed the A* algorithm for finding a shortest path from a start vertex to a target vertex [FormalBasisHeuristic'Hart.Nilsson.ea'1968]. Although A* works quite well for many problems, especially in early game AI, it falls short in cases where the assumptions of A* do not hold. For example, A* requires a heuristic, and an optimal solution under stochastic environments. It is computationally infeasible on large problems. To address this problem, RealtimeHeuristicSearch'Korf'1990

framed the problem of **Real-Time Heuristic Search**, where the agent has to make a decision in each timestep with bounded computation, and developed the **Real-Time-A*** algorithm as a modified version of A* with bounded computation per step [RealtimeHeuristicSearch`Korf`1990]. Tree-based search algorithms such as **MiniMax** and **Alpha-Beta Pruning** were developed to play and solve two-player games [AnalysisAlphaBetaPriming`Knuth.Moore`1975]. Monte Carlo techniques are designed to handle complex environments.

2.2 Monte Carlo Methods

In 1873, Joseph Jagger observed the bias in roulette wheels at the Monte Carlo Casino. He studied the bias by recording the results of roulette wheels and won over 2 million francs over several days by betting on the most favorably biased wheel [MonteCarloCasino`2022]. Therefore, **Monte Carlo (MC)** methods gained their name as a class of algorithms based on random samplings.

MC methods are used in many domains but in this thesis we will primarily focus on its usage in search. In a game where terminal states are usually unreachable by the limited search depth, evaluation has to be performed on the leaf nodes that represent intermediate game states. One way of obtaining an evaluation on a state is by applying a heuristic function. Heuristic functions used this way are usually hand-crafted by human based on expert knowledge, and hence are prone to human error. The other way of evaluating the state is to perform a rollout from that state to a terminal state by selecting actions randomly. This evaluation process is called **random rollout** or **Monte Carlo rollout**.

2.3 Monte Carlo Tree Search (MCTS)

BanditBasedMonteCarlo`Kocsis.Szepesvari`2006 developed the **Upper Confidence Bounds applied to Trees (UCT)** method as an extension of the **Upper Confidence Bound (UCB)** algorithm employed in multi-armed bandit problems [BanditBasedMonteCarlo`Kocsis.Szepesvari`2006]. Rémi Coulom developed the general idea of **Monte Carlo Tree Search** that combines Monte Carlo rollouts with tree search [EfficientSelectivityBackup`Coulom`2007] for his Go program CrazyStone. Shortly afterwards, **ModificationUCTPatterns`Gelly.Wang.ea`2006** implemented another Go program MoGo, that uses the UCT selection formula [ModificationUCTPatterns`Gelly.Wang.ea`2006]. MCTS was generalized by **MonteCarloTreeSearch`Chaslot.Bakkes.ea`2008** as a framework for game AI [MonteCarloTreeSearch`Chaslot.Bakkes.ea`2008]. This framework requires less domain knowledge than classic approaches to game AI while giving better results. The core idea of this framework is to gradually build the search tree by iteratively applying four steps: **selection**, **expansion**, **evaluation**, and **backpropagation**. The search tree built in this way emphasizes more promising moves and game states based on collected statistics in rollouts. More promising states are visited more often, have more children, have deeper subtrees, and rollout results are aggregated to yield more accurate values. Here we detail the four steps in the MCTS framework by **MonteCarloTreeSearch`Chaslot.Bakkes.ea`2008** (see Figure ??).

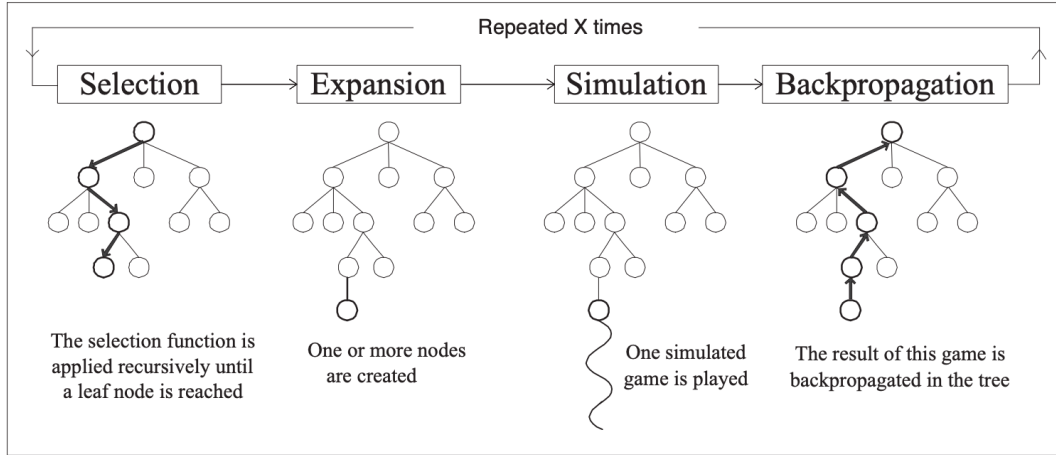


Figure 2.1: The Monte Carlo Tree Search Framework. Source: [MonteCarloTreeSearch'Chaslot.Bakkes.ea'2008]

2.3.1 Selection

The selection process starts at the root node and repeats until a leaf node in the current tree is reached. At each level of the tree, a child node is selected based on a selection formula such as UCT or PUCT. A selection formula usually has two parts: the exploitation part based on the evaluation function E , and the exploration bonus function B . For actions $(s, a), a \in \mathcal{A}$ of a parent state s , the selection $I(s)$ is defined as

$$I(s) = \operatorname{argmax}_{a \in \mathcal{A}} [E(s, a) + B(s, a)]$$

The evaluation function E can be based on the value of the child, the accumulated reward of the child, or the prior selection probability based on the policy $\pi(a | s)$. The exploration bonus function B is usually based on the visit count of the child and the parent. The more visits a child has, the smaller the exploration bonus becomes. For example, the UCT algorithm uses

$$E(s, a) = \frac{V(s)}{N(s, a)}$$

$$B(s, a) = \sqrt{\frac{2 * \log(\sum_{b \in \mathcal{A}} N(s, b))}{N(s, a)}}$$

where $V(s)$ is the value of the node, and $N(s, a)$ is the visit count of the edge. This **ModificationUCTPatterns'Gelly.Wang.ea'2006** used this selection rule in their implementation of MoGo, the first computer Go program that uses UCT [ModificationUCTPatterns'Gelly.Wang.ea'2006].

MultiarmedBanditsEpisode'Rosin'2011 developed the PUCB and the PUCT algorithm that utilize a predictor $P(s, a)$ that estimates the prior probability of the action a being selected from state s and later being used in AlphaGo (??, [MultiarmedBanditsEpisode'Rosin'2011]).

2.3.2 Expansion

The selected leaf node is expanded by adding one or more children. Each child represents a successor game state reached by playing the associated legal move.

2.3.3 Evaluation

The expanded node is evaluated, either by playing a game with a rollout policy, or by using an evaluation function, or by using a blend of both approaches. Many MCTS algorithms use a randomized policy as the rollout policy and the game result as the evaluation. Early work on evaluation functions focused on hand-crafted or machine learned heuristic functions based on expert knowledge. Recently, evaluation functions use deep neural networks specifically trained for the problems (??).

2.3.4 Backpropagation

After the expanded nodes are evaluated, the nodes on the path from the expanded nodes back to the root are updated. The statistics updated usually include visit count, estimated value and accumulated reward of the nodes.

2.3.5 MCTS Iteration and Move Selection

The four MCTS steps are repeated until the budget runs out. The budget is usually a limited number of simulations or a period of time. After the search, the agent acts by selecting the action associated with the most promising child of the root node. This could be the most visited child, the child with the greatest value, or the child with the greatest lower confidence bound value [FreshMaxLcb`RoyJonathan`2019, AcceleratingSelfPlayLearning`Wu`2020].

2.4 AlphaGo

In **MasteringGameGo`Silver.Schrittwieser.ea`2017**, **MasteringGameGo`Silver.Schrittwieser.ea`2017** developed **AlphaGo**, the first Go program that beat a human Go champion on even terms [MasteringGameGo`Silver.Schrittwieser.ea`2017]. AlphaGo was trained with a machine learning pipeline with multiple stages. For the first stage of training, a supervised learning policy (or SL policy) is trained to predict expert moves using a neural network. This SL policy p is parametrized by weights σ , denoted p_σ . The input of the policy network is a representation of the board state, denoted s . The network outputs a probability distribution over all legal moves a through the last softmax layer. During the training of the network, randomly sampled expert moves are used as training targets. The weights σ are then updated through gradient ascent to maximize the probability of matching human expert move:

$$\Delta\sigma \propto \frac{\partial \log p_\sigma(a | s)}{\partial \sigma}$$

For the second stage of training, the supervised policy p_σ is used as the starting point for training with reinforcement learning. This reinforcement learning trained policy (or RL policy) is parametrized by weights ρ and is initialized $p_\rho = p_\sigma$. Training data is generated in form of self-play games using p_ρ as the rollout policy. For each game,

the game outcome $z_t = \pm r(s_T)$, where s_T is the terminal state, $z_T = +1$ for winning, $z_T = -1$ for losing from the perspective of the current player. Weights ρ are updated using gradient ascent to maximize the expected outcome using the update formula:

$$\Delta\rho \propto \frac{\partial \log p_\rho(a_t | s_t)}{\partial \rho} z_t$$

Finally, a value function is trained to evaluate board positions. This value function is modeled with a neural network with weights θ , denoted V_θ . Given a state s , $V_\theta(s)$ predicts the outcome of the game if both players act according to the policy p_ρ . This neural network is trained with stochastic gradient descent to minimize the mean squared error (MSE) between the predicted value $V_\theta(s)$ and the outcome z .

$$\Delta\theta \propto \frac{\partial V_\theta(s)}{\partial \theta} (z - V_\theta(s))$$

AlphaGo combines the policy network p_ρ and the value network V_θ with MCTS for acting. AlphaGo uses a MCTS variant called PUCT similar to that described in ???. In the search tree, each edge (s, a) stores an action value $Q(s, a)$, a visit count $N(s, a)$, and a prior probability $P(s, a)$. At each time step, the search starts at the root node and simulates until the budget runs out. In the select phase of each simulation, an action is selected for each traversed node using the same base formula (???). In AlphaGo, the exploitation score of the selection formula is the estimated average value of the next state after taking the actions, namely $Q(s, a)$. In AlphaGo's PUCT formula, The exploration bonus of edge (s, a) is based on the prior probability P and decays as its visit count N grows. As before, the action taken at time t maximizes the sum of the exploitation score and the exploration bonus

$$\begin{aligned} I(s) &= \operatorname{argmax}_{a \in \mathcal{A}} [E(s, a) + B(s, a)] \\ E(s, a) &= Q(s, a) \\ B(s, a) &\propto \frac{P(s, a)}{1 + N(s, a)} \end{aligned}$$

AlphaGo evaluates a leaf node state s_L by blending both the value network estimation $V_\theta(s_L)$ and the game result z_L obtained by the rollout policy p_π . The mixing parameter $\lambda \in [0, 1]$ is used to balance these two types of evaluations into the final evaluation $V(s_L)$

$$V(s_L) = (1 - \lambda)V_\theta(s_L) + \lambda z_L$$

2.5 AlphaGo Zero

AlphaGo Zero is a successor of AlphaGo that beat AlphaGo by 100-0 in 100 games [MasteringGameGo`Silver.Schrittwieser.ea`2017]. In contrast, AlphaGo Zero learns to play Go from *tabula rasa*. This means it learns solely by reinforcement learning from self-play, starting from random play, without supervision from human expert data.

Central to AlphaGo Zero is a deep neural network f_θ with parameters θ . Given a state s as an input, the network outputs both move probabilities \mathbf{p} and value

estimation v

$$(\mathbf{p}, v) = f_\theta(s)$$

To generate self-play games s_1, \dots, s_T , MCTS is performed at each state s using the latest neural network f_θ . To select a move for a parent node p in the search tree, a variant of the PUCT algorithm is used:

$$\begin{aligned} I(s) &= \operatorname{argmax}_{a \in \mathcal{A}} (E(s, a) + B(s, a)) \\ E(s, a) &= Q(s, a) \\ B(s, a) &\propto P(s, a) \frac{\sqrt{\sum_{b \in \mathcal{A}} N(s, b)}}{1 + N(s, a)} \end{aligned}$$

Self-play games are processed into training targets to update the network parameters θ through gradient descent on the loss function l

$$\mathcal{L}(\theta) = (z - v)^2 - \boldsymbol{\pi}^\top \log \mathbf{p} + c \|\theta\|^2$$

Here $(z - v)^2$ is the mean squared error of the prediction value, $-\boldsymbol{\pi}^\top \log \mathbf{p}$ is the cross-entropy loss of the move probabilities, and $c \|\theta\|^2$ is a L_2 weight regularization. Many other components of this system are similar to those in AlphaGo.

2.6 AlphaZero

AlphaZero reduces game specific knowledge of AlphaGo Zero even further so that the same algorithm can be also applied to Shogi and chess [MasteringChessShogi·Silver.Hubert.ea'2017]. One generalization is that in AlphaZero the game result is no longer either winning or losing ($z \in \{-1, +1\}$), but can also be a draw ($z \in \{-1, 0, +1\}$).

2.7 MuZero

In MasteringAtariGo·Schrittwieser.Antonoglou.ea'2020, MasteringAtariGo·Schrittwieser.Antonoglou.ea'2020 developed **MuZero**, an even more algorithm that learns to play Atari, Go, chess and Shogi at superhuman level. Compared to the AlphaGo and AlphaZero, MuZero has no access to a perfect model of the game. MuZero plans with a neural network that learns the game dynamics through experience. Therefore, MuZero can be applied to games where either the perfect model is not known or is infeasible to compute with.

MuZero defines three main functions. The **representation function** h encodes a history of observations o_1, o_2, \dots, o_t and actions a_1, a_2, \dots, a_{t-1} into a hidden state \mathbf{x}_t^0 . This hidden state is learned, and is the main conceptual change from AlphaZero. The **dynamics function** g implements action execution in the representation. Given a hidden state \mathbf{x}^k and action a^k , produces an immediate reward r^k and the next hidden state \mathbf{x}^{k+1} . The **prediction function** f corresponds to the one network in AlphaZero. Given a hidden state \mathbf{x}^k , it produces a probability distribution p^k of actions and a value v^k associated to that hidden state. Three

functions f, g, h are approximated jointly in a neural network with weights θ

$$\mathbf{x}_t^0 = h_\theta(o_1, \dots, o_t, a_1, \dots, a_{t-1}) \quad (2.1)$$

$$(\mathbf{x}^{k+1}, \hat{r}^{k+1}) = g_\theta(\mathbf{x}^k, a^k) \quad (2.2)$$

$$(v^k, \mathbf{p}^k) = f_\theta(\mathbf{x}^k) \quad (2.3)$$

The superscripts of \mathbf{x}, a, v denote the depth of such values in the search tree, and depth 0 is at the search tree's root. Equivalently, the superscripts also mean the number of recurrent inferences (through the dynamics function g) the algorithm performs to obtain that value.

MuZero plans with a search method based on the MCTS framework (discussed in ??). Due to the lack of access to a perfect model, MuZero's MCTS differs from a standard one in numerous ways. The nodes are no longer perfect representations of the board states. Instead, each node is associated with a hidden state \mathbf{x} as a learned representation of the board state. The transition is no longer made by the perfect model but by the dynamics function g . Moreover, since the dynamics function also predicts a reward, edges created through inferencing with the dynamics function also contribute to the Q value estimation.

To act in the environment, MuZero plans following the MCTS framework described in section ??. At each timestep t , \mathbf{x}_t^0 is created using (?). A variant of PUCT is used to select an action during the search

$$\begin{aligned} I(s) &= \underset{a \in \mathcal{A}}{\operatorname{argmax}} (E(s, a) + B(s, a)) \\ E(s, a) &= Q(s, a) \\ B(s, a) &\propto P(s, a) \frac{\sqrt{\sum_{b \in \mathcal{A}} N(s, b)}}{1 + N(s, a)} \left[c_1 + \log \left(\frac{\sum_{b \in \mathcal{A}} N(s, b) + c_2 + 1}{c_2} \right) \right] \end{aligned}$$

Where c_1 and c_2 are two constants that adjust the exploration bonus. The selected edge (\mathbf{x}^k, a^k) at depth k is expanded using (?) and evaluated using (?). At the end of the simulation, the statistics of the nodes along the search path are updated. We denote the updated prior action probabilities \mathbf{p}^* , and the updated value estimation v^* . Notice since the transitions of the nodes are approximated by the neural network, the search is performed over hypothetical trajectories without using a perfect model. Finally, the action a^0 of the most visited edge (\mathbf{x}^0, a^0) of the root node is selected as the action to take in the environment.

Experience generated is stored in a replay buffer and processed into training targets. The three functions of the model are trained jointly using the loss function

$$\mathcal{L}_t(\theta) = \underbrace{\sum_{k=0}^K \mathcal{L}^p(p_{t+k}^*, p_t^k)}_{\textcircled{1}} + \underbrace{\sum_{k=0}^K \mathcal{L}^v(z_{t+k}, v_t^*)}_{\textcircled{2}} + \underbrace{\sum_{k=1}^K \mathcal{L}^r(r_{t+k}, \hat{r}^k)}_{\textcircled{3}} + \underbrace{c \|\theta\|^2}_{\textcircled{4}} \quad (2.4)$$

where K is the rollout depth, $\textcircled{1}$ is the loss of the predicted prior move probabilities and move probabilities improved by the search, $\textcircled{2}$ is the loss of the predicted value and experienced N -step return, $\textcircled{3}$ is the loss of the predicted reward and the experienced reward, and finally $\textcircled{4}$ is the L_2 regularization.

2.7.1 MuZero Reanalyze

MasteringAtariGo [Schrittwieser.Antonoglou.ea'2020] also developed **MuZero Reanalyze**, a sample efficient variant of MuZero [MasteringAtariGo.Schrittwieser.Antonoglou.ea'2020]. This method generates training targets in addition to those generated through game play through re-executing search on old games using the latest parameters. **MuZero Unplugged** and **Efficient Zero** also use similar mechanisms to generate new data by updating search statistics of old data [OnlineOfflineReinforcement.Schrittwieser.Hubert.ea'2021b, MasteringAtariGames.Ye.Liu.ea'2021]. In Efficient Zero, experiments are ran with a reanalyze ratio of 0.99, which means only 1% of the training data are generated through interacting with the environment, and the other 99% are generated by re-running search on old trajectories. In our project, we also implement a reanalyze worker to perform this task (see section ??, ??).

2.8 Atari Games Playing

2.8.1 Atari Learning Environment

The **Atari 2600** gaming console was developed by *Atari, Inc.* and was released in 1977. Over 30 million copies of the console sold over its 15 years on the market [Atari2600''2022]. The most popular game, **PacMan**, was sold over 8 million copies and was the all-time best-selling video game back then. **Stella** is a multi-platform Atari 2600 emulator released under the GNU General Public License (GPL) [StellaMultiPlatformAtari'']. Stella was ported to popular operating systems such as Linux, MacOS, and Windows, providing Atari 2600 experiences to users without physical copies of the equipment. In **ArcadeLearningEnvironment** [Bellemare.Naddaf.ea'2013a, ArcadeLearningEnvironment.Bellemare.Naddaf.ea'2013a] introduced the **Arcade Learning Environment (ALE)** and the library has been publicly available since [ArcadeLearningEnvironment.Bellemare.Naddaf.ea'2013a]. ALE provides interfaces of over a hundred of Atari game environments using Stella as the backend. Each ALE environment has specifications on its visual representation, action space, and reward signals. ALE environments are suitable for controlled machine learning research, because data are well-represented and evaluation metrics are clearly defined. Moreover, ALE environments are diverse in their characteristics: while some environments require more mechanical mastery of the agent, others require more long-term planning. This makes solving multiple ALE environments using the same algorithm a good general game playing problem (??).

2.8.2 Deep Q-Networks

PlayingAtariDeep [Mnih.Kavukcuoglu.ea'2013] pioneered the study of using deep neural networks to learn in ALE environments [PlayingAtariDeep.Mnih.Kavukcuoglu.ea'2013]. They developed the algorithm **Deep Q-Networks (DQN)** that learned to play seven of the Atari games and reached human-level performance. The DQN agent has a neural network that approximates the Q function, parametrized by weights θ , denoted Q_θ . Experiences are generated through interacting with the environment by taking the action that

maximizes the immediate Q value

$$\pi(a_t \mid (o_{t-L+1}, \dots, o_t)) = \underset{a}{\operatorname{argmax}} Q_\theta(o_{t-L+1}, \dots, o_t, a)$$

where L is the length of history, and o_t is the “frame”, a partial observation of the game state at timestep t (also see ??). Generated experience is stored in an experience replay buffer implemented as a FIFO queue. For each training step, a batch of uniformly sampled experience is drawn from the experience replay, and the loss is computed using

$$\mathcal{L}(\theta) \propto \mathbb{E}_\pi \left[r + \gamma \max_{a'} Q_{\theta'}(s', a') - Q_\theta(s, a) \right] .$$

The network parameters θ' are updated less frequently than θ .

2.8.3 Double Q Learning

DoubleQLearning Hasselt'2010 analyzed the overestimation problem of Q values in Q -learning and developed **double Q learning**, where a double Q update replaces the traditional Q update [**DoubleQLearning** Hasselt'2010]. Double Q learning reduces the overestimation problem by introducing an additional Q estimator's and updating two estimator using each other

$$\begin{aligned} Q^A(s, a) &\leftarrow Q^A(s, a) + \alpha \left(r + \gamma Q^B \left(s', \underset{a'}{\operatorname{argmax}} Q^A(s', a') \right) - Q^A(s, a) \right) \\ Q^B(s, a) &\leftarrow Q^B(s, a) + \alpha \left(r + \gamma Q^A \left(s', \underset{a'}{\operatorname{argmax}} Q^B(s', a') \right) - Q^B(s, a) \right) \end{aligned}$$

where Q^A and Q^B are two different Q estimators updated alternately. **DeepReinforcementLearning** Hasselt.Guez.ea'2016 applied the double Q learning in DQN [**DeepReinforcementLearning** Hasselt.Guez.ea'2016]. Similar to the double Q update above, a double Q update for neural networks is formulated as

$$\begin{aligned} \mathcal{L}(\theta^A) &\propto \mathbb{E}_\pi \left[r + \gamma Q_{\theta^B} \left(s', \underset{a'}{\operatorname{argmax}} Q_{\theta^A}(s', a') \right) - Q_{\theta^A}(s, a) \right] \\ \mathcal{L}(\theta^B) &\propto \mathbb{E}_\pi \left[r + \gamma Q_{\theta^A} \left(s', \underset{a'}{\operatorname{argmax}} Q_{\theta^B}(s', a') \right) - Q_{\theta^B}(s, a) \right] . \end{aligned}$$

Here Q_{θ^A} and Q_{θ^B} are two sets of parameters of the same neural network architecture.

This paragraph looks redundant. Should I trim this into two sentences instead?

2.8.4 Experience Replay

PrioritizedExperienceReplay Schaul.Quan.ea'2016 studied the role of a experience replay in DQN and developed the **prioritized experience replay** method [**PrioritizedExperienceReplay** Schaul.Quan.ea'2016]. In the original work of DQN, all samples were drawn from the experience replay uniformly. In prioritized experience replay, however, samples are drawn according to a

distribution based on their calculated priority

$$P(i) = \frac{p_i^\alpha}{\sum_k p_k^\alpha}$$

where $P(i)$ is the probability of the i -th sample being drawn, α is a constant, and p_i is the priority of the sample. **PrioritizedExperienceReplay** Schaul.Quan.ea'2016 developed two approaches to compute priorities of samples. In **proportional sampling**, the priority p of sample i is calculated by

$$p_i = |\delta_i| + \epsilon$$

where δ_i is the temporal-difference error of the sample, and ϵ is a small constant to give all samples a non-zero probability to be drawn. In **rank-based sampling**, the same temporal differences are calculated, but the final priority is computed based on the rank of the error,

$$\begin{aligned} \text{score}(i) &= |\delta_i| + \epsilon \\ p_i &= \frac{1}{\text{rank}(\text{score}(i))} \end{aligned}$$

This formula looks redundant. Should I cut it?

DistributedPrioritizedExperience Horgan.Quan.ea'2018 followed up by implementing a distributed version of prioritized experience replay [**DistributedPrioritizedExperience** Horgan.Quan.ea'2018]. **RECURRENTEXPERIENCEREPLAY** Kapturowski.Ostrovski.ea'2019 investigated the challenges of using experience replays for RNN-based agents and developed **Recurrent Replay Distributed DQN** [**RECURRENTEXPERIENCEREPLAY** Kapturowski.Ostrovski.ea'2019].

2.8.5 Network Architectures

DuelingNetworkArchitectures Wang.Schaul.ea' studied an alternative neural network architecture for ALE learning [**DuelingNetworkArchitectures** Wang.Schaul.ea']. **Dueling Q-network** retains the input and output specifications of the Q-network used in DQN and structurally represented the learning of the advantage function $A(s, a)$ defined as

$$A(s, a) = Q(s, a) - V(s)$$

The Q-network has three parts: θ , the shared trunk of the network; Λ , the advantage head; and Y , the value head. The network approximates the value function internally through the shared trunk and the value head, denoted $V_{\theta, Y}$, and the advantage function, denoted $A_{\theta, \Lambda}$. The values computed by the two heads are combined to form the Q-value as follows

$$Q_{\theta, Y, \Lambda}(s, a) = V_{\theta, Y}(s) + \left(A_{\theta, \Lambda}(s, a) - \frac{1}{|\mathcal{A}|} \sum_{a'} A_{\theta, \Lambda}(s, a') \right)$$

Similar to DQN, the dueling Q-network is trained through fitting to empirical data generated by interacting with the environment. Experiments show that this architecture encourages the network to learn to differentiate between the values of

states and the values of state-action pairs, and leads to better performance of the agent.

2.8.6 Scalar Transformation

ObserveLookFurther’Pohlen.Piot.ea’2018 introduced enhancements to achieve more stable training in Atari games [**ObserveLookFurther’Pohlen.Piot.ea’2018**]. We focus on discussing the **transformed Bellman Operator** since both MuZero and MooZi use it. For different Atari games, reward signals can vary drastically both in density and scale. This leads to high variance in training targets during training of the algorithms, causing algorithms to have difficulty converging. In DQN, rewards are clipped the reward signal to a range of $[-1, 1]$ to reduce such variance [**PlayingAtariDeep’Mnih.Kavukcuoglu.ea’2013**]. However, this clipping discards the scale of rewards and consequently changes the set of optimal policies. The transformed Bellman Operator was developed to address this problem. The Q update of the new operator is as follows

$$Q(s, a) \leftarrow Q(s, a) + \alpha \phi \left(r + \gamma \max_{a' \in \mathcal{A}} \phi^{-1} (Q(s', a')) \right)$$

where ϕ is an invertible transformation that contracts. One example of such a transformation is

$$\begin{aligned} \phi(x) &= \text{sign}(x) \left(\sqrt{|x| + 1} - 1 \right) + \varepsilon x \\ \phi^{-1}(x) &= \text{sign}(x) \left(\left(\frac{\sqrt{1 + 4\varepsilon(|x| + 1 + \varepsilon)} - 1}{2\varepsilon} \right)^2 - 1 \right) \end{aligned}$$

Both MuZero and MooZi use this specific ϕ definition for both value transformations and reward transformations (??).

2.8.7 MinAtar

MinAtar, developed by **MinAtarAtariInspiredTestbed’Young.Tian’2019**, is an open-source project that offers RL environments inspired by ALE [**MinAtarAtariInspiredTestbed’Young.Tian’2019**]. MinAtar offers five environments that pose similar challenges to ALE environments: learning representation from raw pixels, and learning behaviors that associate actions and delayed rewards. MinAtar environments are implemented in pure Python, have simpler environment dynamics, and are visually less rich than ALE environments. This makes MinAtar perfect test environment for university research.

2.8.8 Consistency Loss

One interesting characteristic of Atari-like games is that the environment frames are usually temporally consistent. For example, given the position of the player avatar for the last few frames, it is not difficult for a human to guess the position of the avatar in the next frame. To take advantage of this property, one common approach is to enforce temporal consistency in the loss function. **VisualizingMuZeroModels’deVries.Voskuil.ea’2021** visualized the latent space of a learned model of MuZero in a 3D space, in which a hidden state is a point in the

space [VisualizingMuZeroModels'deVries.Voskuil.ea'2021]. As MuZero applies recurrent inferences to a hidden state, the transitions can be traced as a 1-D path in the 3D space. The consistency loss they developed creates a smoother path in the 3D space and improves performance. MasteringAtariGames'Ye.Liu.ea'2021 developed a project-then-predict structure similar to a Siamese network to enforce consistency [MasteringAtariGames'Ye.Liu.ea'2021, SiameseNeuralNetworks'Koch.Zemel.ea'].

2.9 Deep Reinforcement Learning Systems

Deep reinforcement learning systems involve irregular computation patterns and complicated hardware interactions between CPUs and AI accelerators. Designing such systems efficiently is a great challenge. Decisions the designer has to make include but are not limited to (1) Where and how to generate experience? (2) Where and how to store generated experience? (3) Where to store the model and copies of it? (4) Where is the gradient computation carried out? (4) How to orchestrate processes for stable training? Here we briefly review popular deep reinforcement learning system designs that utilize parallelization to achieve faster and more efficient training.

2.9.1 AsynchronousMethodsDeep'Mnih.Badia.ea'2016's Asynchronous Methods Framework

AsynchronousMethodsDeep'Mnih.Badia.ea'2016 developed asynchronous variants for four popular RL algorithms with a parallelization structure uses actor-learner processes [AsynchronousMethodsDeep'Mnih.Badia.ea'2016]. Each actor-learner process holds a local copy of the model, generates experience locally using the model, and accumulates gradients locally. Once in a while, all local gradients are aggregated to update the global model. Delaying and aggregating updates to neural network parameters reduces gradient variance among processes and achieves a more stable learning. Among the asynchronous algorithm variants, **Asynchronous Advantage Actor Critic (A3C)** had the best performance and achieved the state-of-the-art at the time using only half the training time.

2.9.2 The IMPALA Architecture

IMPALAScalableDistributed'Espeholt.Soyer.ea'2018 developed IMPALA, a scalable distributed deep reinforcement learning agent [IMPALAScalableDistributed'Espeholt.Soyer.ea'2018]. IMPALA deploys two types of computation workers: *actor* and *learner*. A actor holds a copy of the neural network parameters and the environment. It performs model inferences locally to interact with its environments and generates experiences. Generated experiences are saved in a local storage and subsequently pushed into the learner's local storage. The learner holds the master copy of the neural network parameters. Once the learner receives enough experiences from the actors, it samples experiences from its local queue and performs batched forward pass and back-propagation steps using its model. Figure ?? shows two variants of this structure.

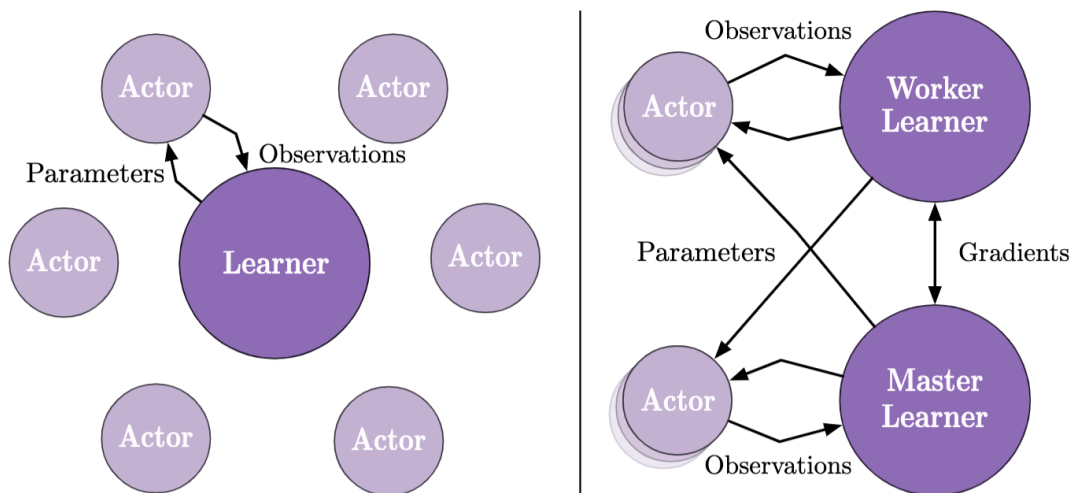


Figure 2.2: IMPALA Architecture, from [IMPALAScalableDistributed’Espeholt.Soyer.ea’2018].
 Left: a single learner computes all gradients; Right: multiple worker learners compute gradients and one master learner collects and aggregates gradients.

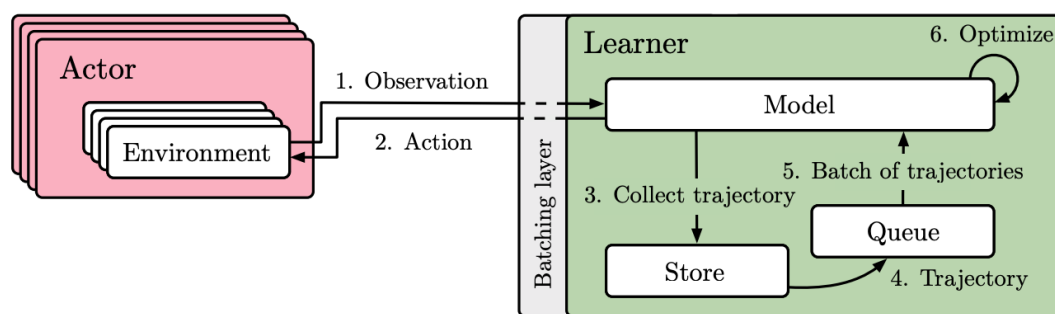


Figure 2.3: The SEED Architecture, from [SEEDRLScalable’Espeholt.Marinier.ea’2020].
 All inferences are computed on the learner and actors act through querying the learner.

2.9.3 The SEED Architecture

SEEDRLScalable’Espeholt.Marinier.ea’2020 developed the **Scalable, Efficient Deep-RL (SEED)** architecture to effectively utilize accelerators using a centralized inference server [SEEDRLScalable’Espeholt.Marinier.ea’2020]. Similar to IMPALA, SEED also uses two main types of workers: actors and learners. However, in SEED, actors do not hold copies of the model. Instead, SEED actors interact with their environments through querying the learner. The learner not only computes gradients and stores trajectories as in IMPALA, but also has a batching layer that batches actor queries and efficiently performs batched inference with the model. Since actors no longer need to pull neural network parameters from the learner, the IO overhead from serializing and messaging parameters is eliminated. Moreover, since the learner batches queries from all actors, the IO overhead from moving inputs and outputs to accelerators (GPUs or TPUs) is also reduced, increasing the overall inferencing throughput. One downside of the SEED architecture is that actors have to wait for a response from the learner to take an action, and thus have a higher latency for taking a step. Figure ?? illustrates a distributed SEED agent.

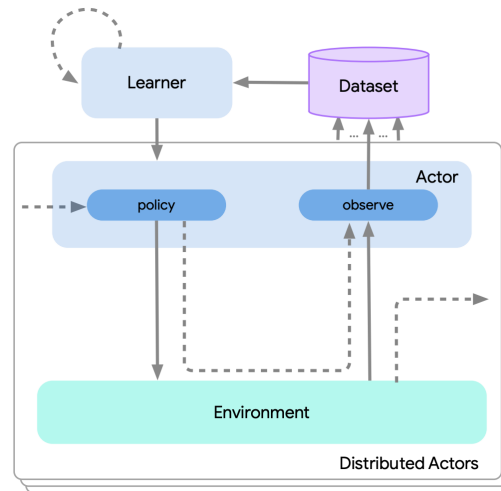


Figure 2.4: Example of a distributed asynchronous agent with Acme, from [AcmeResearchFramework'Hoffman.Shahriari.ea'2020].

2.9.4 The Acme Framework

AcmeResearchFramework'Hoffman.Shahriari.ea'2020 developed the Acme research framework [AcmeResearchFramework'Hoffman.Shahriari.ea'2020]. Acme similar to IMPALA: processes that interact with the environment are actors, and processes that collect experience and update gradients are learners. Additionally, Acme has a *dataset* component, which is synonymous to the replay buffer used in DQN. This component uses **Reverb**, a high-performance library developed by ReverbFrameworkExperience'Cassirer.Barth-Maron.ea'2021 for storing and sampling collected experiences [ReverbFrameworkExperience'Cassirer.Barth-Maron.ea'2021]. Figure ?? illustrates a distributed asynchronous agent in Acme.

2.9.5 Ray and RLlib

RayDistributedFramework'Moritz.Nishihara.ea'2018 designed and implemented Ray, a framework for scalable distributed computing [RayDistributedFramework'Moritz.Nishihara.ea'2018]. Ray enables both task-level and actor-level parallelization through a unified interface. **Ray Core** was designed with AI applications in mind and has powerful primitives for building distributed AI systems. For example, Ray uses shared memory to store inputs and outputs of tasks, allowing zero-copy data sharing among tasks. This is useful for DRL systems in which generated experiences are stored and sampled in a separate process. **RLlibAbstractionsDistributed'Liang.Liaw.ea'2018** developed **RLlib**, an industrial-grade deep reinforcement learning library. RLlib is built on top of Ray Core and provides abstractions for a broad range of DRL systems could make use of. Figure ?? illustrates RLlib's abstraction layers. As of the writing of this thesis, RLlib implemented 24 popular DRL algorithms using its abstractions. One major difference between RLlib agents and other DRL agents is that RLlib deploys a hierarchical control over the worker processes. Our project uses Ray Core to implement its worker processes and deploys a hierarchical control paradigm similar to RLlib (see ??).

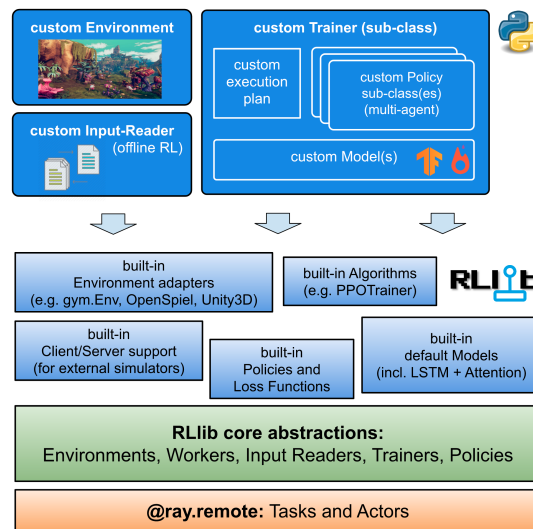


Figure 2.5: RLLib Abstraction Layers, from [RLLibAbstractionsDistributed’Liang.Liaw.ea’2018].

2.9.6 JAX and Podracer Architecture

CompilingMachineLearning’Frostig.Johnson.ea’2019 designed **JAX**, a just-in-time (JIT) compiler that compiles computations expressed in Python code into high-performance accelerators code [CompilingMachineLearning’Frostig.Johnson.ea’2019]. JAX is compatible with **Autograd**, so computation procedures expressed and compiled with JAX can be automatically differentiated. JAX also supports control flow, allowing more sophisticated logic to be expressed while taking advantage of accelerators. Our project uses JAX for both neural networks and search. As a result, we are able to compile the entire policy in rollout workers, including history stacking, planning, and neural networks inferencing, into a single optimized program that can be hardware-accelerated.

PodracerArchitecturesScalable’Hessel.Kroiss.ea’2021 designed two paradigms to efficiently use JAX for DRL systems [PodracerArchitecturesScalable’Hessel.Kroiss.ea’2021]. In the **Anakin** architecture, the environment is implemented with JAX and the entire agent-environment loop is compiled using JAX and computed with accelerators. **Gymnax**, developed by **GymnaxJAXbasedReinforcement’RobertTjarkoLange’2022**, provides environment implementations in native JAX, and is compatible with the Anakin architecture [GymnaxJAXbasedReinforcement’RobertTjarkoLange’2022]. However, pure JAX implemented environments are not always feasible, especially when environments involve external services, such as Stella or Unity in their backend. Alternatively, in the **Sebulba** architecture, environments run on CPUs, but policies could be compiled and computed on accelerators. Generated experiences in both architectures can be used to compute gradients directly on accelerators. Figure ?? illustrates the Sebulba architecture.

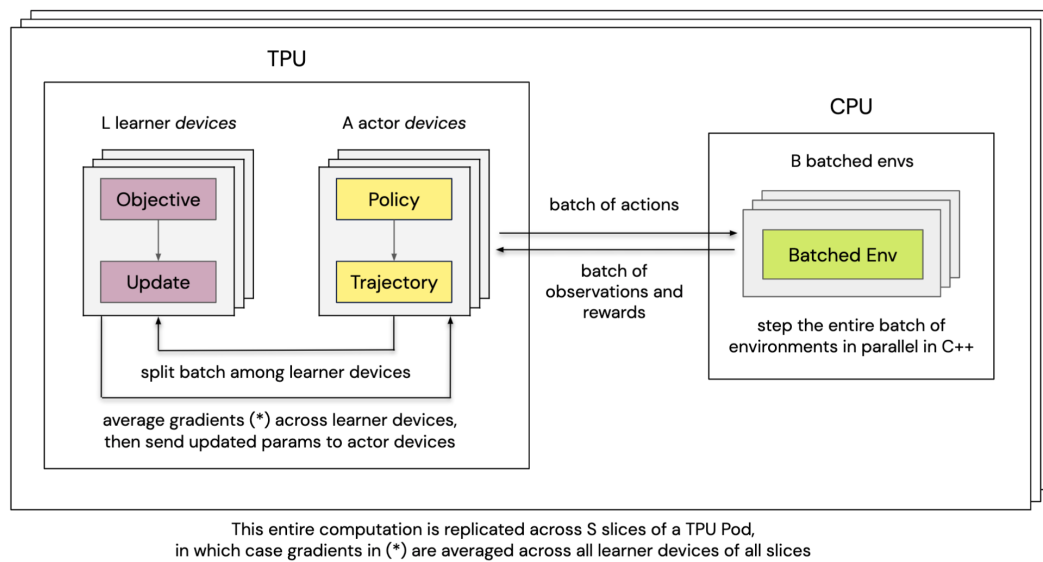


Figure 2.6: Sebulba architecture, from [PodracerArchitecturesScalable'Hessel.Kroiss.ea'2021]. The environments runs on CPUs. Inferences and gradient computations are compiled, optimized and executed on TPUs.

Chapter 3

Problem Definition

3.1 Markov Decision Process and Agent-Environment Interface

A RL problem is usually represented as a **Markov Decision Process (MDP)**. MDP is four-tuple $(\mathcal{S}, \mathcal{A}, R, P)$. \mathcal{S} is a set of states that forms the **state space**. \mathcal{A} is a set of actions that forms the **action space**; $P(s, a, s') = \Pr[s_{t+1} = s' \mid s_t = s, a_t = a]$ is the **transition probability function**. $R(s, a, s')$ is the **reward function**. We use the **agent-environment interface** to solve a problem formulated as an MDP. Figure ?? illustrates the agent-environment interface. The MDP is represented as the **environment**. The decision maker that interacts with the environment is called the **agent**. At each time step t , the agent starts at state $s_t \in \mathcal{S}$, takes an action $a_t \in \mathcal{A}$, transitions to state $s_{t+1} \in \mathcal{S}$ based on the transition probability function $P(s_{t+1} \mid s_t, a_t)$ and receives a reward $R(s_t, a_t, s_{t+1})$. These interactions yield a sequence of actions, states, and rewards $s_0, a_0, r_1, s_1, a_1, r_2, \dots$. We call this sequence a **trajectory**. When a trajectory ends at a terminal state s_T at time $t = T$, this sequence is completed and we called it an **episode**. Figure ?? illustrates the interaction between the agent and the environment.

my notation is messed up but I will fix them later for sure.

3.2 Policies and Value Functions

At each state s , then agent takes an action based on a **policy** $\pi(a \mid s)$. This policy represents the conditional probability of the agent taking an action given a state $\pi(a \mid s) = \Pr[a_t = a \mid s_t = s]$. The agent maximizes the expected discounted sum of

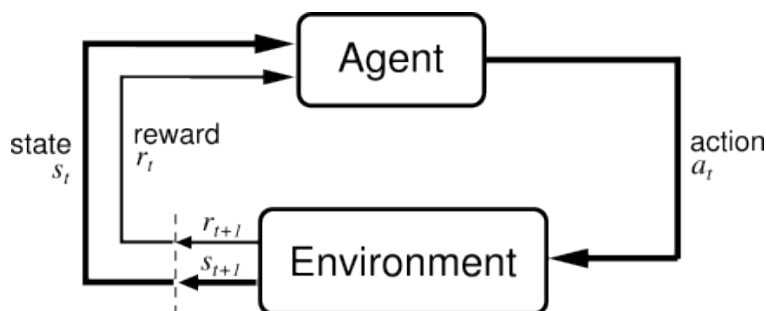


Figure 3.1: The Agent-Environment Interface, from [ReinforcementLearningIntroduction'Sutton.Barto'2018].

rewards from the current state s_t following the policy π

$$\begin{aligned} \text{maximize } & \mathbb{E}_\pi [G_t \mid s_t = s], \quad \forall s \in \mathcal{S} \\ G_t = & \sum_{k=0}^T \gamma^k r_{t+k+1} \quad . \end{aligned}$$

Here γ is the discount factor to favor short-term rewards. G is the discounted sum of rewards, or, equivalently, the discounted return. We represent the maximization target above as the **value function** V

$$V_\pi(s) = \mathbb{E}_\pi [G_t \mid s_t = s] \quad .$$

The value function tells us how good a state is following the policy π . Similarly, we define the **state-action value function**

$$Q_\pi(s, a) = \mathbb{E}_\pi [G_t \mid s_t = s, a_t = a]$$

that tells us how good a state and action pair is.

Chapter 4

Method

(20 - 25 pages)

4.1 Design Philosophy

4.1.1 Use of Pure Functions

One of the most notable differences of the MooZi implementation compared to other implementations is the use of pure functions. In MooZi, we separate the storage of data and the handling of data whenever possible, especially for the parts with heavy computations. We use **JAX** and **Haiku** to implement neural network related modules (??, [HaikuSonnetJAX'Hennigan.Cai.ea'2020, JAXComposableTransformations'JamesBradbury.RoyFrostig.ea'2018]). These libraries separate the **specification** and the **parameters** of a neural network. The **specification** of a neural network is a pure function that is internally represented by a fixed computation graph. The **parameters** of a neural network includes all learned variables that could be used with the specification to perform a forward pass. For example, say we have a simple neural network with a single dense layer that does the following

$$\mathbf{y} = \tanh(\mathbf{Ax} + \mathbf{b})$$

where \mathbf{x} is the input vector of shape $(n, 1)$, \mathbf{y} is the output vector of shape $(m, 1)$, \mathbf{A} is the learned weights of shape (m, n) , and \mathbf{b} is the learned bias of shape $(m, 1)$. We demonstrate how to build this simple network using JAX and Haiku in Algorithm ???. We visualize the computation graph of it in Figure ??.

```

import haiku as hk
import jax
import jax.numpy as jnp

m = 3
n = 2

# specify the computations to be performed
class Model(hk.Module):
    def __call__(self, x):
        A = hk.get_parameter('A', shape=(m, n), init=jnp.zeros)
        b = hk.get_parameter('b', shape=(m, 1), init=jnp.zeros)

        return jax.nn.tanh(A @ x + b)

# haiku transforms the object-oriented model into a functional one
model = hk.without_apply_rng(hk.transform(lambda x: Model()(x)))

# construct a concrete input
x = jnp.ones((n, 1))

# initialize the parameters
params = model.init(jax.random.PRNGKey(0), x)

# perform the forward pass
out = model.apply(params, x)

```

Algorithm 1: A simple dense layer implemented in JAX and Haiku.

The model in the code is the specification of the neural network. The params in the code is the parameters of the neural network. Only params contains concrete data.

Using these pure functions separates the *algorithm* of the agent and the *state* of the agent both conceptually and in implementations. The *algorithm* part of the agent could be abstracted into a computation graph that could be compiled and optimized using a specialized compiler, such as XLA, for hardware acceleration ?? . The *state* part of the agent could be efficiently handled by tools specialized in data manipulations and transferring such as Ray (??). This way, our system efficiently performs inferences on accelerators and transfers data on CPUs.

4.1.2 Training Efficiency

In section ?? we reviewed common DRL systems in which developers gave training efficiency the highest priority in their system designs. We also designed our system so that it's efficient and scalable. Here we describe key features our system has to improve its efficiency. The first one is system parallelization. The computation throughput of a single process is simply not enough for DRL systems. In the published results of MuZero by [MasteringAtariGo'Schrittwieser.Antonoglou.ea'2020], the agent generated over 20 billion environment frames for training. Let's do a quick back-of-envelope-calculation for what this means for a non parallelized system. Consider Gymnax's efficient MinAtar implementation where each environment step takes about 1 millisecond [GymnaxJAXbasedReinforcement'RobertTjarkoLange'2022]. With a single process, it would take more than 200 days just to step the environment. As a result, we have to build a distributed system to increase total throughput through parallelism. The second one is the environment transition speed. In Atari games, especially Atari games in ALE (??), taking one environment step invokes a

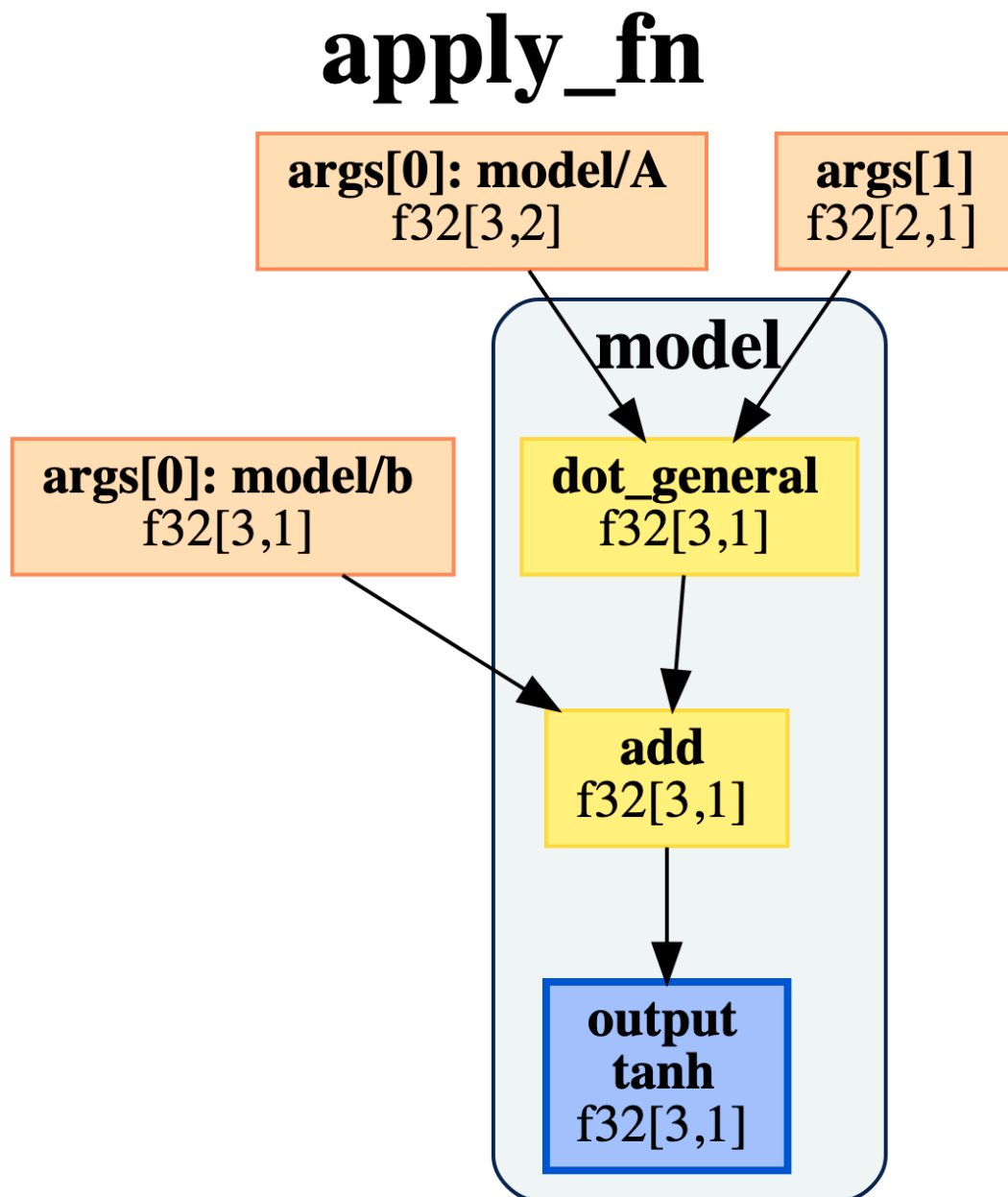


Figure 4.1: Computation graph of the simple dense layer in Algorithm ??. This computation graph shows no concrete data, but the data types, shapes, and operators of the layer (f32 stands for *single-precision float*). To complete a forward pass, we need both concrete neural network parameters (**A**, **b**) and concrete input value (**x**).

full-fledged Atari emulator in the backend and is much more time consuming than neural network inferences. Board games, especially those are implemented in performance focus languages, are much faster. We use MinAtar (??) for simpler variants of Atari games, and OpenSpiel for efficient implementations of board games to reduce the time cost on environment transitions. The third one is neural network inferences used in acting. DRL systems like IMPALA assume that the policy output can be computed by a single forward pass of a neural network. However, MuZero’s policy not only requires dozens inferences per action taken, but also requires a planner that prepares inputs for the inferences and initiates inferences. Our system, utilizing JAX and MCTX, handles planning with multiple inferences per action efficiently.

4.1.3 Understanding is Important

Machine learning algorithms, especially those involve neural networks, have interpretability issues and sometimes could only be used as “black boxes” [ExplainableAIReview’Linardatos.Papastefanopoulos.ea’2021]. We believe that having a system that we can understand is much more useful for future research than having a system that “just works”. Therefore, our project studies the behavior of the system through extensive logging and visualization utilities.

4.2 Architecture Overview

In MooZi, we use the **Ray** library designed by **RayDistributedFramework’Moritz.Nishihara.ea’2018** for orchestrating distributed processes [RayDistributedFramework’Moritz.Nishihara.ea’2018]. We also adopt the terminology used by Ray. In a distributed system with **centralized control**, a single **driver** process is responsible for operating all other processes. Other processes are either **tasks** or **actors**. **Tasks** are stateless functions that take inputs and return outputs. **Actors** are stateful objects that can perform multiple tasks. In the RL literature, **actor** is also a commonly used term for describing the process that holds a copy of the network weights and interacts with an environment [SEEDRLScalable’Espeholt.Marinier.ea’2020, IMPALAScalableDistributed’Espeholt.Soyer.ea’2018]. Even though MooZi does not adopt this concept of a RL actor, we will use the terms **Ray task** and **Ray actor** to avoid confusion. In contrast to distributed systems with **distributed control**, ray tasks and actors are reactive and do not have busy loops. The driver controls when a ray task or actor is activated, what data is used as inputs, and where the outputs goes. The driver orchestrates the data and control flow of the entire system. Ray tasks and actors merely responded to instructions, process input and return output on command. We illustrate MooZi’s architecture in Figure (??).

4.3 Components

4.3.1 Environment Bridges

Environment bridges unify environments which are defined in different libraries into a shared interface. In software engineering terms, environment bridges follow the **bridge design pattern** [BridgePattern’2022]. In our project we implement environment bridges for three types of environments that are commonly used in RL

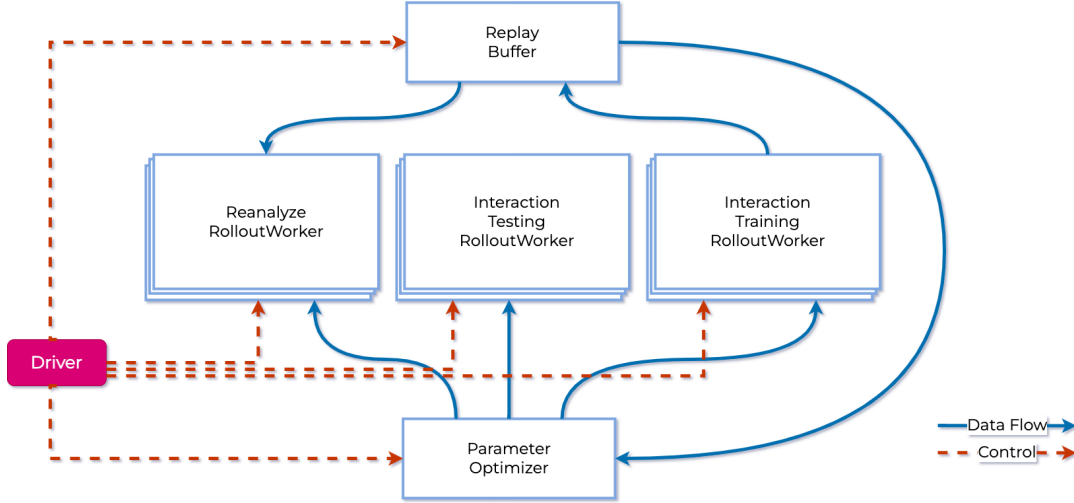


Figure 4.2: The MooZi Architecture. The *driver* is the entrance point of the program and is responsible for setting up configurations, The *parameter server* stores the latest copy of the network weights and performs batched updates to them (??). The *replay buffer* stores generated trajectories and processes these trajectories into training targets (??). A *training worker* is a ray actor responsible for generating experiences by interacting with the environment (??). A *testing worker* is a ray actor responsible for evaluating the system by interacting with the environment (??). A *reanalyze worker* is a ray actor that updates search statistics for history trajectories (??).

research: OpenAI Gym, OpenSpiel, and MinAtar [OpenAIGym`Brockman.Cheung.ea`2016, OpenSpielFrameworkReinforcement`Lanctot.Lockhart.ea`2020, MinAtarAtariInspiredTestbed`Young.Tian`2019]. The bridges wrap these environments into the The DeepMind RL Environment API [DmEnvDeepMind`2022]. In this format, each environment step outputs a four-tuple $(\text{step_type}, r, \gamma, o)$. r, γ, o are reward, discount, and partial observation respectively. The step_type is an enumerated value which indicates the type of timestep. Three possible values of step_type are (1) *first*, indicating the start of an episode, (2) *mid*, indicating an intermediate step, and (3) *last* indicating the last step of an episode. Our bridges wrap these environments again to produce a flat dictionary used by MooZi.

The final environments share the same signature as follows:

- Inputs

b_t^{last} : A boolean indicating the episode end.

a_t : An integer encoding of the action taken.

- Outputs

o_t : An N-dimensional array representing the observation of the current timestep as an image in the shape (H, W, C_e) . H is the height, W is the width, and C_e is the number of channels.

b_t^{first} : A boolean indicating the episode start.

b_t^{last} : A boolean indicating the episode end.

r_t : A float indicating the reward of taking the given action.

$m_t^{A^a}$: A bit mask indicating legal action indices. Valid action indices are 1 and invalid actions indices are 0 (see ??).

All environments are generalized to continuous tasks by passing an additional input b_t^{last} to the environment stepping argument. For an episodic task, the environment is reset internally when b_t^{last} is True. The policy still executes for the last environment step, but the resulting action is discarded. For a continuous task, the environment always step with the latest action and the b_t^{last} input is ignored. Algorithm ?? demonstrates the unified main loop interface.

```
# interact with the environment with a policy indefinitely
def main_loop(env, policy):
    action = 0
    while True:
        result = step(env, action, result.is_last)
        action = policy(result.observation)

def step(env, action, is_last):
    if env.type == "episodic":
        if is_last:
            return env.reset()
        else:
            return env.step(action)
    elif env.type == "continuous":
        return env.step(action)
```

Algorithm 2: Environment Adapter Interface. Both *episodic* environments and *continuous* environments are handled with the same main loop.

We also implement a mock environment using the same interface [MockObject’2021]. A mock environment is initialized with a **trajectory sample** \mathcal{T} , and simulates the environment by outputting step samples one at a time. An agent can interact with this mock environment as if it were a real environment. However, the actions taken by the agent do not affect state transitions since they are predetermined by the given trajectory from initialization. This mock environment is used by the reanalyze workers in section ??.

4.3.2 Vectorized Environment

We also implement a vectorized environment supervisor that stack multiple individual environments to form a single vectorized environment. The resulting vectorized environment takes inputs and produces outputs similar to an individual environment but with an additional batch dimension. For example, an individual environment produces a single frame of shape (H, W, C) , while the vectorized environment produces a batched frame of shape (B, H, W, C) . Previously scalar outputs such as reward are also stacked into vectors of size B . Since environment bridges generalize episodic tasks as continuous tasks, we do not need special handling for the first and the last timesteps in the vectorized environment and its main loop looks exactly like that in Algorithm ?. Using vectorized environments increases the communication bandwidth between the environment and agent and facilitates designing an vectorized agent that processes batched inputs and returns batched actions in one call.

The mock environment described in section ?? is less trivial to vectorize. Each mock environment has to be initialized with a trajectory sample. To vectorize B mock environments, at least B trajectories have to be tracked at the same time. These B trajectories usually have different length and therefore terminate at

different timesteps. Once one of the mocked trajectories reaches its termination, another trajectory has to fill the slot. We create a trajectory buffer to address this problem. When a new trajectory is needed by one of the mocked environments, the buffer replenish it, so the vectorized mocked environment can process batched interactions like a regular vectorized environment until the trajectory buffer runs out of trajectories. An external process has to refill the buffer once in a while. The driver pulls the latest trajectories from the replay buffer and supplies the mock environment's trajectory buffer

reference training

4.3.3 Action Space Augmentation

We augment the action space by adding a dummy action a^{dummy} indexed at 0. This dummy action is used to construct history observations when the horizon extends beyond the current timestep. For example, if the history horizon is 3, we need the last three frames and actions to construct the input observation to the policy. However, if the current timestep is 0, the agent hasn't taken any actions yet. We use zeroed frames with the same shape as history frames, and the augmented dummy action as history actions. Moreover, MooZi's planner (??) does not have access to a perfect model, and it does not know when a node represents a terminal state. Node expansions do not stop at terminal states and the tree search could simulate multiple steps beyond the end. Search performed in these invalid subtrees not only wastes precious search budget, but also back-propagates value and reward estimates that are not learned from generated experience. We address this issue by letting the model learn a policy that always takes the dummy action beyond a terminal state. This learned dummy action acts as a switch that, once taken, treats all nodes in its subtree as absorbing states and edges that have zero values and rewards respectively. This discourages the planner to search in invalid regions and improves search performance for near-end game states. To formally differentiate these two types of action spaces, we denote the original environment action space \mathcal{A}^e and the augmented action space \mathcal{A}^a , and

$$\begin{aligned}\mathcal{A}^a &= \mathcal{A}^e \cup a^{\text{dummy}} \\ a_i &= a^{\text{dummy}} \quad \forall i < 0 && \text{(before the first timestep)} \\ a_i &= a^{\text{dummy}} \quad \forall i \geq T && \text{(after the last timestep)}\end{aligned}$$

Notice that the environment terminates at timestep T so the last effective action taken by the agent is a_{T-1} .

4.3.4 History Stacking

In fully observable environments, the state s_t at timestep t observed by the agent entails sufficient information about the future state distribution. However, for partially observable environments, this does not hold. The optimal policy might not be representable by a policy $\pi(a \mid o_t)$ that only takes into account the most recent partial observation o_t . Most Atari games are such partially observable environments. In DQN, **Playing Atari Deep** Mnih.Kavukcuoglu.ea'2013 alleviated this problem by augmenting the inputs of the policy network from a single frame

observation to a stacked history of four frames so that the policy network had a signature of $\pi(a \mid o_{t-3}, o_{t-2}, o_{t-1}, o_t)$ (??, [PlayingAtariDeep'Mnih.Kavukcuoglu.ea'2013]). AlphaZero and MuZero use not only a stacked history of environment frames, but also a history of past actions. MooZi uses the last L environment frames and taken actions, so the signature of the learned model through the policy head of the prediction function is $\mathbf{p} = f(a \mid o_{t-L+1}, \dots, o_t, a_{t-L}, \dots, a_{t-1})$. The greater L is, the better the stacked observation represents a full state. In a deterministic environment with a fixed starting state, the stacked history represents a full environment state when $L = \infty$. On the other hand, $L = 1$ is sufficient for fully-observable perfect information environments.

The exact process of creating the model input by stacking history frames and actions is as follows:

1. Prepare L saved environment frames of shape (L, H, W, C_e) .
2. Stack the L dimension with the environment channels dimension C_e , resulting in shape $(H, W, L * C_e)$
3. Prepare saved L past actions of shape (L) , encoded as integers.
4. One-hot encode the actions as shape $(L, |\mathcal{A}^a|)$.
5. Normalize the action planes by the number of actions $|\mathcal{A}^a|$, shape remains the same.
6. Stack the L axis with the action axis, now shape $(L * |\mathcal{A}^a|)$.
7. Tile action planes $(L * |\mathcal{A}^a|)$ along the H and W dimensions, now shape $(H, W, L * |\mathcal{A}^a|)$
8. Stack the environment planes and actions planes, now shape $(H, W, L * (C_e + |\mathcal{A}^a|))$
9. The history is now represented as an image with height of H , width of W , and $L * (C_e + |\mathcal{A}^a|)$ channels

To process batched inputs from vectorized environments described in ??, all operations above are performed with an additional batch dimension B , yielding the final output with the shape $(B, H, W, L * (C_e + |\mathcal{A}^a|))$. We denote the channels of the final stacked history $C_h = L * (C_e + |\mathcal{A}^a|)$, where the subscript h means the channel dimension for the representation function h . Figure ?? illustrates this process with an example.

4.3.5 Planner

The planner component \mathcal{P} takes a stacked history as its input (??), performs a search, collects search statistics, and outputs a action and search statistics

$$a_t, v_t^*, \mathbf{p}_t^* = \mathcal{P}(o_{t-L+1}, \dots, o_t, a_{t-L}, \dots, a_{t-1}) .$$

Here v_t^* is the search-updated value estimate of the root, \mathbf{p}_t^* is the search-updated action visits at the root, and a_t is action to take. Training workers (??) use the full planner. Testing workers only use the action output from the planner. The planner used this way is synonymous to a policy π . Reanalyze workers only use the

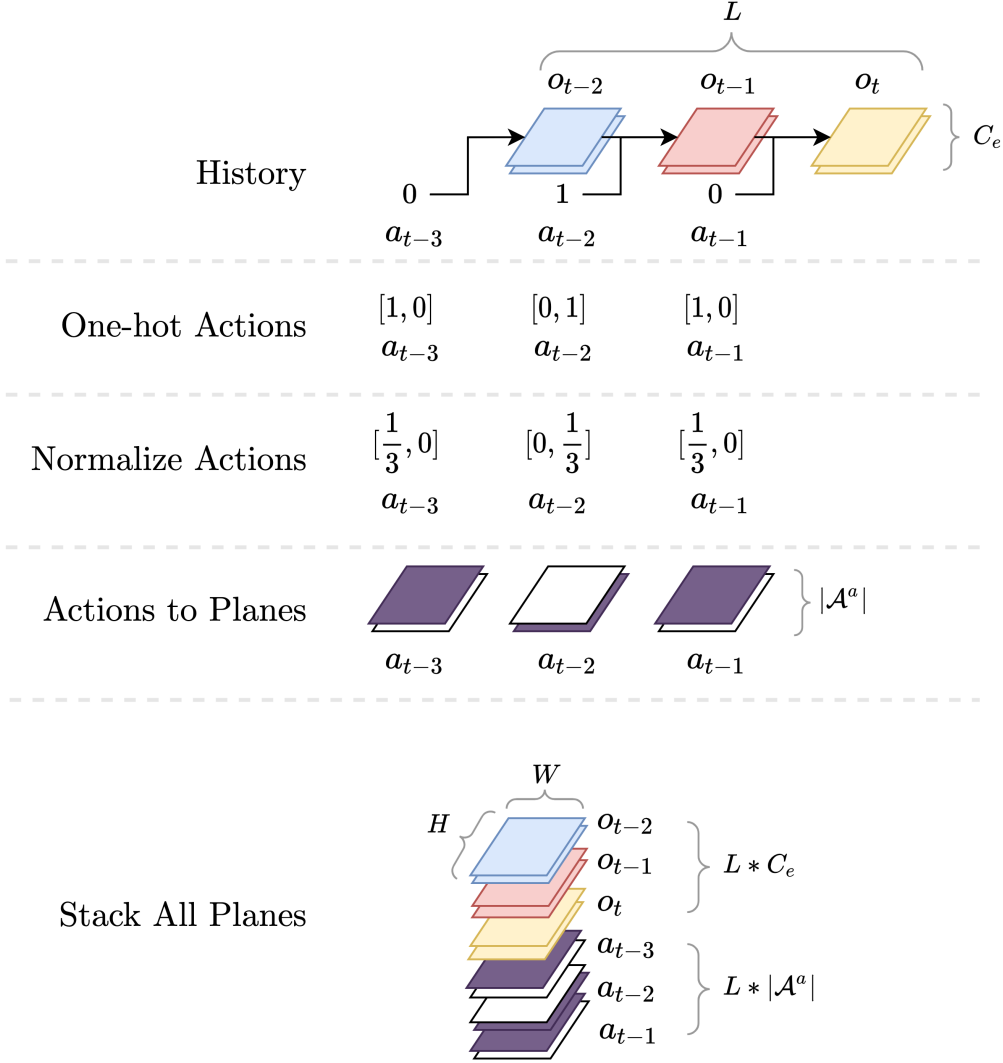


Figure 4.3: An example of history stacking. *History:* Partial observations and actions from the last 3 timesteps ($L = 3$). Actions are integers and observations are images with 2 channels each. *One-hot Actions:* One-hot encodes L history actions into vectors. *Normalize Actions:* Divide the resulting one-hot encoded actions by the size of the action space. *Actions to Planes:* One-hot encodes actions into feature planes that has the same resolution (i.e., same width and height) as the observations, $|\mathcal{A}^a| = 2$. *Stack Planes:* Stack all planes together, creating an image with 12 channels and the same resolution as the observations.

statistics outputs from the planner. Inside of the planner, we use the MuZero variant of MCTS described in ?? ?? with the help from MCTX by POLICYIMPROVEMENTPLANNING`Danilhelka.Guez.ea`2022 [POLICYIMPROVEMENTPLANNING`Danilhelka.Guez.ea`2022].

4.3.6 MooZi Neural Network

MooZi use JAX, and Haiku to build the neural network [HaikuSonnetJAX`Hennigan.Cai.ea`2020, CompilingMachineLearning`Frostig.Johnson.ea`2019, JAXComposableTransformations`JamesBradbury.RoyFrostig.ea`2018]. We consulted other open-source projects that use neural networks to play games [MuZeroGeneral`Duvaud.AureleHainaut`2022, MasteringAtariGames`Ye.Liu.ea`2021, AcceleratingSelfPlayLearning`Wu`2020]. MooZi implements two neural network variants, one is based on multilayer-perceptrons (contribution of Jiuqi Wang) and the other one based residual blocks [DeepResidualLearning`He.Zhang.ea`2016]. These two implementations shared the same interface and can be used interchangeably.

Similar to MuZero described in section ??, the model has a representation function h , a dynamics function g , and a prediction function f . Additionally, MooZi has a **projection function** ϱ for training with the self-consistency loss ?. The learned model is used to construct the root node of a tree search using the representation function h and the prediction function f . We call this process the **initial inference**. The learned model is also used to create edges and child nodes using the dynamics function g and the prediction function f . We call this process the **recurrent inference**. For convenience, the initial inference and the recurrent inference both produce a tuple of $(\mathbf{x}, v, \hat{r}, \mathbf{p})$, where \mathbf{x} is the hidden state, v is the value prediction, \hat{r} is the reward prediction, and \mathbf{p} is the policy prediction. The reward prediction \hat{r} is set to 0 for the initial inference. During training, v and \hat{r} are logits with size $|Z|$. During acting, the logits of v and \hat{r} are first converted to softmax distributions, then converted to scalars using Φ ?.

This terminology is aligned with MuZero's pseudo-code and MCTX's real code. Should I mention this?

Moreover, we applied the invertible transformation ϕ described in section ?? to both the scalar reward targets and scalar value targets to create categorical representations with the same support size. Scalars were first transformed using ϕ , then converted to a linear combination of the nearest two integers in the support. For example, for scalar $\phi(x) = 1.3$, the nearest two integers in the support are 1 and 2, and the linear combination is $\phi(x) = 1 * 0.7 + 2 * 0.3$, which means the target of this scalar is 0.7 for the category 1, and 0.3 for the category 2. We denote Φ for this process of applying ϕ then categorizing the resulting value into a support Z . Using the same example that $\phi(x) = 1.3$, assume the support is $Z = [-2, -1, 0, 1, 2]$, $|Z| = 5$, then $\Phi(x) = [0, 0, 0, 0.7, 0.3]$, and $\Phi(x) \cdot Z = \phi(x) = 1.3$. For training, the value head and the reward head first produced estimations as logits of size $|Z|$. These logits were aligned with the scalar targets to produce categorization loss as described in the ??. For acting, the neural network additionally applied the softmax function to the logits to generated a distribution over the support. The linear combination of the distribution and their corresponding integer values were computed and fed through the inverse of the transformation, namely ϕ^{-1} , to produce scalar values. This means from the

perspective of the planner (??), the scalar estimations made by the model were in same shape and scale as those produced by the environment.

4.3.7 Training Targets Generation

At each timestep t , the environment provides a tuple of data as described in section (??). The agent interacts with the environment by performing a tree search and taking action a_t . The search statistics of the tree search were also saved, including the updated value estimate of the root action \hat{v}_t , and the updated action probability distribution \hat{p}_t . These completes one **step sample** \mathcal{T}_t for timestep t , which is a tuple of $(o_t, a_t, b_t^{\text{first}}, b_t^{\text{last}}, r_t, m_t^{A_a}, \hat{v}_t, \hat{p}_t)$. Once an episode concludes ($b_T^{\text{last}} = 1$), all recorded step samples are gathered and stacked together. This yields a final trajectory sample \mathcal{T} that has a similar shape to a step sample but with an extra batch dimension with the size of T . For example, o_t is stacked from shape (H, W, C_e) to shape (T, H, W, C_e) . The training workers described in ?? generate trajectories this way. The reanalyze workers generate trajectories with the same signature, but through statistics update described in using a vectorized mocked environment (see ?? and ??).

Each trajectory sample with T step samples were processed into T training targets. For each training target at timestep i , we create a training target as follows:

- Observations $o_{i-L+1}, \dots, o_{i+1}$ where H is the history stacking size. The first H observations were used to create policy inputs as described in ??, and the pair of observation o_i, o_{i+1} were used to compute self-consistency loss described in ??.
- Actions $a_{i-L}, \dots, a_{i+K-1}$. Similarly, The first H actions were used for policy input and the pair of actions at (a_{i-1}, a_i) were used for self-consistency loss. The actions a_i, \dots, a_{i+K-1} were used to unroll the model during the training for K steps.
- Rewards r_{i+1}, \dots, r_{i+K} as targets of the reward head of the dynamics function.
- Action probabilities $\mathbf{p}_i^*, \dots, \mathbf{p}_{i+K}^*$ from the statistics of $K + 1$ search trees.
- Root values v_i^*, \dots, v_{i+K}^* , similarly, from the statistics of $K + 1$ search trees.
- N-step return G_i^N, \dots, G_{i+K}^N . Each N-step return was computed based on the formula

$$G_t^N = \sum_{i=0}^{N-1} \gamma^i r_{t+i+1} + \gamma^N v_{t+N}^*$$

- Importance sampling ratio $\rho = 1$. Placeholder value for future override based on replay buffer sampling weights (see ??).

Training targets were computed with minimum information necessary to be used in the loss function (??) so that the precomputed training targets take up the least memory.

4.3.8 Loss Computation

Our loss function is similar to that of ??, but with additional self-consistency loss, terminal action loss, and value loss coefficient

$$\begin{aligned}
 \mathcal{L}_t(\theta) = & \underbrace{\left[\mathcal{L}^p(\mathbf{p}_t^*, \mathbf{p}_t^0) + \frac{1}{K} \sum_{k=1}^K \mathcal{L}^p(\mathbf{p}_{t+k}^*, \mathbf{p}_t^k) \right]}_{\textcircled{1}} \\
 & + c^v \underbrace{\left(\mathcal{L}^v(G_t^N, v_t^*) + \frac{1}{K} \sum_{k=1}^K \mathcal{L}^v(G_{t+k}^N, v_{t+k}^*) \right)}_{\textcircled{2}} \\
 & + \underbrace{\sum_{k=1}^K \mathcal{L}^r(\hat{r}_t^k, r_{t+k})}_{\textcircled{3}} + \underbrace{c^s \mathcal{L}^s(\mathbf{x}_t^1, \mathbf{x}_{t+1}^0)}_{\textcircled{4}} \\
 & + \underbrace{c^{L_2} \|\theta\|^2}_{\textcircled{5}} \Big] \cdot \rho
 \end{aligned}$$

To compute terms used in the loss function, we use the history observations o_{t-L+1}, \dots, o_t and history actions a_{t-L}, \dots, a_{t-1} to reconstruct the stacked frames as the input of the initial inference (??). We apply the initial inference to obtain $\mathbf{p}_t^0, v_t^0, \mathbf{x}_t^0$. We apply K consecutive recurrent inferences using actions a_t, \dots, a_{t+K-1} to obtain $\mathbf{p}_t^1, \dots, \mathbf{p}_t^K, v_t^1, \dots, v_t^K, \mathbf{x}_t^1, \dots, \mathbf{x}_t^K$. The policy loss $\textcircled{1}$ is the standard categorization loss using cross-entropy

$$\mathcal{L}^p(\mathbf{p}, \mathbf{q}) = - \sum_{p \in \mathbf{p}, q \in \mathbf{q}} p \log q$$

The policy targets \mathbf{p}_{t+i}^* ($i = 0, 1, \dots, K$) are action visits at the root of $K + 1$ searches performed in the game (??). To compute the value loss $\textcircled{2}$ and the reward loss $\textcircled{3}$, we apply the scalar transformation Φ (??) that converts scalar values to categorizations, and use the same cross-entropy categorization loss

$$\mathcal{L}^v(p, q) = \mathcal{L}^r(p, q) = - \sum_{p \in \Phi(p), q \in \Phi(q)} p \log q$$

MooZi also trains with a self-consistency loss similar to that described by **MasteringAtariGames`Ye.Liu.ea`2021** and

VisualizingMuZeroModels`deVries.Voskuil.ea`2021

[**MasteringAtariGames`Ye.Liu.ea`2021**,

VisualizingMuZeroModels`deVries.Voskuil.ea`2021].

To compute the self-consistency loss $\textcircled{4}$, we reconstruct the initial inference for the next timestep $o_{t-L+2}, \dots, o_{t+1}, a_{t-L+1}, \dots, a_t$, and compute the cosine distance between the projected one-step hidden state $\varrho(\mathbf{x}_t^1)$ of timestep t and the initial hidden state \mathbf{x}_{t+1}^0

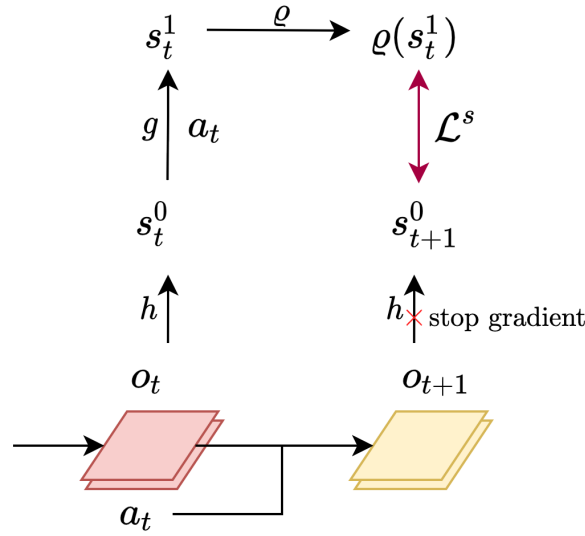


Figure 4.4: Self-consistency Loss Computation. The hidden state \mathbf{x}_t^1 after projection should be similar to the hidden state \mathbf{x}_{t+1}^0 . We assume the next timestep has more information, so we stop gradient from \mathbf{x}_{t+1}^0 to push the representation of the previous timestep towards the next timestep.

of the next timestep $t + 1$. Formally,

$$\text{cosine distance } (\mathbf{a}, \mathbf{b}) = 1 - \frac{\mathbf{a} \cdot \mathbf{b}}{\|\mathbf{a}\| \|\mathbf{b}\|}$$

$$\textcircled{4} = \mathcal{L}^s(\mathbf{x}_t^1, \mathbf{x}_{t+1}^0) = 1 - \frac{\varrho(\mathbf{x}_t^1) \cdot \mathbf{x}_{t+1}^0}{\|\varrho(\mathbf{x}_t^1)\| \|\mathbf{x}_{t+1}^0\|}$$

Figure ?? illustrates the intuition behind this loss.

make this image smaller

$\textcircled{5}$ is a standard L_2 regularization loss to prevent network from overfitting, and coefficient c^{L_2} is used to control the strength of this regularization. The overall loss of a training target is scaled by its importance sampling ratio. We also use the gradient scaling described by [MasteringAtariGo·Schrittwieser.Antonoglou.ea'2020](#) that halves the gradient at the beginning of each dynamics function call [[MasteringAtariGo·Schrittwieser.Antonoglou.ea'2020](#)].

4.3.9 Updating the Parameters

We use a standard Adam optimizer developed by [AdamMethodStochastic·Kingma.Ba'2017](#) [[AdamMethodStochastic·Kingma.Ba'2017](#)]. We also clip the gradient as described by [DifficultyTrainingRecurrent·Pascanu.Mikolov.ea'](#) [[DifficultyTrainingRecurrent·Pascanu.Mikolov.ea'](#)]. The dynamics function g in our learned model is essentially an RNN, so we expect this gradient clipping trick to have a similar effect in our model. **Optax**, developed by [OptaxComposableGradient·MatteoHessel.DavidBudden.ea'2020](#), is a library for

gradient manipulations implemented in JAX [OptaxComposableGradient'MatteoHessel.DavidBudden.ea'2020]. We use Optax's implementation for both the Adam optimizer and the gradient clipper. Moreover, we also use a target network that was used in DQN to stabilize training [PlayingAtariDeep'Mnih.Kavukcuoglu.ea'2013].

4.3.10 Reanalyze

In ??, we reviewed **MuZero Reanalyze**. In our project, we also implement a type of worker process that re-runs search on old trajectories with the latest neural network parameters. Given a trajectory sample \mathcal{T} , for each timestep t in the trajectory, the reanalyze process is as follows

- Use observations (o_{t-T+1}, \dots, o_t) and actions $(a_{t-T}, \dots, a_{t-1})$ to reconstruct the planner input.
- Feed the planner \mathcal{P} with the reconstructed input, obtaining the update action \tilde{a}_t , the updated policy target at the root $\tilde{\mathbf{p}}_t^*$, and the updated value target at the root \tilde{v}_t^* .
- Discard the updated action \tilde{a}_t since the action that got executed in the environment has to be the old action a_t to keep the trajectory consistent.
- Replace the old policy target \mathbf{p}_t^* with the updated policy target $\tilde{\mathbf{p}}_t^*$.
- Replace the old value target v_t^* with the updated policy target \tilde{v}_t^* .

Once the entire trajectory \mathcal{T} is processed, we obtain an updated trajectory $\tilde{\mathcal{T}}$ in which only the value targets and policy targets are replaced.

4.3.11 Training Worker

The main goal of **training workers** is to generate trajectories by interacting with environments for training purposes. For each worker, a vectorized environment is created as described in ??, a history stacker is created as described in ??, and a planner was created using MCTS configurations as described in ??. Each worker also has a delayed copy of the parameters similar to that in IMPALA (?? and [IMPALAScalableDistributed'Espeholt.Soyer.ea'2018]). Step samples and trajectory samples are collected as the planners giving actions and the vectorized environments taking the actions. Each worker is allocated with one CPU and a fraction of a GPU (usually 10% to 20% of a GPU) so neural network inferences could be done on GPU. Collected trajectory samples are returned as the final output of one run of the worker.

4.3.12 Testing Worker

The main goal of **testing workers** is to generate trajectories by interacting with environments for evaluation. These workers are similar to training workers and they hold the same type of data. The differences are: testing workers only use a single environment, have less GPU allocation, and only ran once every other n training steps, where n is a configurable number (usually 5).

4.3.13 Reanalyze Worker

The main goal of **reanalyze workers** is to update search statistics using the reanalyze process described in ??, and push updated trajectories to the replay buffer.

4.3.14 Replay Buffer

The **replay buffer** processes trajectories into training targets and samples trajectories or training targets. Since most training targets are expected to be sampled more than once, the replay buffer precomputes the training targets for all received trajectory samples in the replay buffer with the process described in ??. The replay buffer also computes the value difference δ for each target, which is the difference between the predicted value from the search, and the bootstrapped N-step return (??)

$$\delta_i = |v_i^* - G_i^N|$$

We implemented three modes of sampling: **uniform**, **proportional**, and **rank-based**. In uniform sampling, every training target has equal probability of being drawn. The proportional sampling and rank-based sampling follows the same formula described by **PrioritizedExperienceReplay** [Schaul.Quan.ea'2016] [PrioritizedExperienceReplay' Schaul.Quan.ea'2016]. However, instead of one-step temporal difference error, we use the δ error we described above. For each training target i , the replay buffer also computes the importance sampling ratio $\rho(i)$ based on the probability $P(i)$ of it being drawn

$$\rho_i = \frac{1}{N \cdot P(i)}$$

Since the probabilities of targets depends on other targets as well, the importance sampling ratio of targets are not static, and have to be recomputed each time a batch is sampled from the replay buffer.

4.3.15 Parameter Server

The parameter server holds the central copy of the neural network parameters, performs updates the parameters. Once a batch of training targets is received by the parameter server, the loss and gradients are computed as described in ??, and

4.4 System in Action

Algorithm ?? is the driver that manages all others processes and dataflow. The driver starts by initializing all rollout workers, a parameter server, and a replay buffer. At the beginning of a training step, the driver performs lightweight tasks of all processes such as synchronizing parameters. During the training step, all processes perform their heavyweight tasks. Rollout workers interact with environments, the parameter server computes gradients, and the replay buffer process trajectories into training targets. The method calls made by the driver do not block. They schedule call events and return immediately rather than waiting for the methods to finish. The immediate return values of the calls are *promises* managed by Ray [FuturesPromises'2022]. Actors execute their scheduled method calls sequentially once their concrete inputs are ready. Figure ?? illustrates how tasks are executed in parallel over time.

```

start_training = False
trajectory_samples = []
parameter_server = make_parameter_server()
replay_buffer = make_replay_buffer()
training_worker = [make_train_worker() for i in range(num_train_workers)]
testing_worker = make_test_worker()
reanalyze_workers = [make_reanalyze_worker() for i in range(num_reanalyze_worker)]

for epoch in range(num_epochs):
    if not start_training:
        if start_training_condition_met():
            start_training = True

    # synchronize all workers parameters if conditions are met
    if (epoch % train_worker_update_period) == 0:
        for worker in training_worker:
            worker.set_parameters(parameter_server.get_parameters())

    if (epoch % reanalyze_worker_update_period) == 0:
        for worker in reanalyze_workers:
            worker.set_parameters(parameter_server.get_parameters())

    if (epoch % test_worker_eval_period) == 0:
        testing_worker.set_parameters(parameter_server.get_parameters())

    # process trajectories into training targets
    replay_buffer.process_trajectory_samples(trajectory_samples)

    # update parameters by sampling from replay buffer
    if start_training:
        for i in range(num_updates_per_epoch):
            batch = replay_buffer.sample_batch(batch_size)
            parameter_server.update(batch)

    # generate train targets with training workers
    trajectory_samples.clear()
    for worker in training_worker:
        sample = worker.run()
        trajectory_samples.append(sample)

    # test with testing worker
    if (epoch % test_worker_eval_period) == 0:
        test_result = testing_worker.run()

    # update search statistics with reanalyze workers
    if start_training:
        for worker in reanalyze_workers:
            trajs_to_update = replay_buffer.get_trajectory_samples()
            worker.refill_trajectory(trajs_to_update)
            updated_trajs = worker.run()
            replay_buffer.add_trajs(updated_trajs)

```

Algorithm 3: The driver.

4.5 Logging and Visualization

MooZi incorporates extensive logging and visualization utilities to help users understand its behavior better. All distributed process contains a dedicated log file that records all events within the process. Figure

add figure here

shows an example of the log files. MooZi uses **TensorBoard** to log informative scalars and vectors [TensorFlowLargeScaleMachine'Abadi.Agarwal.ea']. Figure ?? shows the MooZi TensorBoard dashboard. MooZi also provides utilities to visualize

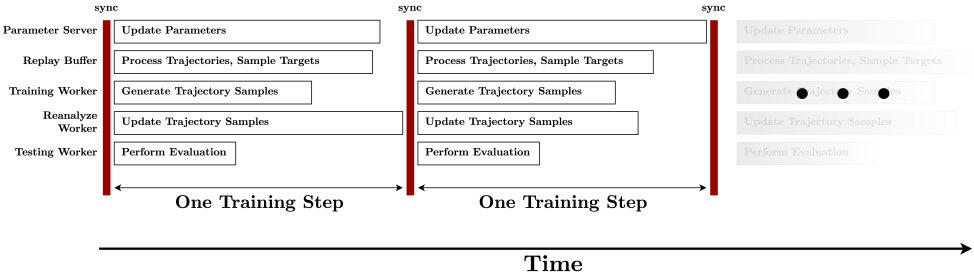


Figure 4.5: Timeline of training steps. The red bar indicates a synchronization barrier. The duration of each training step is decided by the last finished task.

the behavior of the algorithm. Testing workers use the GIF maker tool to create animated records of evaluations. Figure

add figure here

shows an example of the GIF tiled as a sprite image.

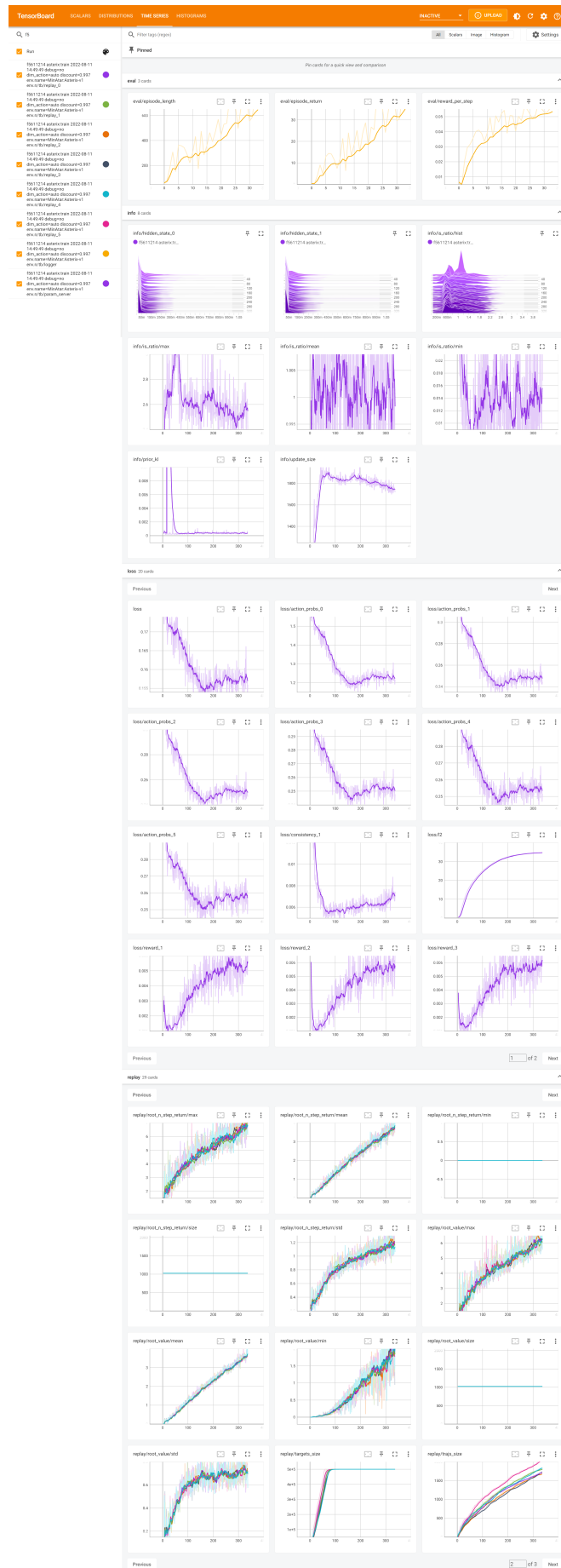


Figure 4.6: MooZi Tensorboard dashboard. MooZi logs over 50 different metrics into this dashboard. These include 29 metrics from the training workers, 3 metrics from the testing logger, 20 metrics from decomposing the loss function, 8 metrics from the parameter optimizer, and optionally shows distributions of parameters from all layers.

Chapter 5

Experiments

5.1 Experiment Setup

5.1.1 Basics

can't think of a better title.

For all of our experiments, we use unrolled steps $K = 5$, history length $L = 4$, bootstrap steps $N = 10$, a discount of 0.997, and a support Z from the interval $[-30, 30]$ ($|Z| = 61$). We use a single computer with Intel Xeon CPUs (72×2.3 GHz), Nvidia Tesla V100 GPUs (8×32 GB), and 500 Gigabytes of system memory. A running MooZi system roughly uses 50% to 75% of the CPUs, 40% to 60% of the GPUs, and 25% of the memory.

5.1.2 Neural Network Configurations

We use the residual-blocks-based variant of the network for all of our experiments.

Residual Block

We follow the residual block definition by **DeepResidualLearning** [He.Zhang.ea'2016]. One residual block is defined as follows:

- input x
- save a copy of x to x'
- apply a 2-D padded convolution on x , with kernel size 3 by 3, same channels
- apply batch normalization on x
- apply relu activation on x
- apply a 2-D padded convolution on x , with kernel size 3 by 3, same channels
- apply batch normalization on x
- add x' to x
- apply relu activation on x

The Representation Function

The representation function h is defined as follows

- input stacked history ψ of shape (H, W, C_h)
- apply a 2-D padded convolution on ψ , with kernel size 1 by 1, 32 channels
- apply 6 residual blocks with 32 channels on ψ
- output the hidden state \mathbf{x} of shape $(H, W, 32)$

The Prediction Function

The prediction function f is parametrized as follows

- input hidden state \mathbf{x} of shape $(H, W, 32)$
- apply 1 residual block with 32 channels on \mathbf{x}
- flatten \mathbf{x} , now shape $(H * W * 32)$
- apply 1 dense layer with output size of 128 on flattened \mathbf{x} to obtain the value head \mathbf{x}_v , now shape (128)
- apply batch normalization and relu activation on \mathbf{x}_v
- apply 1 dense layer with output size of $|Z|$ on \mathbf{x}_v , now shape (1)
- output the value head \mathbf{x}_v as the value prediction v
- apply 1 dense layer with output size of 128 on flattened \mathbf{x} to obtain the policy head \mathbf{x}_p , now shape (128)
- apply batch normalization and relu activation on \mathbf{x}_p
- apply 1 dense layer with output size equals to the actin space size on \mathbf{x}_p , now shape $(|\mathcal{A}^a|)$
- output the policy head as the policy prediction \mathbf{p}

The Dynamics Function

The dynamics function g is parametrized as follows:

- input hidden state \mathbf{x} of shape $(H, W, 32)$, action a as an integer
- encode a as action planes of shape $(H, W, |\mathcal{A}^a|)$ (same procedure as ??)
- stack \mathbf{x} on top of the encoded action, now shape $(H, W, 32 + |\mathcal{A}^a|)$
- apply a 2-D padded convolution on ψ , with kernel size 1 by 1, 32 channels, now shape $(H, W, 32)$
- apply 1 residual block with 32 channels on \mathbf{x} to obtain the hidden state head \mathbf{x}_s
- apply 1 residual block with 32 channels on the hidden state head \mathbf{x}_s
- output the hidden state head \mathbf{x}_s as the next hidden state \mathbf{x}'

Worker	c_1	c_2	Temperature	Dirichlet Noise	Simulations
Training	2.25	19652	1.0	0.2	25
Reanalyze	1.75	19652	-	0.2	50
Testing	1.75	19652	0.25	0.1	40

Table 5.1: Planner configurations. Reanalyze workers do not sample actions to act so the temperature parameter does not affect their behavior. c_1 is the exploration constant in the PUCT formula form ???. A greater c_1 favors the less visited actions more. c_2 is also an exploration constant in the PUCT formula. We use the same value as [MasteringAtariGo'Schrittwieser.Antonoglou.ea'2020] because we find it sufficient to tune c_1 to balance exploration. The *temperature* controls how action is selected from the distribution of action visit counts from the root nodes. A temperature of 0 means select most visited actions at the root nodes. A temperature of ∞ means select actions at random. The *dirichlet noise* controls the exploitation noise added to the actions at the root nodes. A greater dirichlet noise adds more prior probabilities to less explored actions. The *simulations* is the number of simulations the planner performs at each timestep ???. The training workers favor exploration and generate data quickly with less simulations. The testing workers favor exploitation and spend more time on simulations to get better average return. The reanalyze workers do not need to interact with the environment so they spend more time performing simulations.

- apply 1 dense layer with output size of 128 on \mathbf{x} to obtain the reward head \mathbf{x}_r , now shape of (128)
- apply batch normalization and relu activation on \mathbf{x}_r
- apply 1 dense layer with output size of Z on \mathbf{x}_r , now shape of ($|Z|$)
- output reward head \mathbf{x}_r as the reward prediction \hat{r}

The Projection Function

The projection function q is simply one residual block.

Network Training

In the loss function (??), we use $c^v = 0.25$, $c^s = 2.0$, and $c^{L_2} = 1.0 \times 10^{-4}$. We use a batch size of 1024, a learning rate of 1.0×10^{-2} , and a global norm clipping of 5.0. We perform gradient updates with samples four times of the number of step samples generated each training step. For example, if we have 30 training workers, each with 16 environments, and each performs 100 environment steps per training step, then the total number step samples generated is $30 * 16 * 100 = 48000$. This means we update the gradient using $4 * 48000 = 192000$ sampled training targets from the replay buffer. That is $\frac{192000}{1024} \approx 188$ gradient updates from mini-batches of size 1024 per training step. We use a target network similar to DQN (??) to smooth gradient updates, and we overwrites the target network every other 500 gradient updates.

5.1.3 Planners Configurations

Table ?? shows the configurations of planners from different type of workers.

5.1.4 Driver Configuration

We use 30 training workers, each with 16 copies of the environment. For every training step, each training worker performs 100 environment steps. Freeway environment has much longer episodes, so each training worker uses 2 environments and performs 2500 environment steps per training step. We use 1 testing worker with 1 copy of the environment. For every 10 training steps, the testing worker collects 10 trajectories and logs the average return of the trajectories. We do not use the reanalyze except for the experiments in ???. All workers sync with the latest neural network parameters per 10 training steps.

5.2 MinAtar Games vs PPO

We run MooZi using four environments in MinAtar [MinAtarAtariInspiredTestbed`Young.Tian`2019]. All environments produce frames of resolution 10×10 , with four to seven environment channels C_e . In **Breakout**, the player moves a paddle left or right on the bottom of the screen to bounce a ball to hit the bricks. Reward is +1 for each brick broken and 0 in all other situations. The game ends when the paddle fails to catch the ball. In **Space Invaders**, the player controls a cannon that can move left, move right, or fire the cannon. A cluster of enemies move across the screen and fire at the player. Reward is +1 for each enemy hit by the player and 0 in all other situations. The game ends when the play is hit by an enemy's bullet. In **Freeway**, the player starts at the bottom of the screen, moves up or down once every three frames to travel across a road with traffic. Cars are spawn randomly and travel horizontally across the screen, and when they hit the player the player is moved back to its starting position. Reward is +1 for each time the player successfully travel across the road and 0 in all other situations. The game ends after 2500 timesteps. In **Asterix**, the player moves in four directions, and enemies and treasures spawn randomly on the edges. Reward is +1 if the player obtains a treasure and 0 in all other situations. The game ends when the player bump into an enemy. Enemies have different speed indicated by the color of their tail. MuZero is reported to have much better performance in more deterministic environments [VectorQuantizedModels`Ozair.Li.ea`2021]. To compare MooZi in more deterministic environments, we compare the results reported by Gymnax, in which algorithms are evaluated in environments with no sticky action. We discuss this difference in section ???. Gymnax benchmarked the performance of PPO with three set of hyper-parameters in each of these four games [Gymnax]AXbasedReinforcement`RobertTjarkoLange`2022, ProximalPolicyOptimization`Schulman.Wolski.ea`2017]. We compare our results with the best performing PPO in each of these games. The average return uses the planner setting of a test worker (??). Figure ?? shows the result comparisons.

Do I need to show what the environments look like? Those will be static images so the readers will be clueless still...

5.3 Sticky Actions in MinAtar

MinAtar environments have a default sticky action probabilities of 10%. This means one out of ten timesteps, the environment uses the last taken action to step

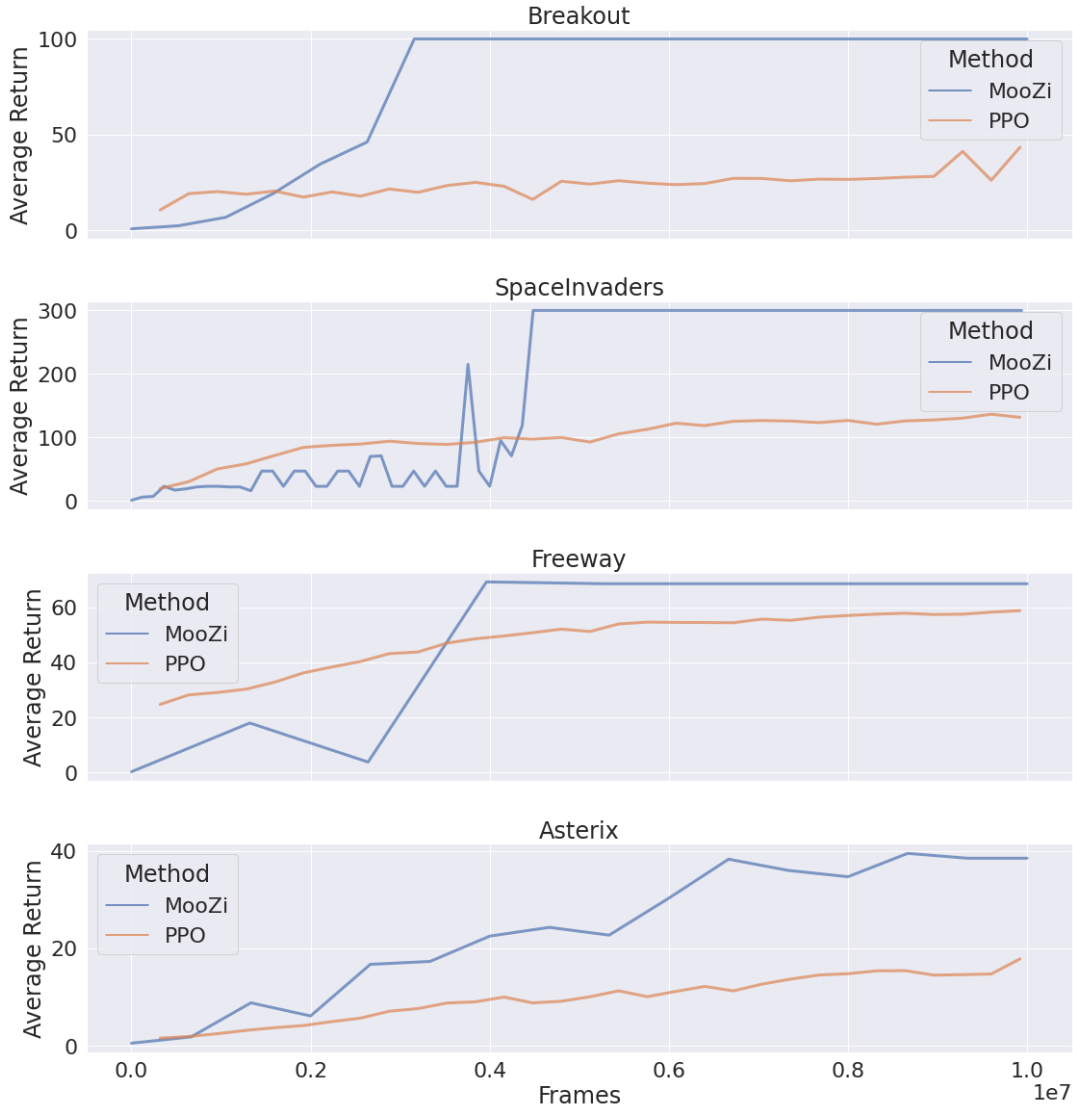


Figure 5.1: MinAtar games. MooZi vs PPO. In *Breakout*, MooZi obtained a near-optimal strategy and the paddle almost never fails to catch the ball. There is no default step limits in the environments so we cap the return at 100. Similarly in *Space Invaders*, MooZi obtained a near-optimal strategy and we cap the return at 300. In *Freeway*, MooZi also obtained a near-optimal strategy. Since episode in this environment is much longer than other environments, MooZi cycles through less training steps and save less checkpoints. Hence MooZi has much less data points in the line plot for this environment. In *Asterix*, both MooZi and PPO do not obtain an optimal strategy, but MooZi have twice as much average return as PPO.

the environment instead of using the agent output action. For example, at time t , the agent outputs action $a_t = \text{MoveLeft}$ and moves to the left. At time $t + 1$, the agent outputs action $a_{t+1} = \text{MoveRight}$. However, this time, the environment applies the sticky action and overrides the action $a_{t+1} = \text{MoveLeft}$, and the agent moves to the left even further. The presence of a non-zero sticky action probability adds stochasticity to environments and changes the set of optimal policies. For example, in Space Invaders, if sticky action probability is 0, then the agent can move away from enemy bullets one frame before the agent is about to get hit. However, with a sticky action probability 10%, moving away one frame in advance means there's a 10% chance that the agent will die, and moving away two frames in advance means there's a $(10\% * 10\%) = 1\%$ chance that the agent will die two frames later. We observe that a MooZi agent trained in Space Invaders with a sticky actions probability moves away from an enemy bullet right after the bullet is visible on the screen. **MinAtarAtariInspiredTestbed`Young.Tian`2019** shipped the MinAtar environments with four algorithms including two variants of DQN and two variants of Actor-Critic (AC) method. MinAtar's paper only reports Breakout results with sticky action probability of 10% and Gymnax only reports Breakout results without sticky action probability. We compare with MinAtar and Gymnax in their respective testing environments using the same MooZi agent configuration. Figure ?? shows the comparison.

5.4 Search Budget and Testing Strength

We use the trained MooZi model in Space Invaders environment to evaluate its strength using different number of simulations. We only use model checkpoints from the first 3 million environment frames because after that the agent behaviors optimally even just using the prior to act. For each of these checkpoints, we runs a testing worker to collect 30 episodes and we calculate the average return of these episodes. The testing worker uses the same planner configuration as the testing worker in ?? except for the number of simulations. Figure ?? shows the result.

Gymnax's results were based on environments variants without sticky actions, while the original MinAtar environments have a default sticky action probability of 10%.

5.4.1 Sample Efficiency with Reanalyze

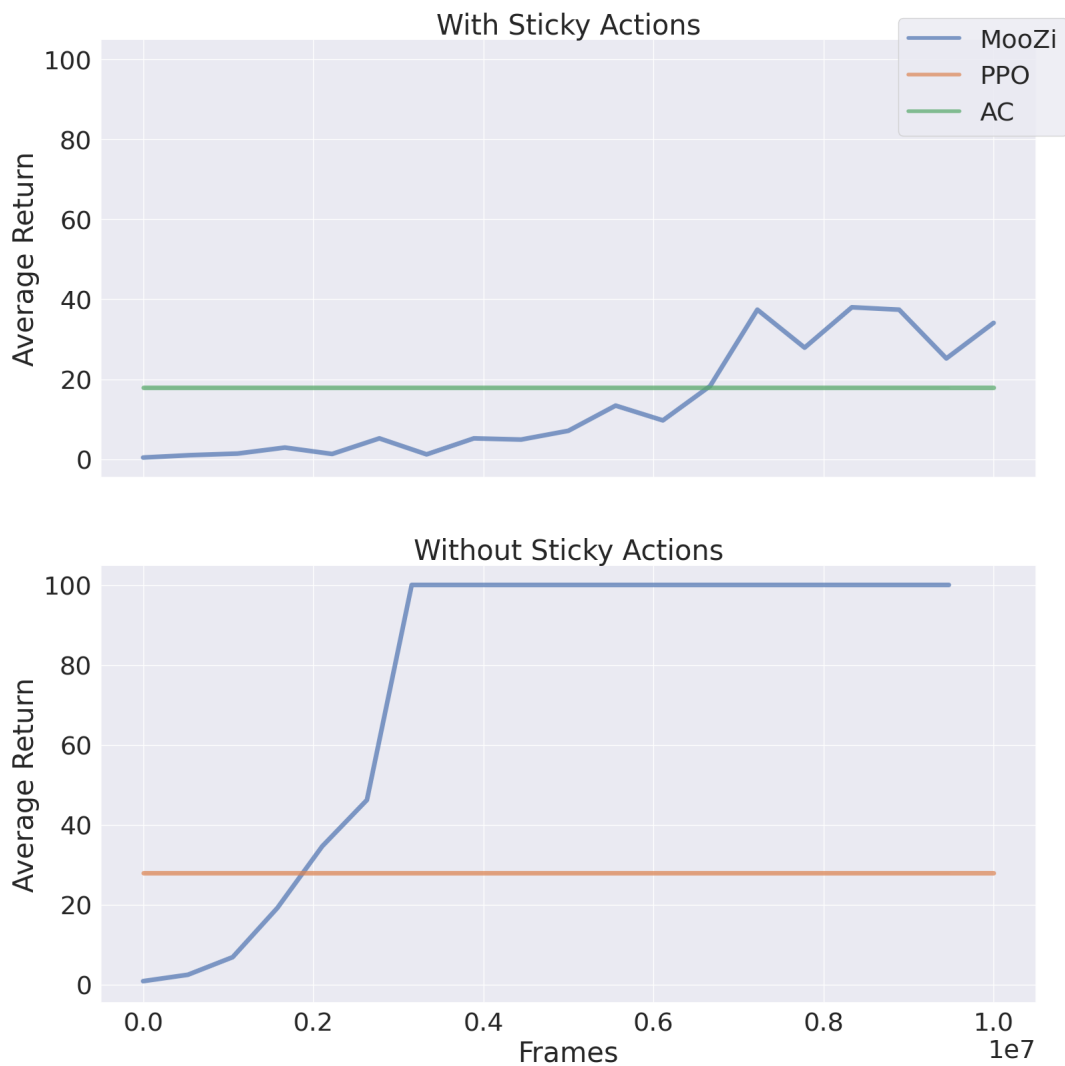


Figure 5.2: MinAtar Breakout with or without sticky actions. MooZi vs PPO vs AC. We use the final average returns of PPO and AC reported in Gymnax and MinAtar’s paper respectively. In Breakout with sticky actions, MooZi learns slower and the final average return is much lower. In Breakout without sticky actions, MooZi quickly learns a near-optimal policies and never fails to return the ball. In both environments, MooZi out-performs the other algorithm.

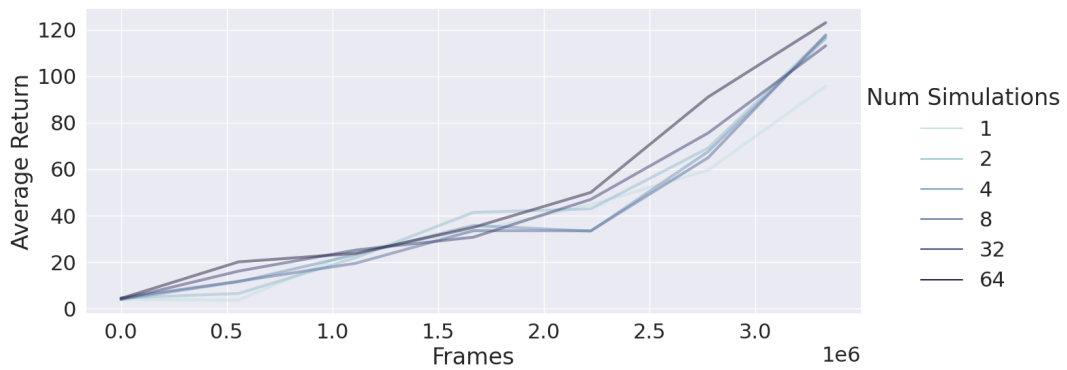


Figure 5.3: Agent strength with different number of simulations. We observe that with a greater number of simulations, the agent tends to perform better. With more training, the agent always perform better. The difference of performance due to simulation count seems to be smaller than the difference due to training. For example, with 2.5 million frames of training, acting using 64 simulations gives an average return around 70. With 3 million frames of training, acting according to the prior (1 simulation) does even better and gives an average return around 75. The prior policy, with a bit more training, quickly catches up to or even exceeds the deep search policy with less training. This finding aligns with the analysis by **ROLEPLANNINGMODELBASED’Hamrick.Friesen.ea’2021** [ROLEPLANNINGMODELBASED’Hamrick.Friesen.ea’2021]: the planner contributes more to the algorithm by *generating better data for training the model* rather than *exploiting the model for better testing*.

Chapter 6

Conclusion

(3 pages)

Chapter 7

Future Work

(1 page)