

Reinforcement Learning: An Overview

Kevin P. Murphy

December 3, 2025

Brief Table of Contents

1	Introduction	13
1.1	Sequential decision making	13
1.2	Canonical models	17
1.3	Reinforcement Learning: a high-level summary	22
2	Value-based RL	31
2.1	Basic concepts	31
2.2	Solving for the optimal policy in a known world model	33
2.3	Value function learning using samples from the world model	35
2.4	SARSA: on-policy TD policy learning	38
2.5	Q-learning: off-policy TD policy learning	39
3	Policy-based RL	51
3.1	Policy gradient methods	51
3.2	Actor-critic methods	55
3.3	Policy improvement methods	64
3.4	Off-policy methods	68
3.5	Gradient-free policy optimization	72
3.6	RL as inference	73
4	Model-based RL	83
4.1	Introduction	83
4.2	Decision-time (online) planning	84
4.3	Background (offline) planning	93
4.4	World models	97
4.5	Beyond one-step models: predictive representations	111
5	Multi-agent RL	121
5.1	Games	121
5.2	Solution concepts	126
5.3	Algorithms	133
6	LLMs and RL	147
6.1	Introduction	147
6.2	RL for LLMs	147
6.3	LLMs for RL	160
6.4	Implementation details	165

7	Other topics in RL	173
7.1	Regret minimization	173
7.2	Exploration-exploitation tradeoff	175
7.3	Distributional RL	180
7.4	Intrinsic motivation for reward-free RL	181
7.5	Hierarchical RL	183
7.6	Imitation learning	190
7.7	Offline RL	192
7.8	General RL, AIXI and universal AGI	197
8	Acknowledgements	199

Contents

1	Introduction	13
1.1	Sequential decision making	13
1.1.1	Maximum expected utility principle	13
1.1.2	Episodic vs continual tasks	14
1.1.3	Universal model	15
1.1.4	Further reading	17
1.2	Canonical models	17
1.2.1	Partially observed MDPs	17
1.2.2	Markov decision process (MDPs)	18
1.2.3	Goal-conditioned MDPs	19
1.2.4	Contextual MDPs	19
1.2.5	Contextual bandits	19
1.2.6	Belief state MDPs	20
1.2.7	Optimization problems as decision problems	21
1.2.7.1	Best-arm identification	22
1.2.7.2	Bayesian optimization	22
1.2.7.3	Active learning	22
1.2.7.4	Stochastic Gradient Descent (SGD)	22
1.3	Reinforcement Learning: a high-level summary	22
1.3.1	Value-based RL	23
1.3.2	Policy-based RL	24
1.3.3	Model-based RL	24
1.3.4	State uncertainty (partial observability)	24
1.3.4.1	Optimal solution	25
1.3.4.2	Finite observation history	25
1.3.4.3	Stateful (recurrent) policies	25
1.3.5	Model uncertainty (exploration-exploitation tradeoff)	25
1.3.6	Reward functions	26
1.3.6.1	The reward hypothesis	26
1.3.6.2	Non-Markovian rewards	27
1.3.6.3	Reward hacking	27
1.3.6.4	Sparse reward	27
1.3.6.5	Reward shaping	28
1.3.6.6	Intrinsic reward	28
1.3.7	Best practices for experimental work in RL	28
2	Value-based RL	31
2.1	Basic concepts	31
2.1.1	Value functions	31
2.1.2	Bellman's equations	31

2.1.3	Example: 1d grid world	32
2.2	Solving for the optimal policy in a known world model	33
2.2.1	Value iteration	33
2.2.2	Real-time dynamic programming (RTDP)	34
2.2.3	Policy iteration	34
2.3	Value function learning using samples from the world model	35
2.3.1	Monte Carlo estimation	35
2.3.2	Temporal difference (TD) learning	36
2.3.3	Combining TD and MC learning using $TD(\lambda)$	36
2.3.4	Eligibility traces	38
2.4	SARSA: on-policy TD policy learning	38
2.4.1	Convergence	38
2.4.2	Sarsa(λ)	39
2.5	Q-learning: off-policy TD policy learning	39
2.5.1	Tabular Q learning	39
2.5.2	Q learning with function approximation	41
2.5.2.1	Neural fitted Q	41
2.5.2.2	DQN	41
2.5.2.3	Experience replay	42
2.5.2.4	Prioritized experience replay	42
2.5.2.5	The deadly triad	42
2.5.2.6	Target networks	43
2.5.2.7	Gradient TD methods	44
2.5.2.8	Two time-scale methods	44
2.5.2.9	Layer norm	44
2.5.2.10	Other methods	45
2.5.3	Maximization bias	45
2.5.3.1	Double Q-learning	45
2.5.3.2	Double DQN	45
2.5.3.3	Randomized ensemble DQN	46
2.5.4	DQN extensions	46
2.5.4.1	Q learning for continuous actions	46
2.5.4.2	Dueling DQN	46
2.5.4.3	Noisy nets and exploration	47
2.5.4.4	Multi-step DQN	47
2.5.4.5	$Q(\lambda)$	47
2.5.4.6	Rainbow	48
2.5.4.7	Bigger, Better, Faster	48
2.5.4.8	Other methods	49
2.5.5	Q-learning for GCRL using hindsight relabeling	49
3	Policy-based RL	51
3.1	Policy gradient methods	51
3.1.1	Likelihood ratio estimate	51
3.1.2	Variance reduction using reward-to-go	52
3.1.3	REINFORCE	53
3.1.4	The policy gradient theorem	53
3.1.5	Variance reduction using a baseline	54
3.1.6	REINFORCE with baseline	55
3.2	Actor-critic methods	55
3.2.1	Advantage actor critic (A2C)	55
3.2.2	Generalized advantage estimation (GAE)	57

3.2.3	Two-time scale actor critic algorithms	58
3.2.4	Natural policy gradient methods	58
3.2.4.1	Natural gradient descent	58
3.2.4.2	Natural actor critic	60
3.2.5	Architectural issues	60
3.2.6	Deterministic policy gradient methods	61
3.2.6.1	Deterministic policy gradient theorem	61
3.2.6.2	DDPG	62
3.2.6.3	Twin Delayed DDPG (TD3)	62
3.2.6.4	Wasserstein Policy Optimization (WPO)	62
3.3	Policy improvement methods	64
3.3.1	Policy improvement lower bound	64
3.3.2	Trust region policy optimization (TRPO)	65
3.3.3	Proximal Policy Optimization (PPO)	66
3.3.3.1	Simplified form of the clipping term	66
3.3.3.2	PPO for diffusion policies	67
3.3.3.3	Simple policy optimization	67
3.3.4	Variational Maximum a Posteriori Policy Optimization (VMPO)	67
3.4	Off-policy methods	68
3.4.1	Policy evaluation using importance sampling	68
3.4.2	Off-policy actor critic methods	69
3.4.2.1	Learning the critic using V-trace	69
3.4.2.2	Learning the actor	70
3.4.2.3	Example: IMPALA	71
3.4.2.4	Off-policy learning with deterministic policies	72
3.4.2.5	PGQL: Combining off-policy Q-learning with policy gradient	72
3.4.3	Off-policy policy improvement methods	72
3.5	Gradient-free policy optimization	72
3.6	RL as inference	73
3.6.1	Deterministic case (planning/control as inference)	74
3.6.2	Stochastic case (policy learning as variational inference)	74
3.6.3	EM control	75
3.6.4	KL control (maximum entropy RL)	76
3.6.5	Maximum a Posteriori Policy Optimization (MPO)	76
3.6.6	Sequential Monte Carlo Policy Optimisation (SMC-PO)	77
3.6.7	AWR and AWAC	77
3.6.8	Soft Actor Critic (SAC)	77
3.6.8.1	SAC objective	77
3.6.8.2	Policy evaluation: tabular case	77
3.6.8.3	Policy evaluation: general case	78
3.6.8.4	Policy improvement	79
3.6.8.5	Adjusting the temperature	79
3.6.9	Active inference	81
4	Model-based RL	83
4.1	Introduction	83
4.2	Decision-time (online) planning	84
4.2.1	Receding horizon control	84
4.2.1.1	Forward search	85
4.2.1.2	Branch and bound	85
4.2.1.3	Sparse sampling	86
4.2.1.4	Heuristic search	86

4.2.2	Monte Carlo tree search (MCTS)	86
4.2.2.1	MCTS for 2p0s games: AlphaGo, AlphaGoZero, and AlphaZero	87
4.2.2.2	MCTS with learned world model: MuZero and EfficientZero	88
4.2.2.3	MCTS in belief space	89
4.2.3	Sequential Monte Carlo (SMC) for online planning	89
4.2.4	Model predictive control (MPC), aka open loop planning	90
4.2.4.1	Suboptimality of open-loop planning for stochastic environments	91
4.2.4.2	Trajectory optimization	92
4.2.4.3	LQR	92
4.2.4.4	Random shooting	92
4.2.4.5	CEM	92
4.2.4.6	MPPI	93
4.2.4.7	GP-MPC	93
4.3	Background (offline) planning	93
4.3.1	A game-theoretic perspective on MBRL	93
4.3.2	Dyna	95
4.3.2.1	Tabular Dyna	95
4.3.2.2	Dyna with function approximation	95
4.4	World models	97
4.4.1	World models which are trained to predict observation targets	97
4.4.1.1	Generative world models without latent variables	98
4.4.1.2	Generative world models with latent variables	98
4.4.1.3	Example: Dreamer	98
4.4.1.4	Example: IRIS	101
4.4.1.5	Code world models	101
4.4.1.6	Partial observation prediction	101
4.4.2	World models which are trained to predict other targets	101
4.4.2.1	The objective mismatch problem	102
4.4.2.2	Observation prediction	102
4.4.2.3	Reward prediction	103
4.4.2.4	Value prediction	103
4.4.2.5	Policy prediction	104
4.4.2.6	Self prediction (self distillation)	104
4.4.2.7	Avoiding self-prediction collapse using frozen targets	104
4.4.2.8	Avoiding self-prediction collapse using information-theoretic regularization	105
4.4.2.9	Preventing self-prediction collapse using game-theoretic approaches	106
4.4.2.10	Example: JEPA	107
4.4.2.11	Example: DinoWM	108
4.4.2.12	Example: TD-MPC	108
4.4.2.13	Example: BYOL	109
4.4.2.14	Example: Imagination-augmented agents	110
4.4.3	World models that are trained to help planning	110
4.4.4	Dealing with model errors and uncertainty	110
4.4.4.1	Avoiding compounding errors in rollouts	110
4.4.4.2	Unified model and planning variational lower bound	111
4.4.4.3	Dynamically switching between MFRL and MBRL	111
4.4.5	Exploration for learning world models	111
4.5	Beyond one-step models: predictive representations	111
4.5.1	General value functions	112
4.5.2	Successor representations	112
4.5.3	Successor features	115
4.5.3.1	Generalized policy improvement	116

4.5.3.2	Option keyboard	116
4.5.3.3	Learning SFs	117
4.5.3.4	Choosing the tasks	117
4.5.4	Successor measures	117
4.5.4.1	Learning SMs	118
4.5.4.2	Jumpy models using geometric policy composition	119
4.5.4.3	Other related work	119
4.5.5	Connection between options and successor representations	119
5	Multi-agent RL	121
5.1	Games	121
5.1.1	Normal-form games	121
5.1.2	Stochastic games	123
5.1.3	Partially observed stochastic games (POSG)	123
5.1.3.1	Data generating process	124
5.1.3.2	Objective	124
5.1.3.3	Single agent perspective	125
5.1.3.4	Factored Observation Stochastic Games (FOSG)	125
5.1.4	Extensive form games (EFG)	125
5.1.4.1	Example: Kuhn Poker as EFG	125
5.1.4.2	Converting FOSG to EFG	126
5.2	Solution concepts	126
5.2.1	Notation and definitions	127
5.2.2	Minimax	127
5.2.3	Exploitability	128
5.2.4	Nash equilibrium	128
5.2.5	Approximate Nash equilibrium	128
5.2.6	Entropy regularized Nash equilibria (aka Quantal Response Equilibria)	129
5.2.7	Correlated equilibrium	129
5.2.8	Limitations of equilibrium solutions	130
5.2.9	Pareto optimality	130
5.2.10	Social welfare and fairness	131
5.2.11	No regret	131
5.2.12	Shapley values	132
5.2.13	Stackelberg equilibrium	132
5.3	Algorithms	133
5.3.1	Centralized learning	133
5.3.2	Independent learning	133
5.3.2.1	Independent Q learning	133
5.3.2.2	Independent Actor Critic	134
5.3.2.3	Independent PPO	135
5.3.2.4	Learning dynamics of multi-agent policy gradient methods	135
5.3.3	Centralized training of decentralized policies (CTDE)	135
5.3.3.1	Application to Diplomacy (Cicero)	136
5.3.4	Value decomposition methods for common-reward games	136
5.3.4.1	Value decomposition network (VDN)	137
5.3.4.2	QMIX	137
5.3.5	Policy learning with self-play	137
5.3.6	Policy learning with learned opponent models	138
5.3.7	Best response	138
5.3.7.1	Fictitious play	138
5.3.7.2	Neural fictitious self play (NFSP)	139

5.3.8	Population-based training	139
5.3.8.1	PSRO (policy space response oracle)	139
5.3.8.2	Application to StarCraft (AlphaStar)	140
5.3.9	Counterfactual Regret Minimization (CFR)	140
5.3.9.1	Tabular case	141
5.3.9.2	Deep CFR	141
5.3.9.3	Applications to Poker and other games	141
5.3.10	Regularized policy gradient methods	142
5.3.10.1	Magnetic Mirror Descent (MMD)	142
5.3.10.2	PPO	142
5.3.11	Decision-time planning methods	143
5.3.11.1	Magnetic Mirror Descent Search (MMDS)	143
5.3.11.2	Belief state approximations	144
5.3.11.3	Experiments	144
5.3.11.4	Open questions	145
5.3.12	MARL for LLM agents	145
6	LLMs and RL	147
6.1	Introduction	147
6.2	RL for LLMs	147
6.2.1	RL fine tuning (RLFT)	147
6.2.2	Reward models	148
6.2.2.1	RL with verifiable rewards (RLVR)	148
6.2.2.2	Process vs outcome reward models	148
6.2.2.3	Learning the reward model from human feedback (RLHF)	148
6.2.2.4	Learning the reward model from AI feedback (RLAIF)	149
6.2.2.5	Generative reward models (GRM)	149
6.2.3	Agents which “think”	149
6.2.3.1	Chain of thought prompting	149
6.2.3.2	Training a thinking model using RL	149
6.2.3.3	Thinking as marginal likelihood maximization	150
6.2.3.4	Can we bootstrap a model to think from scratch?	150
6.2.3.5	Agentic AI	150
6.2.4	Algorithms for single-turn RL	150
6.2.4.1	Problem setup	150
6.2.4.2	PPO	151
6.2.4.3	GRPO	151
6.2.4.4	DAPO	152
6.2.4.5	GSPO	152
6.2.4.6	RLOO	153
6.2.4.7	REINFORCE++	153
6.2.4.8	VinePPO	153
6.2.4.9	Adding a KL regularizer	154
6.2.4.10	DPO	154
6.2.4.11	Inference-time scaling using posterior sampling	155
6.2.4.12	RLFT as amortized posterior sampling	156
6.2.5	Algorithms for multi-turn RL	157
6.2.5.1	Example: RAGEN	157
6.2.5.2	Dealing with invalid actions	158
6.2.5.3	Turn-level training	158
6.2.5.4	Self-play for LLM training	159
6.2.6	Alignment and the assistance game	160

6.3	LLMs for RL	160
6.3.1	LLMs for pre-processing the input	160
6.3.1.1	Example: AlphaProof	161
6.3.1.2	VLMs for parsing images into structured data	161
6.3.1.3	Active control of LLM sensor/preprocessor	161
6.3.2	LLMs for rewards	161
6.3.3	LLMs for world models	162
6.3.3.1	LLMs as world models	162
6.3.3.2	LLMs for generating code world models	162
6.3.3.3	LLMs for generating partial code world models	163
6.3.4	LLMs for policies	163
6.3.4.1	LLMs for generating actions	163
6.3.4.2	LLMs for generating code policies	164
6.3.4.3	LLMs for generating code actions	164
6.3.4.4	In-context RL	164
6.3.5	Speeding up LLMs	165
6.3.5.1	Computational complexity of transformer models	165
6.3.5.2	Modern RNNs	165
6.4	Implementation details	165
6.4.1	Policy gradient using Tinker	166
6.4.2	Rolling out episodes	168
6.4.3	Computing the advantages	168
6.4.4	Computing token level loss	169
6.4.5	Computing metrics related to training stability	169
6.4.6	Example	170
7	Other topics in RL	173
7.1	Regret minimization	173
7.1.1	Regret for static MDPs	173
7.1.2	Regret for non-stationary MDPs	174
7.1.3	Minimizing regret vs maximizing expected utility	174
7.2	Exploration-exploitation tradeoff	175
7.2.1	Optimal (Bayesian) approach	175
7.2.1.1	Bandit case (Gittins indices)	176
7.2.1.2	MDP case (Bayes Adaptive MDPs)	176
7.2.2	Thompson sampling	176
7.2.2.1	Bandit case	177
7.2.2.2	MDP case (posterior sampling RL)	177
7.2.3	Upper confidence bounds (UCBs)	178
7.2.3.1	Basic idea	178
7.2.3.2	Bandit case: Frequentist approach	179
7.2.3.3	Bandit case: Bayesian approach	179
7.2.3.4	MDP case	179
7.3	Distributional RL	180
7.3.1	Quantile regression methods	180
7.3.2	Replacing regression with classification	180
7.4	Intrinsic motivation for reward-free RL	181
7.4.1	Knowledge-based intrinsic motivation	181
7.4.1.1	Exploration bonuses	181
7.4.1.2	Random Network Distillation (RND)	181
7.4.1.3	Information-theoretic measures	182
7.4.2	Competence-based intrinsic motivation	182

7.4.2.1	Empowerment	182
7.4.2.2	Curriculum design	183
7.4.2.3	Using an LLM to choose goals	183
7.4.2.4	Go-Explore	183
7.5	Hierarchical RL	183
7.5.1	HRL using Options	183
7.5.1.1	Introduction	183
7.5.1.2	Option hierarchies	185
7.5.1.3	Hierarchical Q learning	185
7.5.1.4	MAXQ	186
7.5.1.5	Option learning using EM	186
7.5.1.6	Skill chaining	186
7.5.1.7	Option critic	186
7.5.1.8	Double actor critic (DAC)	186
7.5.1.9	Avoiding excessive (or insufficient) option switching	187
7.5.1.10	MBRL using options	187
7.5.2	HRL using feudal hierarchies	187
7.5.2.1	Introduction	187
7.5.2.2	Comparison with options	187
7.5.2.3	Feudal Q learning	188
7.5.2.4	Dealing with nonstationarity using hindsight relabeling (HIRO, HAC)	188
7.5.2.5	Learning the goal space and policy	189
7.5.3	Subtask discovery	189
7.5.3.1	Discovery of subgoals	189
7.5.3.2	Discovery of skills	190
7.6	Imitation learning	190
7.6.1	Imitation learning by behavior cloning	191
7.6.2	Imitation learning by inverse reinforcement learning	191
7.6.3	Imitation learning by divergence minimization	192
7.7	Offline RL	192
7.7.1	Offline model-free RL	193
7.7.1.1	Policy constraint methods	193
7.7.1.2	Behavior-constrained policy gradient methods	194
7.7.1.3	Uncertainty penalties	194
7.7.1.4	Conservative Q-learning	195
7.7.2	Offline model-based RL	195
7.7.3	Offline RL using reward-conditioned sequence modeling	196
7.7.4	Offline-to-online methods	196
7.7.4.1	Calibrated Q learning	197
7.7.4.2	Dagger	197
7.8	General RL, AIXI and universal AGI	197
8	Acknowledgements	199

Chapter 1

Introduction

1.1 Sequential decision making

Reinforcement learning or **RL** is a class of methods for solving various kinds of sequential decision making tasks. In such tasks, we want to design an **agent** that interacts with an external **environment**. The agent maintains an internal state z_t , which it passes to its **policy** π to choose an action $a_t = \pi(z_t)$. The environment responds by sending back an observation o_{t+1} , which the agent uses to update its internal state using the state-update function $z_{t+1} = SU(z_t, a_t, o_{t+1})$. See Figure 1.1 for an illustration.

To simplify things, we often assume that the environment is also a Markovian process, which has internal world state w_t , from which the observations o_t are derived. (This is called a POMDP — see Section 1.2.1). We often simplify things even more by assuming that the observation o_t reveals the hidden environment state; in this case, we denote the internal agent state and external environment state by the same letter, namely $s_t = o_t = w_t = z_t$. (This is called an MDP — see Section 1.2.2). We discuss these assumptions in more detail in Section 1.1.3.

RL is more complicated than supervised learning (e.g., training a classifier) or self-supervised learning (e.g., training a language model), because this framework is very general: there are many assumptions we can make about the environment and its observations o_t , and many choices we can make about the form the agent’s internal state z_t and policy π , as well the ways to update these objects as we see more data. We will study many different combinations in the rest of this document. The right choice ultimately depends on which real-world application you are interested in solving.¹

1.1.1 Maximum expected utility principle

The goal of the agent is to choose a policy π so as to maximize the sum of expected rewards:

$$V_\pi(s_0) = \mathbb{E}_{p(a_0, s_1, a_1, \dots, a_T, s_T | s_0, \pi)} \left[\sum_{t=0}^T R(s_t, a_t) | s_0 \right] \quad (1.1)$$

where s_0 is the agent’s initial state, $R(s_t, a_t)$ is the **reward function** that the agent uses to measure the value of performing an action in a given state, $V_\pi(s_0)$ is the **value function** for policy π evaluated at s_0 , and the expectation is wrt

$$p(a_0, s_1, a_1, \dots, a_T, s_T | s_0, \pi) = \pi(a_0 | s_0) p_{\text{env}}(o_1 | a_0) \delta(s_1 = U(s_0, a_0, o_1)) \quad (1.2)$$

$$\times \pi(a_1 | s_1) p_{\text{env}}(o_2 | a_1, o_1) \delta(s_2 = U(s_1, a_1, o_2)) \quad (1.3)$$

$$\times \pi(a_2 | s_2) p_{\text{env}}(o_3 | a_{1:2}, o_{1:2}) \delta(s_3 = U(s_2, a_2, o_3)) \dots \quad (1.4)$$

¹For a list of real-world applications of RL, see e.g., <https://bit.ly/42V7dIJ> from Csaba szepesvari (2024), <https://bit.ly/3EMMYCW> from Vitaly Kurin (2022), and <https://github.com/montrealrobotics/DeepRLInTheWorld>, which seems to be kept up to date.

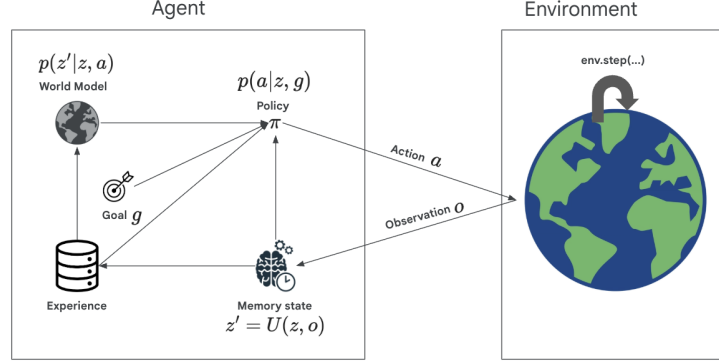


Figure 1.1: A small agent interacting with a big external world. The observation o_t (which, for notational simplicity, includes the previous action a_t) is used to update the internal agent state z_t , which is passed to the policy π which picks the next action a_{t+1} based on the agent’s goal g_t . Rewards are computed internally by the agent, by comparing z_t with its internal goal g_t . The observations, actions and rewards are stored in a replay buffer, which can be used to learn the policy, a value function (not shown), and optionally an internal world model (for use in model-based RL, see Chapter 4).

where p_{env} is the environment’s distribution over observations (which is usually unknown).

We define the optimal policy as

$$\pi^* = \arg \max_{\pi} \mathbb{E}_{p_0(s_0)} [V_{\pi}(s_0)] \quad (1.5)$$

Note that picking a policy to maximize the sum of expected rewards is an instance of the **maximum expected utility** principle. (In Section 7.1, we discuss the closely related concept of choosing a policy which minimizes the **regret**, which can be thought of as the difference between the expected reward of the agent’s policy compared to a reference policy.) There are various ways to design or learn such an optimal policy, depending on the assumptions we make about the environment, and the form of the agent. We will discuss some of these options below.

1.1.2 Episodic vs continual tasks

If the agent can potentially interact with the environment forever, we call it a **continual task** [Nai+21]. In this case, we replace the sum of rewards (when defining the value function) with the **average reward** [WNS21].

Alternatively, we say the agent is in an **episodic task** if its interaction terminates once the system enters a **terminal state** or **absorbing state**, which is a state which transitions to itself with 0 reward. After entering a terminal state, we may start a new **episode** from a new initial world state $z_0 \sim p_0$. (The agent will typically also reinitialize its own internal state s_0 .) The episode length is in general random. (For example, the length of an interaction with a chatbot may be quite variable, depending on the decisions taken by the chatbot agent and the randomness in the environment (i.e., the responses from the user)). Finally, if the trajectory length T in an episodic task is fixed and known, it is called a **finite horizon problem**.

We define the **return** for a state at time t to be the sum of expected rewards obtained going forwards, where each reward is multiplied by a **discount factor** $\gamma \in [0, 1]$:

$$G_t \triangleq r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^{T-t-1} r_{T-1} \quad (1.6)$$

$$= \sum_{k=0}^{T-t-1} \gamma^k r_{t+k} = \sum_{j=t}^{T-1} \gamma^{j-t} r_j \quad (1.7)$$

where $r_t = R(s_t, a_t)$ is the reward, and G_t is the **reward-to-go**. For episodic tasks that terminate at time T , we define $G_t = 0$ for $t \geq T$. Clearly, the return satisfies the following recursive relationship:

$$G_t = r_t + \gamma(r_{t+1} + \gamma r_{t+2} + \dots) = r_t + \gamma G_{t+1} \quad (1.8)$$

Furthermore, we define the value function to be the expected reward-to-go:

$$V_\pi(s_t) = \mathbb{E}[G_t | \pi] \quad (1.9)$$

The discount factor γ plays two roles. First, it ensures the return is finite even if $T = \infty$ (i.e., infinite horizon), provided we use $\gamma < 1$ and the rewards r_t are bounded. Second, it puts more weight on short-term rewards, which generally has the effect of encouraging the agent to achieve its goals more quickly. (For example, if $\gamma = 0.99$, then an agent that reaches a terminal reward of 1.0 in 15 steps will receive an expected discounted reward of $0.99^{15} = 0.86$, whereas if it takes 17 steps it will only get $0.99^{17} = 0.84$.) However, if γ is too small, the agent will become too greedy. In the extreme case where $\gamma = 0$, the agent is completely **myopic**, and only tries to maximize its immediate reward. In general, the discount factor reflects the assumption that there is a probability of $1 - \gamma$ that the interaction will end at the next step. (If $\gamma = 1 - \frac{1}{T}$, the agent expects to live on the order of T steps; for example, if each step is 0.1 seconds, then $\gamma = 0.95$ corresponds to 2 seconds.) For finite horizon problems, where T is known, we can set $\gamma = 1$, since we know the life time of the agent a priori.

1.1.3 Universal model

A generic representation for sequential decision making problems is shown in Figure 1.2. This is an extended version of the “universal modeling framework” proposed in [Pow19; Pow22], and is related to the “common model of the intelligent decision maker” discussed in [Sut22]. This common model assumes the environment can be modeled by a **controlled Markov process**² with hidden state w_t , which gets updated at each step in response to the agent’s action a_t . To allow for non-deterministic dynamics, we write this as $w_{t+1} = M(w_t, a_t, \epsilon_t^w)$, where M is the environment’s state transition function (which is usually not known to the agent) and ϵ_t^w is random system noise.³ The agent does not see the world state w_t , but instead sees a potentially noisy and/or partial observation $o_{t+1} = O(w_{t+1}, \epsilon_{t+1}^o)$ at each step, where ϵ_{t+1}^o is random observation noise. For example, when navigating a maze, the agent may only see what is in front of it, rather than seeing everything in the world all at once; furthermore, even the current view may be corrupted by sensor noise. Any given image, such as one containing a door, could correspond to many different locations in the world (this is called **perceptual aliasing**), each of which may require a different action.

Thus the agent needs use these observations to main an internal **belief state** about the world, denoted by z . This gets updated using the state update function

$$z_{t+1} = SU(z_t, a_t, o_{t+1}) \quad (1.10)$$

In the simplest setting, the internal z_t can just store all the past observations, $\mathbf{h}_t = (\mathbf{o}_{1:t}, \mathbf{a}_{1:t-1})$, but such non-parametric models can take a lot of time and space to work with, so we will usually consider parametric approximations. The agent can then pass its state to its policy to pick actions, using $a_{t+1} = \pi_t(z_{t+1})$.

We can further elaborate the behavior of the agent by breaking the state-update function into two parts. First the agent predicts its own next state, $z_{t+1|t} = P(z_t, a_t)$, using a **prediction function** P , and then it updates this prediction given the observation using **update function** U , to give $z_{t+1} = U(z_{t+1|t}, o_{t+1})$. Thus the SU function is defined as the composition of the predict and update functions

$$z_{t+1} = U(P(z_t, a_t), o_{t+1}) \quad (1.11)$$

²The Markovian assumption is without loss of generality, since we can always condition on the entire past sequence of states by suitably expanding the Markovian state space.

³Representing a stochastic function as a deterministic function with some noisy inputs is known as a functional causal model, or structural equation model. This is standard practice in the control theory and causality communities.

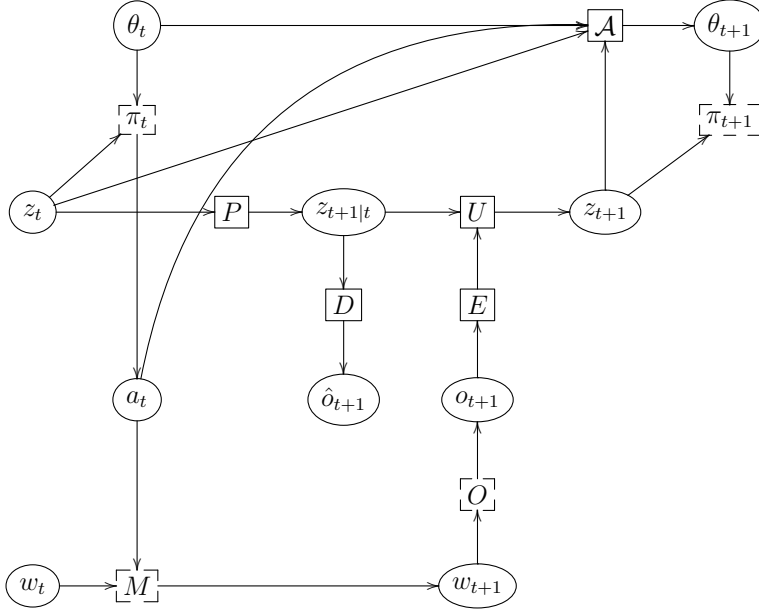


Figure 1.2: Detailed illustration of the interaction of an agent in an environment. The agent has internal state z_t , and chooses action a_t based on its policy π_t using $a_t \sim \pi_t(z_t|\theta_t)$. It then predicts its next internal states, $z_{t+1|t}$, via the predict function P , and optionally predicts the resulting observation, \hat{o}_{t+1} , via the observation decoder D . The environment has (hidden) internal state w_t , which gets updated by the environment model M to give the new state $w_{t+1} \sim M(w_t, a_t)$ in response to the agent's action. The environment also emits an observation o_{t+1} via its observation model, $o_{t+1} \sim O(w_{t+1})$. This gets encoded to $e_{t+1} = E(o_{t+1})$ by the agent's observation encoder E , which the agent uses to update its internal state using $z_{t+1} = U(z_t, a_t, e_{t+1})$. The policy is parameterized by θ_t , and these parameters may be updated (at a slower time scale) by an RL algorithm denoted by A . Square nodes are functions, circles are variables (either random or deterministic), dashed square nodes are stochastic functions that take an extra source of randomness (not shown).

If the observations are high dimensional (e.g., images), the agent may choose to encode its observations into a low-dimensional embedding e_{t+1} using an encoder, $e_{t+1} = E(o_{t+1})$; this can encourage the agent to focus on the relevant parts of the sensory signal. In this case, the state update becomes

$$z_{t+1} = U(P(z_t, a_t), E(o_{t+1})) \quad (1.12)$$

Optionally the agent can also learn to invert this encoder by training a decoder to predict the next observation using $\hat{o}_{t+1} = D(z_{t+1}|t)$; this can be a useful training signal, as we will discuss in Chapter 4. Finally, the agent needs to learn the action policy $\pi_t(z_t) = \pi(z_t; \theta_t)$. We can update the policy parameters using a learning algorithm, denoted

$$\theta_t = \mathcal{A}(o_{1:t}, a_{1:t}, r_{1:t}) = \mathcal{A}(\theta_{t-1}, a_t, z_t, r_t) \quad (1.13)$$

See Figure 1.2 for an illustration.

We see that, in general, there are three interacting stochastic processes we need to deal with: the environment’s states w_t (which are usually affected by the agents actions); the agent’s internal states z_t (which reflect its beliefs about the environment based on the observed data); and the agent’s policy parameters θ_t (which are updated based on the information stored in the belief state and the external observations).

1.1.4 Further reading

In later chapters, we will describe methods for learning the best policy to maximize $V_\pi(s_0) = \mathbb{E}[G_0|s_0, \pi]$. More details on RL can be found in textbooks such as [SB18; KWW22; Pla22; Li23; Sze10], and reviews such as [Aru+17; FL+18; Li18; Wen18a; ID19; JG24]. For a more theoretical treatment, see e.g., [Aga+22a; MMT24; FR23]. For details on how RL relates to **control theory**, see e.g., [Son98; Rec19; Ber19; Mey22]; for connections to operations research, see [Pow22]; for connections to finance, see [RJ22].

1.2 Canonical models

In this section, we describe different forms of model for the environment and the agent that have been studied in the literature.

1.2.1 Partially observed MDPs

The model shown in Figure 1.2 is called a **partially observable Markov decision process** or **POMDP** (pronounced “pom-dee-pee”) [KLC98; LHP22; Sub+22]. Typically the environment’s dynamics model is represented by a stochastic transition function, rather than a deterministic function with noise as an input. We can derive this transition function as follows:

$$p(w_{t+1}|w_t, a_t) = \mathbb{E}_{\epsilon_t^w} [\mathbb{I}(w_{t+1} = W(w_t, a_t, \epsilon_t^w))] \quad (1.14)$$

Similarly the stochastic observation function is given by

$$p(o_{t+1}|w_{t+1}) = \mathbb{E}_{\epsilon_{t+1}^o} [\mathbb{I}(o_{t+1} = O(w_{t+1}, \epsilon_{t+1}^o))] \quad (1.15)$$

Note that we can combine these two distributions to derive the joint world model $p_{WO}(w_{t+1}, o_{t+1}|w_t, a_t)$. Also, we can use these distributions to derive the environment’s non-Markovian observation distribution, $p_{\text{env}}(o_{t+1}|o_{1:t}, a_{1:t})$, used in Equation (1.4), as follows:

$$p_{\text{env}}(o_{t+1}|o_{1:t}, a_{1:t}) = \sum_{w_{t+1}} p(o_{t+1}|w_{t+1}) p(w_{t+1}|a_{1:t}) \quad (1.16)$$

$$p(w_{t+1}|a_{1:t}) = \sum_{w_1} \cdots \sum_{w_t} p(w_1|a_1) p(w_2|w_1, a_1) \cdots p(w_{t+1}|w_t, a_t) \quad (1.17)$$

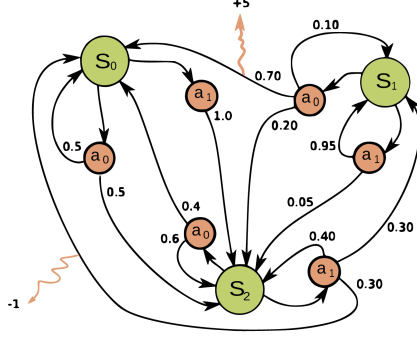


Figure 1.3: Illustration of an MDP as a finite state machine (FSM). The MDP has three discrete states (green circles), two discrete actions (orange circles), and two non-zero rewards (orange arrows). The numbers on the black edges represent state transition probabilities, e.g., $p(s' = s_0 | a = a_0, s' = s_1) = 0.7$; most state transitions are impossible (probability 0), so the graph is sparse. The numbers on the yellow wiggly edges represent expected rewards, e.g., $R(s = s_1, a = a_0, s' = s_0) = +5$; state transitions with zero reward are not annotated. From https://en.wikipedia.org/wiki/Markov_decision_process. Used with kind permission of Wikipedia author waldoalvarez.

If the world model (both $p(o|w)$ and $p(w'|w, a)$) is known, then we can — in principle — solve for the optimal policy. The method requires that the agent’s internal state correspond to the **belief state** $s_t = \mathbf{b}_t = p(w_t | \mathbf{h}_t)$, where $\mathbf{h}_t = (o_{1:t}, a_{1:t-1})$ is the observation history. The belief state can be updated recursively using Bayes rule. See Section 1.2.6 for details. The belief state forms a sufficient statistic for the optimal policy. Unfortunately, computing the belief state and the resulting optimal policy is wildly intractable [PT87; KLC98]. We discuss some approximate methods in Section 1.3.4.

1.2.2 Markov decision process (MDPs)

A **Markov decision process** [Put94] is a special case of a POMDP in which the environment states are observed, so $w_t = o_t = s_t$. We usually define an MDP in terms of the state transition matrix induced by the world model:

$$p_S(s_{t+1} | s_t, a_t) = \mathbb{E}_{\epsilon_t^s} [\mathbb{I}(s_{t+1} = W(s_t, a_t, \epsilon_t^s))] \quad (1.18)$$

In lieu of an observation model, we assume the environment (as opposed to the agent) sends out a reward signal, sampled from $p_R(r_t | s_t, a_t, s_{t+1})$. The expected reward is then given by

$$R(s_t, a_t, s_{t+1}) = \sum_r r \, p_R(r | s_t, a_t, s_{t+1}) \quad (1.19)$$

$$R(s_t, a_t) = \sum_{s_{t+1}} p_S(s_{t+1} | s_t, a_t) R(s_t, a_t, s_{t+1}) \quad (1.20)$$

Note that the field of control theory uses slightly different terminology and notation when describing the same setup: the environment is called the **plant**, the agent is called the **controller**, States are denoted by $\mathbf{x}_t \in \mathcal{X} \subseteq \mathbb{R}^D$, actions are denoted by $\mathbf{u}_t \in \mathcal{U} \subseteq \mathbb{R}^K$, and rewards are replaced by costs $c_t \in \mathbb{R}$.

Given a stochastic policy $\pi(a_t | s_t)$, the agent can interact with the environment over many steps. Each step is called a **transition**, and consists of the tuple (s_t, a_t, r_t, s_{t+1}) , where $a_t \sim \pi(\cdot | s_t)$, $s_{t+1} \sim p_S(s_{t+1} | s_t, a_t)$, and $r_t \sim p_R(r_t | s_t, a_t, s_{t+1})$. Hence, under policy π , the probability of generating a **trajectory** length T , $\boldsymbol{\tau} = (s_0, a_0, r_0, s_1, a_1, r_1, s_2, \dots, s_T)$, can be written explicitly as

$$p(\boldsymbol{\tau}) = p_0(s_0) \prod_{t=0}^{T-1} \pi(a_t | s_t) p_S(s_{t+1} | s_t, a_t) p_R(r_t | s_t, a_t, s_{t+1}) \quad (1.21)$$

In general, the state and action sets of an MDP can be discrete or continuous. When both sets are finite, we can represent these functions as lookup tables; this is known as a **tabular representation**. In this case, we can represent the MDP as a **finite state machine**, which is a graph where nodes correspond to states, and edges correspond to actions and the resulting rewards and next states. Figure 1.3 gives a simple example of an MDP with 3 states and 2 actions.

If we know the world model p_S and p_R , and if the state and action space is tabular, then we can solve for the optimal policy using dynamic programming techniques, as we discuss in Section 2.2. However, typically the world model is unknown, and the states and actions may need complex nonlinear models to represent their transitions. In such cases, we will have to use RL methods to learn a good policy.

1.2.3 Goal-conditioned MDPs

A **goal-conditioned MDP** is one in which the reward is defined as $R(s, a|g) = 1$ iff the goal state is achieved, i.e., $R(s, a|s) = \mathbb{I}(s = g)$. We can also define a dense reward signal using some state abstraction function ϕ , by defining $R(s, a|g) = \text{sim}(s, g)$, where sim is some kind of similarity metric. For example, if s is an image and g is a sentence, we may use cosine similarity

$$\text{sim}(s, g) = \frac{\phi(s)^\top \psi(g)}{\|\phi(s)\| \|\psi(g)\|} \quad (1.22)$$

where $\phi(s)$ is an embedding of the image (state). and $\psi(g)$ is an embedding of the text (goal). Such embeddings can be computed by using a VLM or vision-language model (see Section 6.3.2).

A goal-conditioned policy of the form $\pi(a|s, g)$ is sometimes called a **universal policy** [Sch+15a]. We can learn such policies using **goal-conditioned RL** methods (see e.g., [LZZ22] and Section 2.5.5).

Note that multi-goal RL is different to multi-task RL. The latter refers to the ability to solve different “tasks”, which correspond to entire MDPs (with different dynamics as well as different rewards).

1.2.4 Contextual MDPs

A **Contextual MDP** [HDCM15] is an MDP where the dynamics and rewards of the environment depend on a hidden static parameter referred to as the context. (This is different to a contextual bandit, discussed in Section 1.2.5, where the context is observed at each step.) A simple example of a contextual MDP is a video game, where each level of the game is **procedurally generated**, that is, it is randomly generated each time the agent starts a new episode. Thus the agent must solve a sequence of related MDPs, which are drawn from a common distribution. This requires the agent to **generalize** across multiple MDPs, rather than overfitting to a specific environment [Cob+19; Kir+21; Tom+22]. (This form of generalization is different from generalization within an MDP, which requires generalizing across states, rather than across environments; both are important.)

A contextual MDP is a special kind of POMDP where the hidden variable corresponds to the unknown parameters of the model. In [Gho+21], they call this an **epistemic POMDP**, which is closely related to the concept of belief state MDP which we discuss in Section 1.2.6.

1.2.5 Contextual bandits

A **contextual bandit** is a special case of a POMDP where the world state transition function is independent of the action of the agent and the previous state, i.e., $p(w_t|w_{t-1}, a_t) = p(w_t)$. In this case, we call the world states “contexts”; these are observable by the agent, i.e., $o_t = w_t$. Since the world state distribution is independent of the agents actions, the agent has no effect on the external environment. However, its actions do affect the rewards that it receives. Thus the agent’s internal belief state — about the underlying reward function $R(o, a)$ — does change over time, as the agent learns a model of the world (see Section 1.2.6).

A special case of a contextual bandit is a regular bandit, in which there is no context, or equivalently, s_t is some fixed constant that never changes. When there are a finite number of possible actions, $\mathcal{A} = \{a_1, \dots, a_K\}$,

this is called a **multi-armed bandit**.⁴ In this case the reward model has the form $R(a) = f(\mathbf{w}_a)$, where \mathbf{w}_a are the parameters for arm a .

Contextual bandits have many applications. For example, consider an **online advertising system**. In this case, the state s_t represents features of the web page that the user is currently looking at, and the action a_t represents the identity of the ad which the system chooses to show. Since the relevance of the ad depends on the page, the reward function has the form $R(s_t, a_t)$, and hence the problem is contextual. The goal is to maximize the expected reward, which is equivalent to the expected number of times people click on ads; this is known as the **click through rate** or **CTR**. (See e.g., [Gra+10; Li+10; McM+13; Aga+14; Du+21; YZ22] for more information about this application.) Another application of contextual bandits arises in **clinical trials** [VBW15]. In this case, the state s_t are features of the current patient we are treating, and the action a_t is the treatment the doctor chooses to give them (e.g., a new drug or a **placebo**).

For more details on bandits, see e.g., [LS19; Sli19].

1.2.6 Belief state MDPs

In this section, we describe a kind of MDP where the state represents a probability distribution, known as a **belief state** or **information state**, which is updated by the agent (“in its head”) as it receives information from the environment.⁵ More precisely, consider a contextual bandit problem, where the agent approximates the unknown reward by a function $R(o, a) = f(o, a; \mathbf{w})$. Let us denote the posterior over the unknown parameters by $\mathbf{b}_t = p(\mathbf{w}|\mathbf{h}_t)$, where $\mathbf{h}_t = \{o_{1:t}, a_{1:t}, r_{1:t}\}$ is the history of past observations, actions and rewards. This belief state can be updated deterministically using Bayes’ rule; we denote this operation by $\mathbf{b}_{t+1} = \text{BayesRule}(\mathbf{b}_t, o_{t+1}, a_{t+1}, r_{t+1})$. (This corresponds to the state update SU defined earlier.) Using this, we can define the following **belief state MDP**, with deterministic dynamics given by

$$p(\mathbf{b}_{t+1}|\mathbf{b}_t, o_{t+1}, a_{t+1}, r_{t+1}) = \mathbb{I}(\mathbf{b}_{t+1} = \text{BayesRule}(\mathbf{b}_t, o_{t+1}, a_{t+1}, r_{t+1})) \quad (1.23)$$

and reward function given by

$$p(r_t|o_t, a_t, \mathbf{b}_t) = \int p_R(r_t|o_t, a_t; \mathbf{w})p(\mathbf{w}|\mathbf{b}_t)d\mathbf{w} \quad (1.24)$$

If we can solve this (PO)MDP, we have the optimal solution to the exploration-exploitation problem (see Section 1.3.5).

As a simple example, consider a context-free **Bernoulli bandit**, where $p_R(r|a) = \text{Ber}(r|\mu_a)$, and $\mu_a = p_R(r = 1|a) = R(a)$ is the expected reward for taking action a . The only unknown parameters are $\mathbf{w} = \mu_{1:A}$. Suppose we use a factored beta prior

$$p_0(\mathbf{w}) = \prod_a \text{Beta}(\mu_a|\alpha_0^a, \beta_0^a) \quad (1.25)$$

where $\mathbf{w} = (\mu_1, \dots, \mu_K)$. We can compute the posterior in closed form to get

$$p(\mathbf{w}|\mathcal{D}_t) = \prod_a \text{Beta}(\mu_a|\underbrace{\alpha_0^a + N_t^0(a)}_{\alpha_t^a}, \underbrace{\beta_0^a + N_t^1(a)}_{\beta_t^a}) \quad (1.26)$$

where

$$N_t^r(a) = \sum_{i=1}^{t-1} \mathbb{I}(a_i = a, r_i = r) \quad (1.27)$$

⁴The terminology arises by analogy to a slot machine (sometimes called a “bandit”, because it steals your money) in a casino. If there are K slot machines, each with different rewards (payout rates), then the agent (player) must explore the different machines (by pulling the arms) until they have discovered which one is best, and can then stick to exploiting it.

⁵Technically speaking, this is a POMDP, where we assume the states are observed, and the parameters are the unknown hidden random variables. This is in contrast to Section 1.2.1, where the states were not observed, and the parameters were assumed to be known.

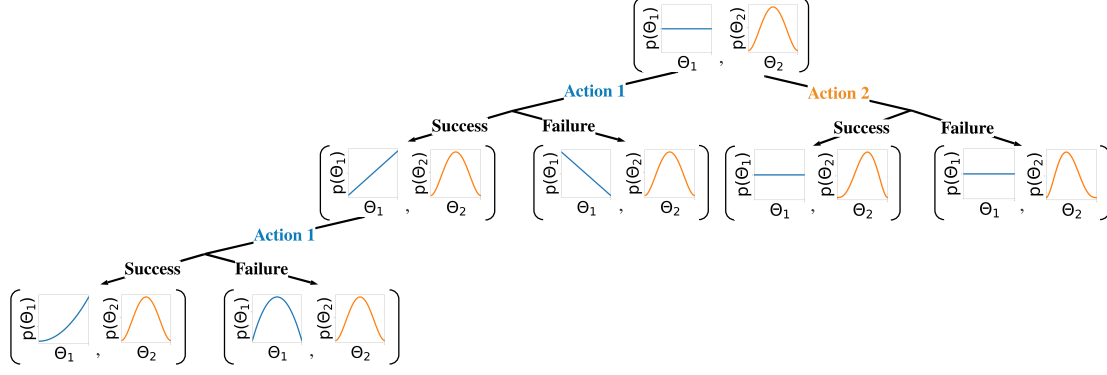


Figure 1.4: Illustration of sequential belief updating for a two-armed beta-Bernoulli bandit. The prior for the reward for action 1 is the (blue) uniform distribution $\text{Beta}(1, 1)$; the prior for the reward for action 2 is the (orange) unimodal distribution $\text{Beta}(2, 2)$. We update the parameters of the belief state based on the chosen action, and based on whether the observed reward is success (1) or failure (0).

This is illustrated in Figure 1.4 for a two-armed Bernoulli bandit. We can use a similar method for a **Gaussian bandit**, where $p_R(r|a) = \mathcal{N}(r|\mu_a, \sigma_a^2)$.

In the case of contextual bandits, the problem is conceptually the same, but becomes more complicated computationally. If we assume a **linear regression bandit**, $p_R(r|s, a; \mathbf{w}) = \mathcal{N}(r|\phi(s, a)^\top \mathbf{w}, \sigma^2)$, we can use Bayesian linear regression to compute $p(\mathbf{w}|\mathcal{D}_t)$ exactly in closed form. If we assume a **logistic regression bandit**, $p_R(r|s, a; \mathbf{w}) = \text{Ber}(r|\sigma(\phi(s, a)^\top \mathbf{w}))$, we have to use approximate methods for approximate Bayesian logistic regression to compute $p(\mathbf{w}|\mathcal{D}_t)$. If we have a **neural bandit** of the form $p_R(r|s, a; \mathbf{w}) = \mathcal{N}(r|f(s, a; \mathbf{w}))$ for some nonlinear function f , then posterior inference is even more challenging (this is equivalent to the problem of inference in Bayesian neural networks, see e.g., [Arb+23] for a review paper for the offline case, and [DMKM22; JCM24] for some recent online methods).

We can generalize the above methods to compute the belief state for the parameters of an MDP in the obvious way, but modeling both the reward function and state transition function.

Once we have computed the belief state, we can derive a policy with optimal regret using the methods like UCB (Section 7.2.3) or Thompson sampling (Section 7.2.2).

1.2.7 Optimization problems as decision problems

The bandit problem is an example of a problem where the agent must interact with the world in order to collect information, but it does not otherwise affect the environment. Thus the agent's internal belief state changes over time, but the environment state does not.⁶ Such problems commonly arise when we are trying to optimize a fixed but unknown function R . We can “query” the function by evaluating it at different points (parameter values), and in some cases, the resulting observation may also include gradient information. The agent’s goal is to find the optimum of the function in as few steps as possible.⁷ We give some examples of this problem setting below.

⁶In the contextual bandit problem, the environment state (context) does change, but not in response to the agent’s actions. Thus $p(o_t)$ is usually assumed to be a static distribution.

⁷If we only care about the final performance of the agent, we can try to minimize the **simple regret**, which is just the regret at the last step, namely l_T . This is the difference between the function value we chose and the true optimum. Minimizing simple regret results in a problem known as **pure exploration** [BMS11], where the agent needs to interact with the environment to learn the underlying MDP; at the end, it can then solve for the resulting policy using planning methods (see Section 2.2). However, in general RL problems, it is more common to focus on the **cumulative regret**, also called the **total regret** or just the **regret**, which is defined as $L_T \triangleq \mathbb{E} \left[\sum_{t=1}^T l_t \right]$.

1.2.7.1 Best-arm identification

In the standard multi-armed bandit problem our goal is to maximize the sum of expected rewards. However, in some cases, the goal is to determine the best arm given a fixed budget of T trials; this variant is known as **best-arm identification** [ABM10]. Formally, this corresponds to optimizing the **final reward** criterion:

$$V_{\pi, \pi_T} = \mathbb{E}_{p(a_{1:T}, r_{1:T} | s_0, \pi)} [R(\hat{a})] \quad (1.28)$$

where $\hat{a} = \pi_T(a_{1:T}, r_{1:T})$ is the estimated optimal arm as computed by the **terminal policy** π_T applied to the sequence of observations obtained by the exploration policy π . This can be solved by a simple adaptation of the methods used for standard bandits.

1.2.7.2 Bayesian optimization

Bayesian optimization is a gradient-free approach to optimizing expensive blackbox functions. That is, we want to find

$$\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmax}} R(\mathbf{w}) \quad (1.29)$$

for some unknown function R , where $\mathbf{w} \in \mathbb{R}^N$, using as few actions (function evaluations of R) as possible. This is essentially an “infinite arm” version of the best-arm identification problem [Tou14], where we replace the discrete choice of arms $a \in \{1, \dots, K\}$ with the parameter vector $\mathbf{w} \in \mathbb{R}^N$. In this case, the optimal policy can be computed if the agent’s state s_t is a belief state over the unknown function, i.e., $s_t = p(R | \mathbf{h}_t)$. A common way to represent this distribution is to use Gaussian processes. We can then use heuristics like expected improvement, knowledge gradient or Thompson sampling to implement the corresponding policy, $\mathbf{w}_t = \pi(s_t)$. For details, see e.g., [Gar23].

1.2.7.3 Active learning

Active learning is similar to BayesOpt, but instead of trying to find the point at which the function is largest (i.e., \mathbf{w}^*), we are trying to learn the whole function R , again by querying it at different points \mathbf{w}_t . Once again, the optimal strategy again requires maintaining a belief state over the unknown function, but now the best policy takes a different form, such as choosing query points to reduce the entropy of the belief state. See e.g., [Smi+23].

1.2.7.4 Stochastic Gradient Descent (SGD)

Finally we discuss how to interpret SGD as a sequential decision making process, following [Pow22]. The action space consists of querying the unknown function R at locations $\mathbf{a}_t = \mathbf{w}_t$, and observing the function value $r_t = R(\mathbf{w}_t)$; however, unlike BayesOpt, now we also observe the corresponding gradient $\mathbf{g}_t = \nabla_{\mathbf{w}} R(\mathbf{w})|_{\mathbf{w}_t}$, which gives non-local information about the function. The environment state contains the true function R which is used to generate the observations given the agent’s actions. The agent state contains the current parameter estimate \mathbf{w}_t , and may contain other information such as first and second moments \mathbf{m}_t and \mathbf{v}_t , needed by methods such as Adam. The update rule (for vanilla SGD) takes the form $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{g}_t$, where the stepsize α_t is chosen by the policy, $\alpha_t = \pi(s_t)$. The terminal policy has the form $\pi(s_T) = \mathbf{w}_T$.

Although in principle it is possible to learn the learning rate (stepsize) policy using RL (see e.g., [Xu+17]), the policy is usually chosen by hand, either using a **learning rate schedule** or some kind of manually designed **adaptive learning rate** policy (e.g., based on second order curvature information).

1.3 Reinforcement Learning: a high-level summary

In this section, we give a brief overview of how to compute optimal policies when the model of the environment is unknown; this is the core problem tackled by RL. We mostly focus on the MDP case, but discuss the POMDP case in Section 1.3.4.

We can categorize RL methods along multiple dimensions, such as the following:

Approach	Method	Functions learned	On/Off	Section
Value-based	SARSA	$Q(s, a)$	On	Section 2.4
Value-based	Q -learning	$Q(s, a)$	Off	Section 2.5
Policy-based	REINFORCE	$\pi(a s)$	On	Section 3.1.3
Policy-based	A2C	$\pi(a s), V(s)$	On	Section 3.2.1
Policy-based	TRPO/PPO	$\pi(a s), \text{Adv}(s, a)$	On	Section 3.3.3
Policy-based	DDPG	$a = \pi(s), Q(s, a)$	Off	Section 3.2.6.2
Policy-based	Soft actor-critic	$\pi(a s), Q(s, a)$	Off	Section 3.6.8
Model-based	MBRL	$p(s' s, a)$	Off	Chapter 4

Table 1.1: Summary of some popular methods for RL. On/off refers to on-policy vs off-policy methods.

- What does the agent learn? Options include the value function, the policy, the model, or some combination of the above.
- How does the agent represent its unknown functions? The two main choices are to use non-parametric or **tabular representations**, or to use parametric representations based on function approximation. If these functions are based on neural networks, this approach is called “**deep RL**”, where the term “deep” refers to the use of neural networks with many layers.
- How are the actions are selected? Options include **on-policy** methods, where actions must be selected by the agent’s current policy), and **off-policy** methods, where actions can be select by any kind of policy, including human demonstrations.

Table 1.1 lists a few common examples of RL methods, classified along these lines. More details are given in the subsequent sections.

1.3.1 Value-based RL

In this section, we give a brief introduction to **value-based RL**, also called **Approximate Dynamic Programming** or **ADP**; see Chapter 2 for more details.

We introduced the value function $V_\pi(s)$ in Equation (1.1), which we repeat here for convenience:

$$V_\pi(s) \triangleq \mathbb{E}_\pi [G_0 | s_0 = s] = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s \right] \quad (1.30)$$

The value function for the optimal policy π^* is known to satisfy the following recursive condition, known as **Bellman’s equation**:

$$V^*(s) = \max_a R(s, a) + \gamma \mathbb{E}_{p_S(s'|s, a)} [V^*(s')] \quad (1.31)$$

This follows from the principle of **dynamic programming**, which computes the optimal solution to a problem (here the value of state s) by combining the optimal solution of various subproblems (here the values of the next states s'). This can be used to derive the following learning rule:

$$V(s) \leftarrow V(s) + \eta[r + \gamma V(s') - V(s)] \quad (1.32)$$

where $s' \sim p_S(\cdot | s, a)$ is the next state sampled from the environment, and $r = R(s, a)$ is the observed reward. This is called **Temporal Difference** or **TD learning** (see Section 2.3.2 for details). Unfortunately, it is not clear how to derive a policy if all we know is the value function. We now describe a solution to this problem.

We first generalize the notion of value function to assigning a value to a state and action pair, by defining the **Q function** as follows:

$$Q_\pi(s, a) \triangleq \mathbb{E}_\pi [G_0 | s_0 = s, a_0 = a] = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s, a_0 = a \right] \quad (1.33)$$

This quantity represents the expected return obtained if we start by taking action a in state s , and then follow π to choose actions thereafter. The Q function for the optimal policy satisfies a modified Bellman equation

$$Q^*(s, a) = R(s, a) + \gamma \mathbb{E}_{p_S(s'|s, a)} \left[\max_{a'} Q^*(s', a') \right] \quad (1.34)$$

This gives rise to the following TD update rule:

$$Q(s, a) \leftarrow r + \gamma \max_{a'} Q(s', a') - Q(s, a) \quad (1.35)$$

where we sample $s' \sim p_S(\cdot|s, a)$ from the environment. The action is chosen at each step from the implicit policy

$$a = \operatorname{argmax}_{a'} Q(s, a') \quad (1.36)$$

This is called **Q learning** (see Section 2.5 for details),

1.3.2 Policy-based RL

In this section we give a brief introduction to **Policy-based RL**; for details see Chapter 3.

In policy-based methods, we try to directly maximize $J(\pi_\theta) = \mathbb{E}_{p(s_0)} [V_\pi(s_0)]$ wrt the parameter's θ ; this is called **policy search**. If $J(\pi_\theta)$ is differentiable wrt θ , we can use stochastic gradient ascent to optimize θ , which is known as **policy gradient** (see Section 3.1).

Policy gradient methods have the advantage that they provably converge to a local optimum for many common policy classes, whereas Q -learning may diverge when approximation is used (Section 2.5.2.5). In addition, policy gradient methods can easily be applied to continuous action spaces, since they do not need to compute $\operatorname{argmax}_a Q(s, a)$. Unfortunately, the score function estimator for $\nabla_\theta J(\pi_\theta)$ can have a very high variance, so the resulting method can converge slowly.

One way to reduce the variance is to learn an approximate value function, $V_w(s)$, and to use it as a baseline in the score function estimator. We can learn $V_w(s)$ using TD learning. Alternatively, we can learn an advantage function, $A_w(s, a)$, and use it as a baseline. These policy gradient variants are called **actor critic** methods, where the actor refers to the policy π_θ and the critic refers to V_w or A_w . See Section 3.2 for details.

1.3.3 Model-based RL

In this section, we give a brief introduction to **model-based RL**; for more details, see Chapter 4.

Value-based methods, such as Q -learning, and policy search methods, such as policy gradient, can be very **sample inefficient**, which means they may need to interact with the environment many times before finding a good policy, which can be problematic when real-world interactions are expensive. In model-based RL, we first learn the MDP, including the $p_S(s'|s, a)$ and $R(s, a)$ functions, and then compute the policy, either using approximate dynamic programming on the learned model, or doing lookahead search. In practice, we often interleave the model learning and planning phases, so we can use the partially learned policy to decide what data to collect, to help learn a better model.

1.3.4 State uncertainty (partial observability)

In an MDP, we assume that the state of the environment s_t is the same as the observation o_t obtained by the agent. But in many problems, the observation only gives partial information about the underlying state of the world (e.g., a rodent or robot navigating in a maze). This is called **partial observability**. In this case, using a policy of the form $a_t = \pi(o_t)$ is suboptimal, since o_t does not give us complete state information. Instead we need to use a policy of the form $a_t = \pi(\mathbf{h}_t)$, where $\mathbf{h}_t = (a_1, o_1, \dots, a_{t-1}, o_t)$ is the entire past history of observations and actions, plus the current observation. Since depending on the entire past is not tractable for a long-lived agent, various approximate solution methods have been developed, as we summarize below.

1.3.4.1 Optimal solution

If we know the true latent structure of the world (i.e., both $p(o|z)$ and $p(z'|z, a)$, to use the notation of Section 1.1.3), then we can use solution methods designed for POMDPs, discussed in Section 1.2.1. This requires using Bayesian inference to compute a belief state, $\mathbf{b}_t = p(w_t|\mathbf{h}_t)$ (see Section 1.2.6), and then using this belief state to guide our decisions. However, learning the parameters of a POMDP (i.e., the generative latent world model) is very difficult, as is recursively computing and updating the belief state, as is computing the policy given the belief state. Indeed, optimally solving POMDPs is known to be computationally very difficult for any method [PT87; KLC98]. So in practice simpler approximations are used. We discuss some of these below. (For more details, see [Mur00].)

Note that it is possible to marginalize out the POMDP latent state w_t , to derive a prediction over the next observable state, $p(o_{t+1}|\mathbf{h}_t, \mathbf{a}_t)$. This can then become a learning target for a model, that is trained to directly predict future observations, without explicitly invoking the concept of latent state. This is called a **predictive state representation** or **PSR** [LS01]. This is related to the idea of **observable operator models** [Jae00], and to the concept of successor representations which we discuss in Section 4.5.2.

1.3.4.2 Finite observation history

The simplest solution to the partial observability problem is to define the state to be a finite history of the last k observations, $\mathbf{s}_t = \mathbf{h}_{t-k:t}$; when the observations \mathbf{o}_t are images, this is often called **frame stacking**. We can then use standard MDP methods. Unfortunately, this cannot capture long-range dependencies in the data.

1.3.4.3 Stateful (recurrent) policies

A more powerful approach is to use a stateful policy, that can remember the entire past, and not just respond to the current input or last k frames. For example, we can represent the policy by an RNN (recurrent neural network), as proposed in the **R2D2** paper [Kap+18], and used in many other papers. Now the hidden state w_t of the RNN will implicitly summarize the past observations, \mathbf{h}_t , and can be used in lieu of the state \mathbf{s}_t in any standard RL algorithm.

RNNs policies are widely used, and this method is often effective in solving partially observed problems. However, they typically will not plan to perform information-gathering actions, since there is no explicit notion of belief state or uncertainty. However, such behavior can arise via meta-learning [Mik+20].

1.3.5 Model uncertainty (exploration-exploitation tradeoff)

In RL problems, we typically assume the underlying transition and reward models are not known. We can either try to explicitly learn these models (as in model-based RL), and then solve for the policy, or just learn the policy directly (as in model-free RL). But in either case, we need to explore the environment in order to collect enough data to figure out what to do. This may involve choosing between actions that the agent knows will yield high reward, vs choosing actions which might not been known to yield high reward but which will be informative about potential future gains. This is called the **exploration-exploitation tradeoff**. In this section, we discuss some simple heuristic solutions to this problem. See Section 7.2 for more sophisticated methods.

If we just want to exploit our current knowledge (without trying to learn new things), we can use the **greedy policy**:

$$a_t = \operatorname{argmax}_a Q(s, a) \quad (1.37)$$

We can add exploration to this by sometimes picking some other, non-greedy action, as we discuss below.

One approach is to use an **ϵ -greedy** policy π_ϵ , parameterized by $\epsilon \in [0, 1]$. In this case, we pick the greedy action wrt the current model, $a_t = \operatorname{argmax}_a \hat{R}_t(s_t, a)$ with probability $1 - \epsilon$, and a random action with probability ϵ . This rule ensures the agent’s continual exploration of all state-action combinations.

$\hat{R}(s, a_1)$	$\hat{R}(s, a_2)$	$\pi_\epsilon(a s_1)$	$\pi_\epsilon(a s_2)$	$\pi_\tau(a s_1)$	$\pi_\tau(a s_2)$
1.00	9.00	0.05	0.95	0.00	1.00
4.00	6.00	0.05	0.95	0.12	0.88
4.90	5.10	0.05	0.95	0.45	0.55
5.05	4.95	0.95	0.05	0.53	0.48
7.00	3.00	0.95	0.05	0.98	0.02
8.00	2.00	0.95	0.05	1.00	0.00

Table 1.2: Comparison of ϵ -greedy policy (with $\epsilon = 0.1$) and Boltzmann policy (with $\tau = 1$) for a simple MDP with 6 states and 2 actions. Adapted from Table 4.1 of [GK19].

Unfortunately, this heuristic can be shown to be suboptimal, since it explores every action with at least a constant probability $\epsilon/|\mathcal{A}|$, although this can be solved by annealing ϵ to 0 over time.

Another problem with ϵ -greedy is that it can result in “dithering”, in which the agent continually changes its mind about what to do. In [DOB21] they propose a simple solution to this problem, known as ϵz -greedy, that often works well. The idea is that with probability $1 - \epsilon$ the agent exploits, but with probability ϵ the agent explores by repeating the sampled action for $n \sim z()$ steps in a row, where $z(n)$ is a distribution over the repeat duration. This can help the agent escape from local minima. (See also [Tre+23], who learn a policy to not only pick an action, but also how long to use that action for, by solving an augmented MDP where the action space is augmented by duration.)

Another simple approach to exploration is to use **Boltzmann exploration**, which assigns higher probabilities to explore more promising actions, taking into account the reward function. That is, we use a policy of the form

$$\pi_\tau(a|s) = \frac{\exp(\hat{R}_t(s, a)/\tau)}{\sum_{a'} \exp(\hat{R}_t(s, a')/\tau)} \quad (1.38)$$

where $\tau > 0$ is a temperature parameter that controls how entropic the distribution is. As τ gets close to 0, π_τ becomes close to a greedy policy. On the other hand, higher values of τ will make $\pi(a|s)$ more uniform, and encourage more exploration. Its action selection probabilities can be much “smoother” with respect to changes in the reward estimates than ϵ -greedy, as illustrated in Table 1.2.

The Boltzmann policy explores equally widely in all states. An alternative approach is to try to explore (state,action) combinations where the consequences of the outcome might be uncertain. This can be achieved using an **exploration bonus** $R_t^b(s, a)$, which is large if the number of times we have tried action a in state s is small. We can then add R_t^b to the regular reward, to bias the behavior in a way that will hopefully cause the agent to learn useful information about the world. This is called an **intrinsic reward** function (Section 7.4).

1.3.6 Reward functions

Sequential decision making relies on the user to define the reward function in order to encourage the agent to exhibit some desired behavior. In this section, we discuss this crucial aspect of the problem.

1.3.6.1 The reward hypothesis

The “**reward hypothesis**” states that “all of what we mean by goals and purposes can be well thought of as maximization of the expected value of the cumulative sum of a received scalar signal (reward)” [Sut04]. (See also the closely related “reward is enough” hypothesis [Sil+21].) Whether this hypothesis is true or not depends on what one means by “goals and purposes”. This can be formalized in terms of preference relations over (state,action) trajectories, as discussed in [Bow+23]. (See also [Boo+23; BKM24] for some related work on reward function design.)

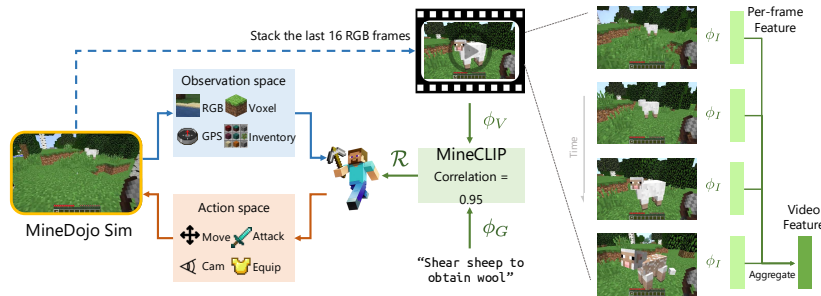


Figure 1.5: Illustration of how the MineCLIP reward function can be used to help train an agent to play Minecraft in the MineDojo simulator. From Figure 4 of [Fan+22]. Used with kind permission of Jim Fan.

1.3.6.2 Non-Markovian rewards

Most of the literature assumes the reward can be defined in terms of the current state and action, $R(s, a)$, or in terms of the most recent state transition, $R(s, a, s')$. In [Bow+23], they discuss when a utility function over trajectories can be converted into a Markovian reward of the form $R(s, a, s')$.

In general, the reward function will need to be non-Markovian. For example, consider training an agent to solve various goals, specified in natural language, inside the Minecraft video game. (For a general discussion of goal-conditioned RL, see Section 1.2.3.) In this case, we do not have access to the underlying world state, and even if we did, it can be hard to determine from a single state, or single state transition pair, whether a generic goal (such as “shear the sheep to obtain wool”) has been satisfied. In the **MineDojo** paper [Fan+22], they tackled this problem by pre-training a reward model of the form $R(o(t - K : t), g)$, where $o(t - K : t)$ are the last K frames, and g is the goal. This model, known as **MineCLIP**, was trained using contrastive learning applied to a large corpus of video-text pairs.⁸

1.3.6.3 Reward hacking

In some cases, the reward function may be misspecified, so even though the agent may maximize the reward, this might turn out not to be what the user desired. For example, suppose the user rewards the agent for making as many paper clips as possible. An optimal agent may convert the whole world into a paper clip factory, because the user forgot to specify various constraints, such as not killing people (which might otherwise be necessary in order to use as many resources as possible for paperclips). In the **AI alignment** community, this example is known as the **paperclip maximizer problem**, and is due to Nick Bostrom [Bos16]. (See e.g., <https://openai.com/index/faulty-reward-functions/> for some examples that have occurred in practice.) This is an example of a more general problem known as **reward hacking** [Ska+22]. For a potential solution, based on the assistance game paradigm, see Section 6.2.6.

1.3.6.4 Sparse reward

Even if the reward function is correct, optimizing it is not always easy. In particular, many problems suffer from **sparse reward**, in which $R(s, a) = 0$ for almost all states and actions, so the agent only every gets feedback (either positive or negative) on the rare occasions when it achieves some unknown goal. This requires **deep exploration** [Osb+19] to find the rewarding states. One approach to this is use to use PSRL (Section 7.2.2.2). However, various other heuristics have been developed, some of which we discuss below.

⁸To make this reward function fast to compute, they computed it using a simple comparison between the embedding of the goal, $\phi_G(g)$, and the aggregated embeddings of each image, $1/K \sum_{k=0}^{K-1} \phi_I(o_{t-k})$. By caching the embeddings of previously seen frames, and using a frozen image encoder which is shared between the reward and the agent, computation could be significantly sped up.

1.3.6.5 Reward shaping

In **reward shaping**, we add prior knowledge about what we believe good states should look like, as a way to combat the difficulties of learning from sparse reward. That is, we define a new reward function $r' = r + F$, where F is called the shaping function. In general, this can affect the optimal policy. For example, if a soccer playing agent is “artificially” rewarded for making contact with the ball, it might learn to repeatedly touch and untouch the ball (toggling between s and s'), rather than trying to win the original game. But in [NHR99], they prove that if the shaping function has the form

$$F(s, a, s') = \gamma\Phi(s') - \Phi(s) \quad (1.39)$$

where $\Phi : \mathcal{S} \rightarrow \mathbb{R}$ is a **potential function**, then we can guarantee that the sum of shaped rewards will match the sum of original rewards plus a constant. This is called **Potential-Based Reward Shaping**.

In [Wie03], they prove that (in the tabular case) this approach is equivalent to initializing the value function to $V(s) = \Phi(s)$. In [TMM19], they propose an extension called potential-based advice, where they show that a potential of the form $F(s, a, s', a') = \gamma\Phi(s', a') - \Phi(s, a)$ is also valid (and more expressive). In [Hu+20], they introduce a reward shaping function z which can be used to down-weight or up-weight the shaping function:

$$r'(s, a) = r(s, a) + z_\phi(s, a)F(s, a) \quad (1.40)$$

They use bilevel optimization to optimize ϕ wrt the original task performance.

1.3.6.6 Intrinsic reward

In Section 7.4, we discuss **intrinsic reward**, which is a set of methods for encouraging agent behavior without the need for any external reward signal. For example, we might want agents to explore their environment just so they can “figure things out”, without any other specific goals in mind. This can be useful even if there is an external reward, but it happens to be sparse.

1.3.7 Best practices for experimental work in RL

Implementing RL algorithms is much trickier than methods for supervised learning, or generative methods such as language modeling and diffusion, all of which have stable (easy-to-optimize) loss functions. Therefore it is often wise to build on existing software rather than starting from scratch. We list some useful libraries in Table 1.3.

Even with good code, RL experiments can be very high variance, making it hard to draw valid conclusions from an experiment. See [Aga+21b; Pat+24; Jor+24] for some recommended experimental practices. For example, when reporting performance across different environments, with different intrinsic difficulties (e.g., different kinds of Atari games), [Aga+21b] recommend reporting the **interquartile mean** (IQM) of the performance metric, which is the mean of the samples between the 0.25 and 0.75 percentiles, (this is a special case of a trimmed mean). Let this estimate be denoted by $\hat{\mu}(\mathcal{D}_i)$, where \mathcal{D} is the empirical data (e.g., reward vs time) from the i ’th run. We can estimate the uncertainty in this estimate using a nonparametric method, such as bootstrap resampling, or a parametric approximation, such as a Gaussian approximation. (This requires computing the standard error of the mean, $\frac{\hat{\sigma}}{\sqrt{n}}$, where n is the number of trials, and $\hat{\sigma}$ is the estimated standard deviation of the (trimmed) data.)

URL	Language	Comments
Stoix	Jax	Mini-library with many methods (including MBRL)
PureJaxRL	Jax	Single files with DQN; PPO, DPO
JaxRL	Jax	Single files with AWAC, DDPG, SAC, SAC+REDQ
Stable Baselines Jax	Jax	Library with DQN, CrossQ, TQC; PPO, DDPG, TD3, SAC
Jax Baselines	Jax	Library with many methods
Rejax	Jax	Library with DDQN, PPO, (discrete) SAC, DDPG
Dopamine	Jax/TF	Library with many methods
Rlax	Jax	Library of RL utility functions (used by Acme)
Acme	Jax/TF	Library with many methods (uses rlax)
CleanRL	PyTorch	Single files with many methods
Stable Baselines 3	PyTorch	Library with DQN; A2C, PPO, DDPG, TD3, SAC, HER
TianShou	PyTorch	Library with many methods (including offline RL)

Table 1.3: Some open source RL software.

Chapter 2

Value-based RL

2.1 Basic concepts

In this section we introduce some definitions and basic concepts.

2.1.1 Value functions

Let π be a given policy. We define the **state-value function**, or **value function** for short, as follows (with $\mathbb{E}_\pi[\cdot]$ indicating that actions are selected by π):

$$V_\pi(s) \triangleq \mathbb{E}_\pi[G_0 | s_0 = s] = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s \right] \quad (2.1)$$

This is the expected return obtained if we start in state s and follow π to choose actions in a continuing task (i.e., $T = \infty$).

Similarly, we define the **state-action value function**, also known as the **Q -function**, as follows:

$$Q_\pi(s, a) \triangleq \mathbb{E}_\pi[G_0 | s_0 = s, a_0 = a] = \mathbb{E}_\pi \left[\sum_{t=0}^{\infty} \gamma^t r_t | s_0 = s, a_0 = a \right] \quad (2.2)$$

This quantity represents the expected return obtained if we start by taking action a in state s , and then follow π to choose actions thereafter.

Finally, we define the **advantage function** as follows:

$$\text{Adv}_\pi(s, a) \triangleq Q_\pi(s, a) - V_\pi(s) \quad (2.3)$$

This tells us the benefit of picking action a in state s then switching to policy π , relative to the baseline return of always following π . Note that $\text{Adv}_\pi(s, a)$ can be both positive and negative, and $\mathbb{E}_{\pi(a|s)}[\text{Adv}_\pi(s, a)] = 0$ due to a useful equality: $V_\pi(s) = \mathbb{E}_{\pi(a|s)}[Q_\pi(s, a)]$.

2.1.2 Bellman's equations

Suppose π^* is a policy such that $V_{\pi^*} \geq V_\pi$ for all $s \in \mathcal{S}$ and all policy π , then it is an **optimal policy**. There can be multiple optimal policies for the same MDP, but by definition their value functions must be the same, and are denoted by V^* and Q^* , respectively. We call V^* the **optimal state-value function**, and Q^* the **optimal action-value function**. Furthermore, any finite MDP must have at least one deterministic optimal policy [Put94].

A fundamental result about the optimal value function is **Bellman’s optimality equations**:

$$V^*(s) = \max_a R(s, a) + \gamma \mathbb{E}_{p_S(s'|s, a)} [V^*(s')] \quad (2.4)$$

$$Q^*(s, a) = R(s, a) + \gamma \mathbb{E}_{p_S(s'|s, a)} \left[\max_{a'} Q^*(s', a') \right] \quad (2.5)$$

Conversely, the optimal value functions are the only solutions that satisfy the equations. In other words, although the value function is defined as the expectation of a sum of infinitely many rewards, it can be characterized by a recursive equation that involves only one-step transition and reward models of the MDP. Such a recursion plays a central role in many RL algorithms we will see later.

Given a value function (V or Q), the discrepancy between the right- and left-hand sides of Equations (2.4) and (2.5) are called **Bellman error** or **Bellman residual**. We can define the **Bellman operator** \mathcal{B} given an MDP $M = (R, T)$ and policy π as a function that takes a value function V and derives a new value function V' that satisfies

$$V'(s) = \mathcal{B}_M^\pi V(s) \triangleq \mathbb{E}_{\pi(a|s)} [R(s, a) + \gamma \mathbb{E}_{T(s'|s, a)} [V(s')]] \quad (2.6)$$

This reduces the Bellman error. Applying the Bellman operator to a state is called a **Bellman backup**. If we iterate this process, we will converge to the optimal value function V_* , as we discuss in Section 2.2.1.

Given the optimal value function, we can derive an optimal policy using

$$\pi^*(s) = \operatorname{argmax}_a Q^*(s, a) \quad (2.7)$$

$$= \operatorname{argmax}_a [R(s, a) + \gamma \mathbb{E}_{p_S(s'|s, a)} [V^*(s')]] \quad (2.8)$$

Following such an optimal policy ensures the agent achieves maximum expected return starting from any state.

The problem of solving for V^* , Q^* or π^* is called **policy optimization**. In contrast, solving for V_π or Q_π for a given policy π is called **policy evaluation**, which constitutes an important subclass of RL problems as will be discussed in later sections. For policy evaluation, we have similar Bellman equations, which simply replace $\max_a \{\cdot\}$ in Equations (2.4) and (2.5) with $\mathbb{E}_{\pi(a|s)} [\cdot]$.

In Equations (2.7) and (2.8), as in the Bellman optimality equations, we must take a maximum over all actions in \mathcal{A} , and the maximizing action is called the **greedy action** with respect to the value functions, Q^* or V^* . Finding greedy actions is computationally easy if \mathcal{A} is a small finite set. For high dimensional continuous spaces, see Section 2.5.4.1.

2.1.3 Example: 1d grid world

In this section, we show a simple example, to make some of the above concepts more concrete. Consider the 1d **grid world** shown in Figure 2.1(a). There are 5 possible states, among them S_{T1} and S_{T2} are absorbing states, since the interaction ends once the agent enters them. There are 2 actions, \uparrow and \downarrow . The reward function is zero everywhere except at the goal state, S_{T2} , which gives a reward of 1 upon entering. Thus the optimal action in every state is to move down.

Figure 2.1(b) shows the Q^* function for $\gamma = 0$. Note that we only show the function for non-absorbing states, as the optimal Q -values are 0 in absorbing states by definition. We see that $Q^*(s_3, \downarrow) = 1.0$, since the agent will get a reward of 1.0 on the next step if it moves down from s_3 ; however, $Q^*(s, a) = 0$ for all other state-action pairs, since they do not provide nonzero immediate reward. This optimal Q -function reflects the fact that using $\gamma = 0$ is completely myopic, and ignores the future.

Figure 2.1(c) shows Q^* when $\gamma = 1$. In this case, we care about all future rewards equally. Thus $Q^*(s, a) = 1$ for all state-action pairs, since the agent can always reach the goal eventually. This is infinitely far-sighted. However, it does not give the agent any short-term guidance on how to behave. For example, in s_2 , it is not clear if it should go up or down, since both actions will eventually reach the goal with identical Q^* -values.

Figure 2.1(d) shows Q^* when $\gamma = 0.9$. This reflects a preference for near-term rewards, while also taking future reward into account. This encourages the agent to seek the shortest path to the goal, which is usually

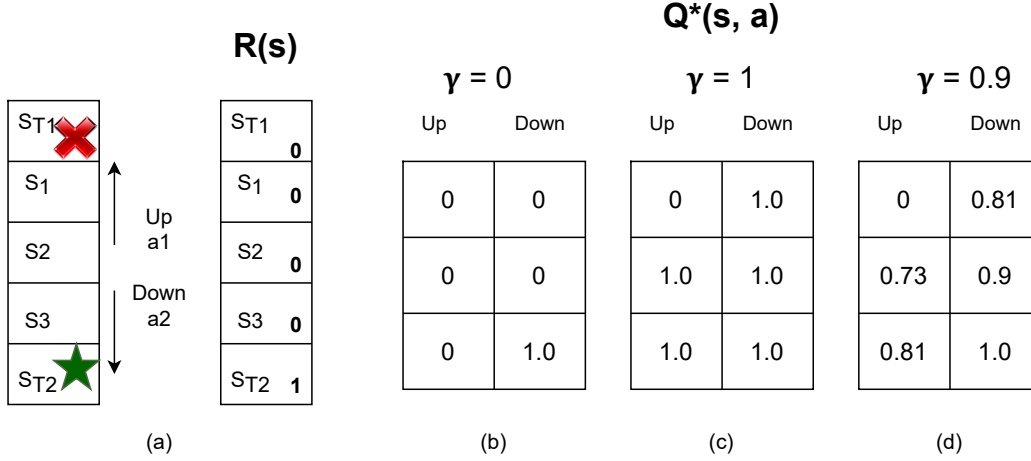


Figure 2.1: Left: illustration of a simple MDP corresponding to a 1d grid world of 3 non-absorbing states and 2 actions. Right: optimal Q-functions for different values of γ . Adapted from Figures 3.1, 3.2, 3.4 of [GK19].

what we desire. A proper choice of γ is up to the agent designer, just like the design of the reward function, and has to reflect the desired behavior of the agent.

2.2 Solving for the optimal policy in a known world model

In this section, we discuss how to compute the optimal value function (the **prediction problem**) and the optimal policy (the **control problem**) when the MDP model is known. (Sometimes the term **planning** is used to refer to computing the optimal policy, given a known model, but planning can also refer to computing a sequence of actions, rather than a policy.) The algorithms we discuss are based on **dynamic programming** (DP) and **linear programming** (LP).

For simplicity, in this section, we assume discrete state and action sets with $\gamma < 1$. However, exact calculation of optimal policies often depends polynomially on the sizes of \mathcal{S} and \mathcal{A} , and is intractable, for example, when the state space is a Cartesian product of several finite sets. This challenge is known as the **curse of dimensionality**. Therefore, approximations are typically needed, such as using parametric or nonparametric representations of the value function or policy, both for computational tractability and for extending the methods to handle MDPs with general state and action sets. This requires the use of **approximate dynamic programming** (ADP) and **approximate linear programming** (ALP) algorithms (see e.g., [Ber19]).

2.2.1 Value iteration

A popular and effective DP method for solving an MDP is **value iteration** (VI). Starting from an initial value function estimate V_0 , the algorithm iteratively updates the estimate by

$$V_{k+1}(s) = \max_a \left[R(s, a) + \gamma \sum_{s'} p(s'|s, a) V_k(s') \right] \quad (2.9)$$

Note that the update rule, sometimes called a **Bellman backup**, is exactly the right-hand side of the Bellman optimality equation Equation (2.4), with the unknown V^* replaced by the current estimate V_k . A fundamental property of Equation (2.9) is that the update is a **contraction**: it can be verified that

$$\max_s |V_{k+1}(s) - V^*(s)| \leq \gamma \max_s |V_k(s) - V^*(s)| \quad (2.10)$$

In other words, every iteration will reduce the maximum value function error by a constant factor.

V_k will converge to V^* , after which an optimal policy can be extracted using Equation (2.8). In practice, we can often terminate VI when V_k is close enough to V^* , since the resulting greedy policy wrt V_k will be near optimal. Value iteration can be adapted to learn the optimal action-value function Q^* .

2.2.2 Real-time dynamic programming (RTDP)

In value iteration, we compute $V^*(s)$ and $\pi^*(s)$ for all possible states s , averaging over all possible next states s' at each iteration, as illustrated in Figure 2.2(right). However, for some problems, we may only be interested in the value (and policy) for certain special starting states. This is the case, for example, in **shortest path problems** on graphs, where we are trying to find the shortest route from the current state to a goal state. This can be modeled as an episodic MDP by defining a transition matrix $p_S(s'|s, a)$ where taking edge a from node s leads to the neighboring node s' with probability 1. The reward function is defined as $R(s, a) = -1$ for all states s except the goal states, which are modeled as absorbing states.

In problems such as this, we can use a method known as **real-time dynamic programming** or **RTDP** [BBS95], to efficiently compute an **optimal partial policy**, which only specifies what to do for the reachable states. RTDP maintains a value function estimate V . At each step, it performs a Bellman backup for the current state s by $V(s) \leftarrow \max_a \mathbb{E}_{p_S(s'|s, a)} [R(s, a) + \gamma V(s')]$. It picks an action a (often with some exploration), reaches a next state s' , and repeats the process. This can be seen as a form of the more general **asynchronous value iteration**, that focuses its computational effort on parts of the state space that are more likely to be reachable from the current state, rather than synchronously updating all states at each iteration.

2.2.3 Policy iteration

Another effective DP method for computing π^* is **policy iteration**. It is an iterative algorithm that searches in the space of deterministic policies until converging to an optimal policy. Each iteration consists of two steps, **policy evaluation** and **policy improvement**.

The policy evaluation step, as mentioned earlier, computes the value function for the current policy. Let π represent the current policy, $\mathbf{v}(s) = V_\pi(s)$ represent the value function encoded as a vector indexed by states, $\mathbf{r}(s) = \sum_a \pi(a|s)R(s, a)$ represent the reward vector, and $\mathbf{T}(s'|s) = \sum_a \pi(a|s)p(s'|s, a)$ represent the state transition matrix. Bellman's equation for policy evaluation can be written in the matrix-vector form as

$$\mathbf{v} = \mathbf{r} + \gamma \mathbf{T} \mathbf{v} \quad (2.11)$$

This is a linear system of equations in $|\mathcal{S}|$ unknowns. We can solve it using matrix inversion: $\mathbf{v} = (\mathbf{I} - \gamma \mathbf{T})^{-1} \mathbf{r}$. Alternatively, we can use value iteration by computing $\mathbf{v}_{t+1} = \mathbf{r} + \gamma \mathbf{T} \mathbf{v}_t$ until near convergence, or some form of asynchronous variant that is computationally more efficient.

Once we have evaluated V_π for the current policy π , we can use it to derive a better policy π' , thus the name policy improvement. To do this, we simply compute a deterministic policy π' that acts greedily with respect to V_π in every state, using

$$\pi'(s) = \underset{a}{\operatorname{argmax}} \{R(s, a) + \gamma \mathbb{E} [V_\pi(s')]\} \quad (2.12)$$

We can guarantee that $V_{\pi'} \geq V_\pi$. This is called the **policy improvement theorem**. To see this, define \mathbf{r}' , \mathbf{T}' and \mathbf{v}' as before, but for the new policy π' . The definition of π' implies $\mathbf{r}' + \gamma \mathbf{T}' \mathbf{v} \geq \mathbf{r} + \gamma \mathbf{T} \mathbf{v} = \mathbf{v}$, where the equality is due to Bellman's equation. Repeating the same equality, we have

$$\mathbf{v} \leq \mathbf{r}' + \gamma \mathbf{T}' \mathbf{v} \leq \mathbf{r}' + \gamma \mathbf{T}' (\mathbf{r}' + \gamma \mathbf{T}' \mathbf{v}) \leq \mathbf{r}' + \gamma \mathbf{T}' (\mathbf{r}' + \gamma \mathbf{T}' (\mathbf{r}' + \gamma \mathbf{T}' \mathbf{v})) \leq \dots \quad (2.13)$$

$$= (\mathbf{I} + \gamma \mathbf{T}' + \gamma^2 \mathbf{T}'^2 + \dots) \mathbf{r}' = (\mathbf{I} - \gamma \mathbf{T}')^{-1} \mathbf{r}' = \mathbf{v}' \quad (2.14)$$

Starting from an initial policy π_0 , policy iteration alternates between policy evaluation (E) and improvement (I) steps, as illustrated below:

$$\pi_0 \xrightarrow{E} V_{\pi_0} \xrightarrow{I} \pi_1 \xrightarrow{E} V_{\pi_1} \dots \xrightarrow{I} \pi^* \xrightarrow{E} V^* \quad (2.15)$$

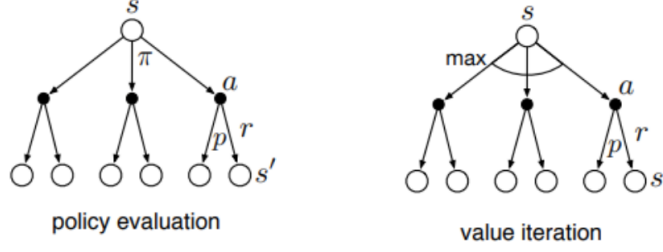


Figure 2.2: Policy iteration vs value iteration represented as backup diagrams. Empty circles represent states, solid (filled) circles represent states and actions. Adapted from Figure 8.6 of [SB18].

The algorithm stops at iteration k , if the policy π_k is greedy with respect to its own value function V_{π_k} . In this case, the policy is optimal. Since there are at most $|\mathcal{A}|^{|\mathcal{S}|}$ deterministic policies, and every iteration strictly improves the policy, the algorithm must converge after finite iterations.

In PI, we alternate between policy evaluation (which involves multiple iterations, until convergence of V_{π}), and policy improvement. In VI, we alternate between one iteration of policy evaluation followed by one iteration of policy improvement (the “max” operator in the update rule). We are in fact free to intermix any number of these steps in any order. The process will converge once the policy is greedy wrt its own value function.

Note that policy evaluation computes V_{π} whereas value iteration computes V^* . This difference is illustrated in Figure 2.2, using a **backup diagram**. Here the root node represents any state s , nodes at the next level represent state-action combinations (solid circles), and nodes at the leaves representing the set of possible resulting next state s' for each possible action. In PE, we average over all actions according to the policy, whereas in VI, we take the maximum over all actions.

2.3 Value function learning using samples from the world model

In the rest of this chapter, we assume the agent only has access to samples from the environment, $(s', r) \sim p(s', r|s, a)$. We will show how to use these samples to estimate optimal value function and Q -function, even without explicitly knowing the MDP dynamics. This is sometimes called “learning” as opposed to “planning”, since the latter requires access to an explicit world model.

2.3.1 Monte Carlo estimation

Recall that $V_{\pi}(s) = \mathbb{E}[G_t | s_t = s]$ is the sum of expected (discounted) returns from state s if we follow policy π . A simple way to estimate this is to rollout the policy, and then compute the average sum of discounted rewards. The trajectory ends when we reach a terminal state, if the task is episodic, or when the discount factor γ^t becomes negligibly small, whichever occurs first. This is called **Monte Carlo estimation**. We can use this to update our estimate of the value function as follows:

$$V(s_t) \leftarrow V(s_t) + \eta [G_t - V(s_t)] \quad (2.16)$$

where η is the learning rate, and the term in brackets is an error term. We can use a similar technique to estimate $Q_{\pi}(s, a) = \mathbb{E}[G_t | s_t = s, a_t = a]$ by simply starting the rollout with action a .

We can use MC estimation of Q , together with policy iteration (Section 2.2.3), to learn an optimal policy. Specifically, at iteration k , we compute a new, improved policy using $\pi_{k+1}(s) = \operatorname{argmax}_a Q_k(s, a)$, where Q_k is approximated using MC estimation. This update can be applied to all the states visited on the sampled trajectory. This overall technique is called **Monte Carlo control**.

To ensure this method converges to the optimal policy, we need to collect data for every (state, action) pair, at least in the tabular case, since there is no generalization across different values of $Q(s, a)$. One way

to achieve this is to use an ϵ -greedy policy (see Section 1.3.5). Since this is an on-policy algorithm, the resulting method will converge to the optimal ϵ -soft policy, as opposed to the optimal policy. It is possible to use importance sampling to estimate the value function for the optimal policy, even if actions are chosen according to the ϵ -greedy policy. However, it is simpler to just gradually reduce ϵ .

2.3.2 Temporal difference (TD) learning

The Monte Carlo (MC) method in Section 2.3.1 results in an estimator for $V(s)$ with very high variance, since it has to unroll many trajectories, whose returns are a sum of many random rewards generated by stochastic state transitions. In addition, it is limited to episodic tasks (or finite horizon truncation of continuing tasks), since it must unroll to the end of the episode before each update step, to ensure it reliably estimates the long term return.

In this section, we discuss a more efficient technique called **temporal difference** or **TD** learning [Sut88]. The basic idea is to incrementally reduce the Bellman error for sampled states or state-actions, based on transitions instead of a long trajectory. More precisely, suppose we are to learn the value function V_π for a fixed policy π . Given a state transition (s_t, a_t, r_t, s_{t+1}) , where $a_t \sim \pi(s_t)$, we change the estimate $V(s_t)$ so that it moves towards the **target value** $y_t = r_t + \gamma V(s_{t+1}) \approx G_{t:t+1}$:

$$V(s_t) \leftarrow V(s_t) + \eta \left[\underbrace{r_t + \gamma V(s_{t+1}) - V(s_t)}_{\delta_t} \right] \quad (2.17)$$

where η is the learning rate. (See [RFP15] for ways to adaptively set the learning rate.) The $\delta_t = y_t - V(s_t)$ term is known as the **TD error**.¹

A more general form of TD update for parametric value function representations is

$$\mathbf{w} \leftarrow \mathbf{w} + \eta [r_t + \gamma V_{\mathbf{w}}(s_{t+1}) - V_{\mathbf{w}}(s_t)] \nabla_{\mathbf{w}} V_{\mathbf{w}}(s_t) \quad (2.18)$$

we see that Equation (2.16) is a special case. The TD update rule for evaluating Q_π is similar, except we replace states with states and actions.

It can be shown that TD learning in the tabular case, Equation (2.16), converges to the correct value function, under proper conditions [Ber19]. However, it may diverge when using nonlinear function approximators, as we discuss in Section 2.5.2.5. The reason is that this update is a “**semi-gradient**”, which refers to the fact that we only take the gradient wrt the value function, $\nabla_{\mathbf{w}} V(s_t, \mathbf{w}_t)$, treating the target U_t as constant.

The potential divergence of TD is also consistent with the fact that Equation (2.18) does not correspond to a gradient update on any objective function, despite having a very similar form to SGD (stochastic gradient descent). Instead, it is an example of **bootstrapping**, in which the estimate, $V_{\mathbf{w}}(s_t)$, is updated to approach a target, $r_t + \gamma V_{\mathbf{w}}(s_{t+1})$, which is defined by the value function estimate itself. This idea is shared by DP methods like value iteration, although they rely on the complete MDP model to compute an exact Bellman backup. In contrast, TD learning can be viewed as using sampled transitions to approximate such backups. An example of a non-bootstrapping approach is the Monte Carlo estimation in the previous section. It samples a complete trajectory, rather than individual transitions, to perform an update; this avoids the divergence issue, but is often much less efficient. Figure 2.3 illustrates the difference between MC, TD, and DP.

2.3.3 Combining TD and MC learning using TD(λ)

A key difference between TD and MC is the way they estimate returns. Given a trajectory $\boldsymbol{\tau} = (s_0, a_0, r_0, s_1, \dots, s_T)$, TD estimates the return from state s_t by one-step lookahead, $G_{t:t+1} = r_t + \gamma V(s_{t+1})$, where the return from time $t+1$ is replaced by its value function estimate. In contrast, MC waits until the end of the episode or until T is large enough, then uses the estimate $G_{t:T} = r_t + \gamma r_{t+1} + \dots + \gamma^{T-t-1} r_{T-1}$. It is possible to

¹The **Rescorla-Wagner rule** is a special case of this, studied in the **classical conditioning** literature. The update is given by $V_{t+1}(s_t) \leftarrow V_t(s_t) + \eta(r_t - V_t(s_t))$, where s_t is a stimulus and r_t is the response associated with it. This is TD with $\gamma = 0$.

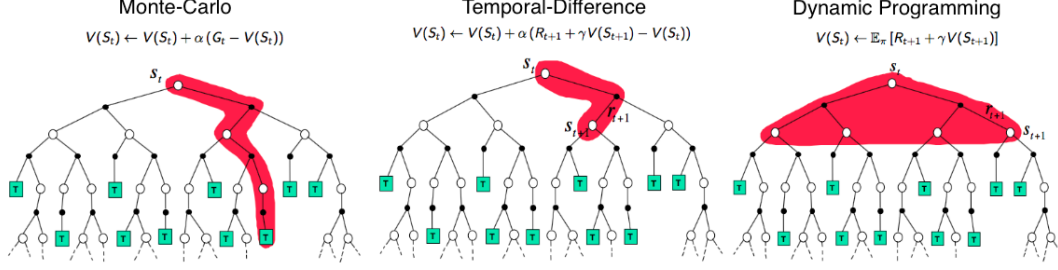


Figure 2.3: Backup diagrams of $V(s_t)$ for Monte Carlo, temporal difference, and dynamic programming updates of the state-value function. Used with kind permission of Andy Barto.

interpolate between these by performing an n -step rollout, and then using the value function to approximate the return for the rest of the trajectory, similar to heuristic search (Section 4.2.1.4). That is, we can use the **n-step return**

$$G_{t:t+n} = r_t + \gamma r_{t+1} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma^n V(s_{t+n}) \quad (2.19)$$

For example, the 1-step and 2-step returns are given by

$$G_{t:t+1} = r_t + \gamma v_{t+1} \quad (2.20)$$

$$G_{t:t+2} = r_t + \gamma r_{t+1} + \gamma^2 v_{t+2} \quad (2.21)$$

The corresponding n -step version of the TD update becomes

$$\mathbf{w} \leftarrow \mathbf{w} + \eta [G_{t:t+n} - V_{\mathbf{w}}(s_t)] \nabla_{\mathbf{w}} V_{\mathbf{w}}(s_t) \quad (2.22)$$

Using this update can help propagate sparse terminal rewards back through many earlier states.

Rather than picking a specific lookahead value, n , we can take a weighted average of all possible values, with a single parameter $\lambda \in [0, 1]$, by using

$$G_t^\lambda \triangleq (1 - \lambda) \sum_{n=1}^{\infty} \lambda^{n-1} G_{t:t+n} \quad (2.23)$$

This is called the **lambda return**. Note that these coefficients sum to one (since $\sum_{t=0}^{\infty} (1 - \lambda) \lambda^t = \frac{1 - \lambda}{1 - \lambda} = 1$, for $\lambda < 1$), so the return is a convex combination of n -step returns. See Figure 2.4 for an illustration. We can now use G_t^λ inside the TD update instead of $G_{t:t+n}$; this is called **TD(λ)**.

Note that, if a terminal state is entered at step T (as happens with episodic tasks), then all subsequent n -step returns are equal to the conventional return, G_t . Hence we can write

$$G_t^\lambda = (1 - \lambda) \sum_{n=1}^{T-t-1} \lambda^{n-1} G_{t:t+n} + \lambda^{T-t-1} G_t \quad (2.24)$$

From this we can see that if $\lambda = 1$, the λ -return becomes equal to the regular MC return G_t . If $\lambda = 0$, the λ -return becomes equal to the one-step return $G_{t:t+1}$ (since $0^{n-1} = 1$ iff $n = 1$), so standard TD learning is often called **TD(0) learning**. This episodic form also gives us the following recursive equation

$$G_t^\lambda = r_t + \gamma [(1 - \lambda) v_{t+1} + \lambda G_{t+1}^\lambda] \quad (2.25)$$

which we initialize with $G_T = v_t$.

2.4.2 Sarsa(λ)

It is possible to apply the eligibility trace idea to Sarsa, since it is an on-policy method. This can help speedup learning in sparse reward scenarios.

The basic idea, in the tabular case, is as follows. We compute an eligibility for each state action pair, denoted $N(s, a)$, representing the visit count. Following Equation (2.27), we perform the update

$$Q(s, a) \leftarrow Q(s, a) + \eta \delta N(s, a) \quad (2.29)$$

where the TD error is

$$\delta = r + \gamma Q(s', a') - Q(s, a) \quad (2.30)$$

Then, following Equation (2.26), we decay all the visit counts (traces) using

$$N(s, a) \leftarrow \gamma \lambda N(s, a) \quad (2.31)$$

This is called Sarsa(λ).

2.5 Q-learning: off-policy TD policy learning

SARSA is an on-policy algorithm, which means it learns the Q -function for the policy it is currently using, which is typically not the optimal policy, because of the need to perform exploration. However, with a simple modification, we can convert this to an off-policy algorithm that learns Q^* , even if a suboptimal or exploratory policy is used to choose actions.

2.5.1 Tabular Q learning

Suppose we modify SARSA by replacing the sampled next action $a' \sim \pi(s')$ in Equation (2.28) with a greedy action: $a' = \operatorname{argmax}_b Q(s', b)$. This results in the following update when a transition (s, a, r, s') happens

$$Q(s, a) \leftarrow Q(s, a) + \eta \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right] \quad (2.32)$$

This is the update rule of **Q-learning** for the tabular case [WD92].

Since it is off-policy, the method can use (s, a, r, s') triples coming from any data source, such as older versions of the policy, or log data from an existing (non-RL) system. If every state-action pair is visited infinitely often, the algorithm provably converges to Q^* in the tabular case, with properly decayed learning rates [Ber19]. Algorithm 1 gives a vanilla implementation of Q-learning with ϵ -greedy exploration.

Algorithm 1: Tabular Q-learning with ϵ -greedy exploration

```

1 Initialize value function  $Q$ 
2 repeat
3   Sample starting state  $s$  of new episode
4   repeat
5     Sample action  $a = \begin{cases} \operatorname{argmax}_b Q(s, b), & \text{with probability } 1 - \epsilon \\ \text{random action}, & \text{with probability } \epsilon \end{cases}$ 
6      $(s', r) = \text{env.step}(a)$ 
7     Compute the TD error:  $\delta = r + \gamma \max_{a'} Q(s', a') - Q(s, a)$ 
8      $Q(s, a) \leftarrow Q(s, a) + \eta \delta$ 
9      $s \leftarrow s'$ 
10  until state  $s$  is terminal;
11 until converged;
```

Q-function episode start		Episode	Time Step	Action	(s, a, r, s')	$r + \gamma Q^*(s', a)$	Q-function episode end	
UP DOWN							UP DOWN	
Q_1	S_1		1	↓	$(S_1, D, 0, S_2)$	$0 + 0.9 \times 0 = 0$	S_1	0
			1	↑	$(S_2, U, 0, S_1)$	$0 + 0.9 \times 0 = 0$		0
	S_2		1	↓	$(S_1, D, 0, S_1)$	$0 + 0.9 \times 0 = 0$	S_2	0
			1	↓	$(S_2, U, 0, S_1)$	$0 + 0.9 \times 0 = 0$		0
	S_3		1	↓	$(S_3, D, 1, S_{T2})$	1	S_3	0
Q_2	S_1		2	↓	$(S_1, D, 0, S_2)$	$0 + 0.9 \times 0 = 0$	S_1	0
	S_2		2	↓	$(S_2, D, 0, S_3)$	$0 + 0.9 \times 1 = 0.9$	S_2	0
	S_3		2	↓	$(S_3, D, 0, S_{T2})$	1	S_3	0
Q_3	S_1		3	↓	$(S_1, D, 0, S_2)$	$0 + 0.9 \times 0.9 = 0.81$	S_1	0
			3	↓	$(S_2, D, 0, S_3)$	$0 + 0.9 \times 1 = 0.9$		0.81
	S_2		3	↑	$(S_3, D, 0, S_2)$	$0 + 0.9 \times 0.9 = 0.81$	S_2	0
			3	↓	$(S_2, D, 0, S_3)$	$0 + 0.9 \times 1 = 0.9$		0.9
	S_3		3	↓	$(S_3, D, 0, S_{T2})$	1	S_3	0.81
Q_4	S_1		4	↓	$(S_1, D, 0, S_2)$	$0 + 0.9 \times 0.9 = 0.81$	S_1	0
			4	↑	$(S_2, U, 0, S_1)$	$0 + 0.9 \times 0.81 = 0.73$		0.81
	S_2		4	↓	$(S_1, D, 0, S_2)$	$0 + 0.9 \times 0.9 = 0.81$	S_2	0.73
			4	↑	$(S_2, U, 0, S_3)$	$0 + 0.9 \times 0.81 = 0.73$		0.9
	S_3		4	↓	$(S_1, D, 0, S_3)$	$0 + 0.9 \times 0.9 = 0.81$	S_3	0.81
			4	↓	$(S_2, D, 0, S_3)$	$0 + 0.9 \times 1 = 0.9$		1
			4	↓	$(S_3, D, 0, S_3)$	1		
Q_5	S_1		5	↑	$(S_1, U, 0, S_{T1})$	0	S_1	0
	S_2						S_2	0.73
	S_3						S_3	0.81

Figure 2.5: Illustration of Q learning for one random trajectory in the 1d grid world in Figure 2.1 using ϵ -greedy exploration. At the end of episode 1, we make a transition from S_3 to S_{T2} and get a reward of $r = 1$, so we estimate $Q(S_3, \downarrow) = 1$. In episode 2, we make a transition from S_2 to S_3 , so S_2 gets incremented by $\gamma Q(S_3, \downarrow) = 0.9$. Adapted from Figure 3.3 of [GK19].

For terminal states, $s \in \mathcal{S}^+$, we know that $Q(s, a) = 0$ for all actions a . Consequently, for the optimal value function, we have $V^*(s) = \max_{a'} Q^*(s, a) = 0$ for all terminal states. When performing online learning, we don't usually know which states are terminal. Therefore we assume that, whenever we take a step in the environment, we get the next state s' and reward r , but also a binary indicator $\text{done}(s')$ that tells us if s' is terminal. In this case, we set the target value in Q-learning to $V^*(s') = 0$ yielding the modified update rule:

$$Q(s, a) \leftarrow Q(s, a) + \eta \left[r + (1 - \text{done}(s')) \gamma \max_{a'} Q(s', a') - Q(s, a) \right] \quad (2.33)$$

For brevity, we will usually ignore this factor in the subsequent equations, but it needs to be implemented in the code.

Figure 2.5 gives an example of Q-learning applied to the simple 1d grid world from Figure 2.1, using $\gamma = 0.9$. We show the Q -function at the start and end of each episode, after performing actions chosen by an ϵ -greedy policy. We initialize $Q(s, a) = 0$ for all entries, and use a step size of $\eta = 1$. At convergence, we have $Q^*(s, a) = r + \gamma Q^*(s', a_*)$, where $a_* = \downarrow$ for all states.

2.5.2 Q learning with function approximation

To make Q learning work with high-dimensional state spaces, we have to replace the tabular (non-parametric) representation with a parametric approximation, denoted $Q_{\mathbf{w}}(s, a)$. We can update this function using one or more steps of SGD on the following loss function

$$\mathcal{L}(\mathbf{w}|s, a, r, s') = \left((r + \gamma \max_{a'} Q_{\mathbf{w}}(s', a')) - Q_{\mathbf{w}}(s, a) \right)^2 \quad (2.34)$$

Since nonlinear functions need to be trained on minibatches of data, we compute the average loss over multiple randomly sampled experience tuples (see Section 2.5.2.3 for discussion) to get

$$\mathcal{L}(\mathbf{w}) = \mathbb{E}_{(s, a, r, s') \sim U(\mathcal{D})} [\mathcal{L}(\mathbf{w}|s, a, r, s')] \quad (2.35)$$

See Algorithm 2 for the pseudocode.

Algorithm 2: Q learning with function approximation and replay buffers

```

1 Initialize environment state  $\mathbf{s}$ , network parameters  $\mathbf{w}_0$ , replay buffer  $\mathcal{D} = \emptyset$ , discount factor  $\gamma$ , step
  size  $\eta$ , policy  $\pi_0(a|s) = \epsilon \text{Unif}(a) + (1 - \epsilon)\delta(a = \arg\max_a Q_{\mathbf{w}_0}(s, a))$ 
2 for iteration  $k = 0, 1, 2, \dots$  do
3   for environment step  $s = 0, 1, \dots, S - 1$  do
4     Sample action:  $a \sim \pi_k(a|s)$ 
5     Interact with environment:  $(s', r) = \text{env.step}(a)$ 
6     Update buffer:  $\mathcal{D} \leftarrow \mathcal{D} \cup \{(s, a, s', r)\}$ 
7    $\mathbf{w}_{k,0} \leftarrow \mathbf{w}_k$ 
8   for gradient step  $g = 0, 1, \dots, G - 1$  do
9     Sample batch:  $B \subset \mathcal{D}$ 
10    Compute error:  $\mathcal{L}(B, \mathbf{w}_{k,g}) = \frac{1}{|B|} \sum_{(s, a, r, s') \in B} [Q_{\mathbf{w}_{k,g}}(s, a) - (r + \gamma \max_{a'} Q_{\mathbf{w}_{k,g}}(s', a'))]^2$ 
11    Update parameters:  $\mathbf{w}_{k,g} \leftarrow \mathbf{w}_{k,g} - \eta \nabla_{\mathbf{w}_{k,g}} \mathcal{L}(B, \mathbf{w}_{k,g})$ 
12   $\mathbf{w}_{k+1} \leftarrow \mathbf{w}_{k,G}$ 

```

2.5.2.1 Neural fitted Q

The first approach of this kind is known as **fitted Q evaluation** (or **FQE**) [EGW05], which was extended in [Rie05] to use neural networks. This corresponds to fully optimizing $\mathcal{L}(\mathbf{w})$ at each iteration (equivalent to using $G = \infty$ gradient steps).

2.5.2.2 DQN

The influential deep Q-network or DQN paper of [Mni+15] also used neural nets to represent the Q function, but performed a smaller number of gradient updates per iteration. Furthermore, they proposed to modify the target value when fitting the Q function in order to avoid instabilities during training (see Section 2.5.2.5 for details).

The DQN method became famous since it was able to train agents that can outperform humans when playing various Atari games from the **ALE** (Atari Learning Environment) benchmark [Bel+13]. Here the input is a small color image, and the action space corresponds to moving left, right, up or down, plus an optional shoot action.²

²For more discussion of ALE, see [Mac+18a], and for a recent extension to continuous actions (representing joystick control), see the CALE benchmark of [FC24]. Note that DQN was not the first deep RL method to train an agent from pixel input; that honor goes to [LR10], who trained an autoencoder to embed images into low-dimensional latents, and then used neural fitted Q learning (Section 2.5.2.1) to fit the Q function.

Since 2015, many more extensions to DQN have been proposed, with the goal of improving performance in various ways, either in terms of peak reward obtained, or sample efficiency (e.g., reward obtained after only 100k steps in the environment, as proposed in the **Atari-100k** benchmark [Kai+19]³), or training stability, or all of the above. We discuss some of these extensions in Section 2.5.4.

2.5.2.3 Experience replay

Since Q learning is an off-policy method, we can update the Q function using any data source. This is particularly important when we use nonlinear function approximation (see Section 2.5.2), which often needs a lot of data for model fitting. A natural source of data is data collected earlier in the trajectory of the agent; this is called an **experience replay** buffer, which stores (s, a, r, s') transition tuples into a buffer. This can improve the stability and sample efficiency of learning, and was originally proposed in [Lin92].

This modification has two advantages. First, it improves data efficiency as every transition can be used multiple times. Second, it improves stability in training, by reducing the correlation of the data samples that the network is trained on, since the training tuples do not have to come from adjacent moments in time. (Note that experience replay requires the use of off-policy learning methods, such as Q learning, since the training data is sampled from older versions of the policy, not the current policy.)

2.5.2.4 Prioritized experience replay

It is possible to replace the uniform sampling from the buffer with one that favors more important transition tuples that may be more informative about Q .

To explain this, we start by discussing the **prioritized sweeping** method of [MA93], which was developed for discrete state spaces using a priority queue. The idea is as follows. Whenever we update the value of a state $V(s)$, we iterate over all state-action pairs (s^-, a^-) that can immediately transition into s (this requires knowing the world model). The priority of any such s^- is then increased to $\mathcal{T}(s|s^-, a^-) \times |V(s) - V^{\text{old}}(s)|$, where $V^{\text{old}}(s)$ is the value before the update. Thus we prioritize updating states which are likely to have lead to states whose values have changed by a lot.

This can be generalized to the non-tabular experience replay setting as described in [Sch+16a], who call the technique **prioritized experience replay**. Consider the TD error for the i 'th tuple τ_i

$$\delta_i = r_i + \gamma \max_{a'} Q_{\bar{w}}(s'_i, a') - Q_w(s_i, a_i) \quad (2.36)$$

Define the priority of i as

$$p_i = (\delta_i + \epsilon)^\alpha \quad (2.37)$$

where $\alpha \geq 0$ determines the degree of prioritization, with $\alpha = 0$ corresponding to no prioritization (uniform sampling). Now define the probability of sampling i as

$$P(i) = \frac{p_i}{\sum_k p_k} \quad (2.38)$$

2.5.2.5 The deadly triad

The problem with the naive Q learning objective in Equation (2.34) is that it can lead to instability, since the target we are regressing towards uses the same parameters w as the function we are updating. So the network is “chasing its own tail”. Although this is fine for tabular models, it can fail for nonlinear models, as we discuss below.

In general, an RL algorithm can become unstable when it has these three components: function approximation (such as neural networks), bootstrapped value function estimation (i.e., using TD-like methods instead

³The Atari-100k benchmark only includes 26 out of 46 games of the ALE that were determined to be “solvable by state-of-the-art model-free deep RL algorithms” at the time of the benchmark’s creation in 2019. This excludes games like Montezuma’s Revenge, which require more exploration and hence more training data.

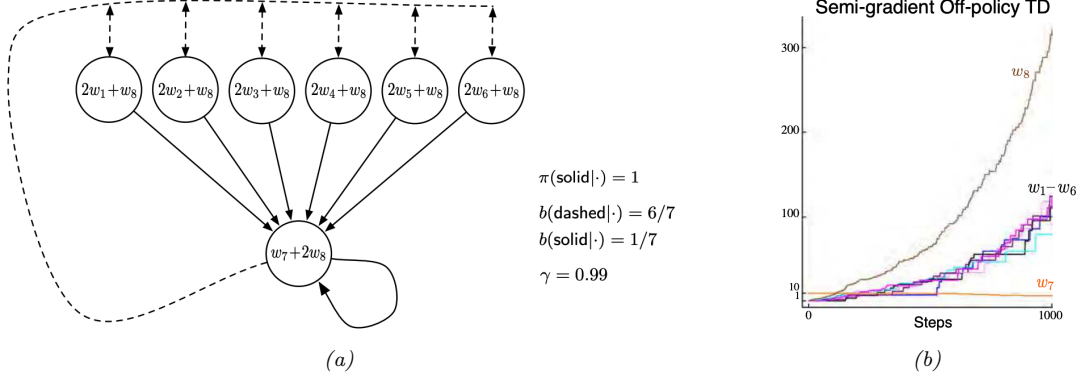


Figure 2.6: (a) A simple MDP. (b) Parameters of the policy diverge over time. From Figures 11.1 and 11.2 of [SB18]. Used with kind permission of Richard Sutton.

of MC), and off-policy learning (where the actions are sampled from some distribution other than the policy that is being optimized). This combination is known as **the deadly triad** [Sut15; van+18]).

A classic example of this is the simple MDP depicted in Figure 2.6a, due to [Bai95]. (This is known as **Baird’s counter example**.) It has 7 states and 2 actions. Taking the dashed action takes the environment to the 6 upper states uniformly at random, while the solid action takes it to the bottom state. The reward is 0 in all transitions, and $\gamma = 0.99$. The value function $V_{\mathbf{w}}$ uses a linear parameterization indicated by the expressions shown inside the states, with $\mathbf{w} \in \mathbb{R}^8$. The target policies π always chooses the solid action in every state. Clearly, the true value function, $V_{\pi}(s) = 0$, can be exactly represented by setting $\mathbf{w} = \mathbf{0}$.

Suppose we use a behavior policy b to generate a trajectory, which chooses the dashed and solid actions with probabilities $6/7$ and $1/7$, respectively, in every state. If we apply TD(0) on this trajectory, the parameters diverge to ∞ (Figure 2.6b), even though the problem appears simple. In contrast, with on-policy data (that is, when b is the same as π), TD(0) with linear approximation can be guaranteed to converge to a good value function approximate [TR97]. The difference is that with on-policy learning, as we improve the value function, we also improve the policy, so the two become self-consistent, whereas with off-policy learning, the behavior policy may not match the optimal value function that is being learned, leading to inconsistencies.

The divergence behavior is demonstrated in many value-based bootstrapping methods, including TD, Q-learning, and related approximate dynamic programming algorithms, where the value function is represented either linearly (like the example above) or nonlinearly [Gor95; TVR97; OCD21]. The root cause of these divergence phenomena is that bootstrapping methods typically are not minimizing a fixed objective function. Rather, they create a learning target using their own estimates, thus potentially creating a self-reinforcing loop to push the estimates to infinity. More formally, the problem is that the contraction property in the tabular case (Equation (2.10)) may no longer hold when V is approximated by $V_{\mathbf{w}}$.

We discuss some solutions to the deadly triad problem below.

2.5.2.6 Target networks

One heuristic solution to the deadly triad, proposed in the DQN paper, is to use a “frozen” **target network** computed at an earlier iteration to define the target value for the DQN updates, rather than trying to chase a constantly moving target. Specifically, we maintain an extra copy the Q -network, $Q_{\bar{\mathbf{w}}}$, with the same structure as $Q_{\mathbf{w}}$. This new Q -network is used to compute bootstrapping targets

$$y(r, s'; \bar{\mathbf{w}}) = r + \gamma \max_{a'} Q_{\bar{\mathbf{w}}}(s', a') \quad (2.39)$$

for training $Q_{\mathbf{w}}$. We can periodically set $\bar{\mathbf{w}} \leftarrow \text{sg}(\mathbf{w})$, usually after a few episodes, where the stop gradient operator is used to prevent autograd propagating gradients back to \mathbf{w} (this is called a **detached target**).

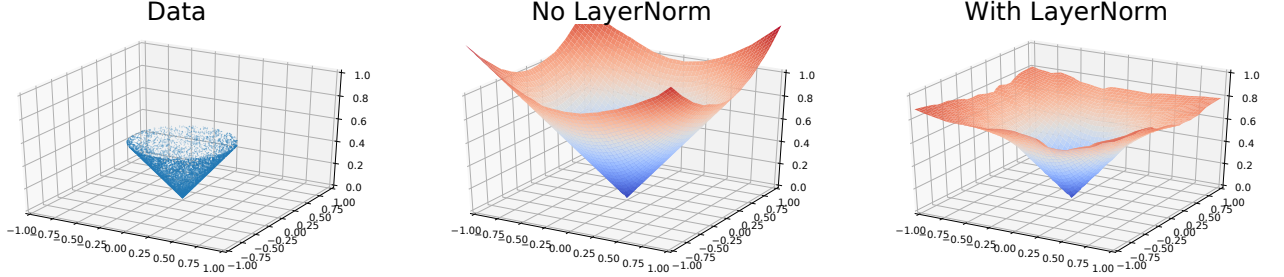


Figure 2.7: We generate a dataset (left) with inputs \mathbf{x} distributed in a circle with radius 0.5 and labels $y = \|\mathbf{x}\|$. We then fit a two-layer MLP without LayerNorm (center) and with LayerNorm (right). LayerNorm bounds the values and prevents catastrophic overestimation when extrapolating. From Figure 3 of [Bal+23]. Used with kind permission of Philip Ball.

Alternatively, we can use an **exponential moving average** or **EMA** of the weights, i.e., at step k we use $\bar{\mathbf{w}}_k = \rho \bar{\mathbf{w}}_{k-1} + (1 - \rho) \mathbf{w}_k$, where the forgetting factor $\rho \approx 0.999$ ensures that $Q_{\bar{\mathbf{w}}}$ slowly catches up with $Q_{\mathbf{w}}$.⁴ The final loss has the form

$$\mathcal{L}(\mathbf{w}) = \mathbb{E}_{(\mathbf{s}, a, r, \mathbf{s}') \sim U(\mathcal{D})} [\mathcal{L}(\mathbf{w} | \mathbf{s}, a, r, \mathbf{s}')] \quad (2.40)$$

$$\mathcal{L}(\mathbf{w} | \mathbf{s}, a, r, \mathbf{s}') = (y(r, \mathbf{s}'; \bar{\mathbf{w}}) - Q_{\mathbf{w}}(\mathbf{s}, a))^2 \quad (2.41)$$

Theoretical work justifying this technique is given in [FSW23; Che+24a].

2.5.2.7 Gradient TD methods

A general way to ensure convergence in off-policy learning is to construct an objective function, the minimization of which leads to a good value function approximation. This is the basis of the **gradient TD method** of [SSM08; Mae+09; Ghi+20].

2.5.2.8 Two time-scale methods

Another approach is to update the target value in the TD update more quickly than the value function itself; this is known as a **two timescale optimization** (see e.g., [Yu17; Zha+19; Hon+23]).

2.5.2.9 Layer norm

More recently, [Gal+24] proved that nonlinear TD learning can be made to converge, even in the off-policy setting, if three simple conditions on the critic (Q network) are satisfied: the final layer weights are bounded (e.g., using ℓ_2 normalization, or using AdamW, which has been shown to solve an L_∞ constrained optimization problem, ensuring that model parameters remain bounded by $1/\lambda$ [XL24; Pet+25]); the penultimate layer is sufficiently wide to represent the true function; and the input to the critic has bounded norm (e.g., using LayerNorm [BKH16], or RMSNorm⁵). In particular, suppose the network has the form $Q(\mathbf{s}, a | \mathbf{w}, \boldsymbol{\theta}) = \mathbf{w}^T \text{ReLU}(\text{LayerNorm}(f(\mathbf{s}, a; \boldsymbol{\theta})))$. Since $\|\text{LayerNorm}(f(\mathbf{s}, a; \boldsymbol{\theta}))\| \leq 1$, we have $\|Q(\mathbf{s}, a | \mathbf{w}, \boldsymbol{\theta})\| \leq \|\mathbf{w}\|$, which means the magnitude of the output is always bounded (assuming $\|\mathbf{w}\|$ is also bounded), as shown in Figure 2.7. In [Gal+24], they prove that this is sufficient to ensure convergence of the value function estimate.

⁴Mathematically, if we have a loss of the form $\mathcal{L}(\theta, \phi) = L(f(\theta) - t(\phi))$, where f is the predictor and t is the target, then using a stop gradient amounts to taking the gradient of \mathcal{L} wrt θ and evaluating the result at $\phi = \theta$. This is the same as using the EMA target $\bar{\phi}$ with $\rho = 0$.

⁵RMSNorm is a simplification of LayerNorm that omits the mean centering step, and has been found to be faster, especially in large models like transformers. In more detail, LayerNorm computes $y = \frac{x - E[x]}{\sqrt{V[x] + \epsilon}} \gamma + \beta$, whereas RMSNorm just computes $y = \frac{x}{\sqrt{\frac{1}{n} \sum_{i=1}^n x_i^2 + \epsilon}} \gamma$.

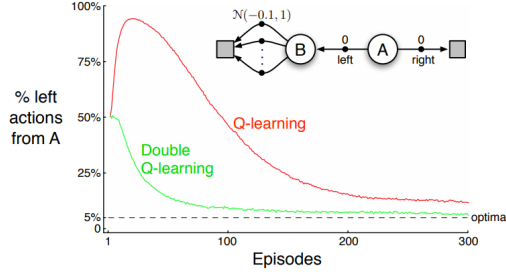


Figure 2.8: Comparison of Q-learning and double Q-learning on a simple episodic MDP using ϵ -greedy action selection with $\epsilon = 0.1$. The initial state is A, and squares denote absorbing states. The data are averaged over 10,000 runs. From Figure 6.5 of [SB18]. Used with kind permission of Richard Sutton.

2.5.2.10 Other methods

A variety of other solutions to the deadly triad have been proposed, including the “chaining value functions” approach of [SSTH22], a combination of target networks and over-parameterized linear function approximation [Che+24a], etc.

2.5.3 Maximization bias

Standard Q-learning suffers from a problem known as the **optimizer’s curse** [SW06], or the **maximization bias**. The problem refers to the simple statistical inequality: $\mathbb{E}[\max_a X_a] \geq \max_a \mathbb{E}[X_a]$, for a set of random variables $\{X_a\}$. Thus, if we pick actions greedily according to their random scores $\{X_a\}$, we might pick a wrong action just because random noise makes it appealing.

Figure 2.8 gives a simple example of how this can happen in an MDP. The start state is A. The right action gives a reward 0 and terminates the episode. The left action also gives a reward of 0, but then enters state B, from which there are many possible actions, with rewards drawn from $\mathcal{N}(-0.1, 1.0)$. Thus the expected return for any trajectory starting with the left action is -0.1 , making it suboptimal. Nevertheless, the RL algorithm may pick the left action due to the maximization bias making B appear to have a positive value.

2.5.3.1 Double Q-learning

One solution to avoid the maximization bias is to use two separate Q-functions, Q_1 and Q_2 , one for selecting the greedy action, and the other for estimating the corresponding Q-value. In particular, upon seeing a transition (s, a, r, s') , we perform the following update for $i = 1 : 2$:

$$Q_i(s, a) \leftarrow Q_i(s, a) + \eta(y_i(s, a) - Q_i(s, a)) \quad (2.42)$$

$$y_i(s, a) = r + \gamma Q_{-i}(s', \arg\max_{a'} Q_{-i}(s', a')) \quad (2.43)$$

So we see that Q_1 uses Q_2 to choose the best action but uses Q_1 to evaluate it, and vice versa. This technique is called **double Q-learning** [Has10]. Figure 2.8 shows the benefits of the algorithm over standard Q-learning in a toy problem.

2.5.3.2 Double DQN

In [HGS16], they combine double Q learning with deep Q networks (Section 2.5.2.2) to get **double DQN**. This modifies Equation (2.43) to its gradient form, and then the current network for action proposals, but the target network for action evaluation. Thus the training target becomes

$$y(r, s'; \mathbf{w}, \bar{\mathbf{w}}) = r + \gamma Q_{\bar{\mathbf{w}}}(s', \arg\max_{a'} Q_{\mathbf{w}}(s', a')) \quad (2.44)$$

In Section 3.2.6.3 we discuss an extension called **clipped double DQN** which uses two Q networks and their frozen copies to define the following target:

$$y(r, s'; \mathbf{w}_{1:2}, \bar{\mathbf{w}}_{1:2}) = r + \gamma \min_{i=1,2} Q_{\bar{\mathbf{w}}_i}(s', \arg\max_{a'} Q_{\mathbf{w}_i}(s', a')) \quad (2.45)$$

where $Q_{\bar{\mathbf{w}}_i}$ is the target network for $Q_{\mathbf{w}_i}$.

2.5.3.3 Randomized ensemble DQN

The double DQN method is extended in the **REDQ** (randomized ensembled double Q learning) method of [Che+20], which uses an ensemble of $N > 2$ Q-networks. Furthermore, at each step, it draws a random sample of $M \leq N$ networks, and takes the minimum over them when computing the target value. That is, it uses the following update (see Algorithm 2 in appendix of [Che+20]):

$$y(r, s'; \mathbf{w}_{1:N}, \bar{\mathbf{w}}_{1:N}) = r + \gamma \max_{a'} \min_{i \in \mathcal{M}} Q_{\bar{\mathbf{w}}_i}(s', a') \quad (2.46)$$

where \mathcal{M} is a random subset from the N value functions. The ensemble reduces the variance, and the minimum reduces the overestimation bias.⁶ If we set $N = M = 2$, we get a method similar to clipped double Q learning. (Note that REDQ is very similar to the **Random Ensemble Mixture** method of [ASN20], which was designed for offline RL.)

2.5.4 DQN extensions

In this section, we discuss various extensions of DQN.

2.5.4.1 Q learning for continuous actions

Q learning is not directly applicable to continuous actions due to the need to compute the argmax over actions. An early solution to this problem, based on neural fitted Q learning (see Section 2.5.2.1), is proposed in [HR11]. This became the basis of the DDPG algorithm of Section 3.2.6.2, which learns a policy to predict the argmax.

An alternative approach is to use gradient-free optimizers such as the cross-entropy method to approximate the argmax. The **QT-Opt** method of [Kal+18] treats the action vector \mathbf{a} as a sequence of actions, and optimizes one dimension at a time [Met+17]. The **CAQL** (continuous action Q-learning) method of [Ryu+20] uses mixed integer programming to solve the argmax problem, leveraging the ReLU structure of the Q-network. The method of [Sey+22] quantizes each action dimension separately, and then solves the argmax problem using methods inspired by multi-agent RL.

2.5.4.2 Dueling DQN

The **dueling DQN** method of [Wan+16], learns a value function and an advantage function, and derives the Q function, rather than learning it directly. This is helpful when there are many actions with similar Q-values, since the advantage $A(s, a) = Q(s, a) - V(s)$ focuses on the differences in value relative to a shared baseline.

In more detail, we define a network with $|A| + 1$ output heads, which computes $A_{\mathbf{w}}(s, a)$ for $a = 1 : A$ and $V_{\mathbf{w}}(s)$. We can then derive

$$Q_{\mathbf{w}}(s, a) = V_{\mathbf{w}}(s) + A_{\mathbf{w}}(s, a) \quad (2.47)$$

However, this naive approach ignores the following constraint that holds for any policy π :

$$\mathbb{E}_{\pi(a|s)} [A^{\pi}(s, a)] = \mathbb{E}_{\pi(a|s)} [Q^{\pi}(s, a) - V^{\pi}(s)] \quad (2.48)$$

$$= V^{\pi}(s) - V^{\pi}(s) = 0 \quad (2.49)$$

⁶In addition, REDQ performs $G \gg 1$ updates of the value functions for each environment step; this high **Update-To-Data** (UTD) ratio (also called **Replay Ratio**) is critical for sample efficiency, and is commonly used in model-based RL.

Fortunately, for the optimal policy $\pi^*(s) = \operatorname{argmax}_{a'} Q^*(s, a')$ we have

$$0 = \mathbb{E}_{\pi^*(a|s)} [Q^*(s, a)] - V^*(s) \quad (2.50)$$

$$= Q^*(s, \operatorname{argmax}_{a'} Q^*(s, a')) - V^*(s) \quad (2.51)$$

$$= \max_{a'} Q^*(s, a') - V^*(s) \quad (2.52)$$

$$= \max_{a'} A^*(s, a') \quad (2.53)$$

Thus we can satisfy the constraint for the optimal policy by subtracting off $\max_a A(s, a)$ from the advantage head. Equivalently we can compute the Q function using

$$Q_{\mathbf{w}}(s, a) = V_{\mathbf{w}}(s) + A_{\mathbf{w}}(s, a) - \max_{a'} A_{\mathbf{w}}(s, a') \quad (2.54)$$

In practice, the max is replaced by an average, which seems to work better empirically.

2.5.4.3 Noisy nets and exploration

Standard DQN relies on the epsilon-greedy strategy to perform exploration. However, this will explore equally in all states, whereas we would like the amount of exploration to be state dependent, to reflect the amount of uncertainty in the outcomes of trying each action in that state due to lack of knowledge (i.e., **epistemic uncertainty** rather than aleatoric or irreducible uncertainty). An early approach to this, known as **noisy nets** [For+18], added random noise to the network weights to encourage exploration which is temporally consistent within episodes. More recent methods for exploration are discussed in Section 1.3.5.

2.5.4.4 Multi-step DQN

As we discussed in Section 2.3.3, we can reduce the bias introduced by bootstrapping by replacing TD(1) updates with TD(n) updates, where we unroll the value computation for n MC steps, and then plug in the value function at the end. We can apply this to the DQN context by defining the target

$$y(s_0, a_0) = \sum_{t=1}^n \gamma^{t-1} r_t + \gamma^n \max_{a_n} Q_{\mathbf{w}}(s_n, a_n) \quad (2.55)$$

This can be implemented for episodic environments by storing experience tuples of the form

$$\tau = (s, a, \sum_{k=1}^n \gamma^{k-1} r_k, s_n, \text{done}) \quad (2.56)$$

where $\text{done} = 1$ if the trajectory ended at any point during the n -step rollout.

Theoretically this method is only valid if all the intermediate actions, $a_{2:n-1}$, are sampled from the current optimal policy derived from $Q_{\mathbf{w}}$, as opposed to some behavior policy, such as epsilon greedy or some samples from the replay buffer from an old policy. In practice, we can just restrict sampling to recent samples from the replay buffer, making the resulting method approximately on-policy.

2.5.4.5 $Q(\lambda)$

Instead of using a fixed n , it is possible to use a weighted combination of returns; this is known as the $Q(\lambda)$ algorithm [PW94; Har+16; Koz+21], and relies on the concept of eligibility traces. Unfortunately it is more complicated than the Sarsa case in Section 2.4.2, since Q learning is off-policy, but the eligibility traces backpropagate information obtained by the exploration policy.

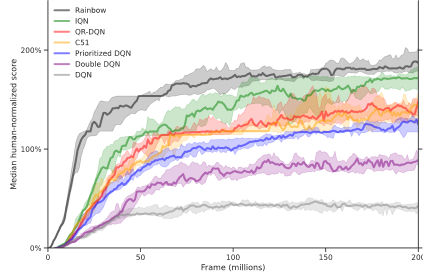


Figure 2.9: Plot of median human-normalized score over all 57 Atari games for various DQN agents. The yellow, red and green curves are distributional RL methods (Section 7.3), namely categorical DQN (C51) (Section 7.3.2) Quantile Regression DQN (Section 7.3.1), and Implicit Quantile Networks [Dab+18]. Figure from https://github.com/google-deepmind/dqn_zoo.

2.5.4.6 Rainbow

The **Rainbow** method of [Hes+18] combined 6 improvements to the vanilla DQN method, as listed below. (The paper is called “Rainbow” due to the color coding of their results plot, a modified version of which is shown in Figure 2.9.) At the time it was published (2018), this produced SOTA results on the Atari-200M benchmark. The 6 improvements are as follows:

- Use double DQN, as in Section 2.5.3.2.
- Use prioritized experience replay, as in Section 2.5.2.4.
- Use the categorical DQN (C51) (Section 7.3.2) distributional RL method.
- Use n -step returns (with $n = 3$), as in Section 2.5.4.4.
- Use dueling DQN, as in Section 2.5.4.2.
- Use noisy nets, as in Section 2.5.4.3.

Each improvement gives diminishing returns, as can be see in Figure 2.9.

Recently the “Beyond the Rainbow” paper [Unk24] proposed several more extensions:

- Use a larger CNN with residual connections, namely the Impala network from [Esp+18] with the modifications (including the use of spectral normalization) proposed in [SS21].
- Replace C51 with Implicit Quantile Networks [Dab+18].
- Use **Munchausen RL** [VPG20], which modifies the Q learning update rule by adding an entropy-like penalty.
- Collect 1 environment step from 64 parallel workers for each minibatch update (rather than taking many steps from a smaller number of workers).

2.5.4.7 Bigger, Better, Faster

At the time of writing this document (2024), the SOTA on the 100k sample-efficient Atari benchmark [Kai+19] is obtained by the **BBF** algorithm of [Sch+23b]. (BBF stands for “Bigger, Better, Faster”.) It uses the following tricks, in order of decreasing importance:

- Use a larger CNN with residual connections, namely a modified version of the Impala network from [Esp+18].

- Increase the **update-to-data** (UTD) ratio (number of times we update the Q function for every observation that is observed), in order to increase sample efficiency [HHA19].
- Use a periodic soft reset of (some of) the network weights to avoid loss of elasticity due to increased network updates, following the **SR-SPR** method of [D’O+22].
- Use n-step returns, as in Section 2.5.4.4, and then gradually decrease (anneal) the n-step return from $n = 10$ to $n = 3$, to reduce the bias over time.
- Add weight decay.
- Add a self-predictive representation loss (Section 4.4.2.6) to increase sample efficiency.
- Gradually increase the discount factor from $\gamma = 0.97$ to $\gamma = 0.997$, to encourage longer term planning once the model starts to be trained.⁷
- Drop noisy nets (which requires multiple network copies and thus slows down training due to increased memory use), since it does not help.
- Use dueling DQN (see Section 2.5.4.2).
- Use distributional DQN (see Section 7.3).

2.5.4.8 Other methods

Many other methods have been proposed to reduce the sample complexity of value-based RL while maintaining performance, see e.g., the **MEME** paper of [Kap+22].

2.5.5 Q-learning for GCRL using hindsight relabeling

In this section, we discuss how to learn a goal-conditioned policy (see Section 1.2.3) using Q learning (or any other off-policy learning method). The basic idea is as follows. We collect various trajectories in the environment, from a starting state s_0 to some terminal state s_T , and then define the goal of each trajectory as being $g = s_T$; this trajectory then serves as a demonstration of how to achieve this goal. This is called **hindsight experience relabeling** or **HER** [And+17], or just **hindsight relabeling**. This can be used to relabel the trajectories stored in the replay buffer. That is, if we have $(s, a, R(s|g), s', g)$ tuples, we replace them with $(s, a, R(s|g'), g')$ where $g' = s_T$. We can then use any off-policy RL method to learn $\pi(a|s, g)$.⁸

In [Eys+20], they show that HER can be viewed as a special case of maximum-entropy inverse RL (see Section 3.6.4), since it is estimating the reward for which the corresponding trajectory was optimal.

⁷The **Agent 57** method of [Bad+20] automatically learns the exploration rate and discount factor using a multi-armed bandit strategy, which lets it be more exploratory or more exploitative, depending on the game. This resulted in super human performance on all 57 Atari games in ALE. However, it required 80 billion frames (environment steps)! This was subsequently reduced to the “standard” 200M frames in the **MEME** method of [Kap+22].

⁸One limitation of HER is that it only works when the reward is Markovian, so $R(s, a) = 1$ iff $s = g$. A method to handle non-Markovian rewards using on-policy RL is discussed in [Gon+24]. The key trick is to use an auto-curriculum (see Section 7.4.2.2) where the goals are chosen to be solvable with the current policy.

Chapter 3

Policy-based RL

In the previous section, we considered methods that estimate the action-value function, $Q(s, a)$, from which we derive a policy. However, these methods have several disadvantages: (1) they can be difficult to apply to continuous action spaces; (2) they may diverge if function approximation is used (see Section 2.5.2.5); (3) the training of Q , often based on TD-style updates, is not directly related to the expected return garnered by the learned policy; (4) they learn deterministic policies, whereas in stochastic and partially observed environments, stochastic policies are provably better [JSJ94].

In this section, we discuss **policy search** methods, which directly optimize the parameters of the policy so as to maximize its expected return. We mostly focus on **policy gradient** methods, that use the gradient of the loss to guide the search (see e.g., [Aga+21a]). As we will see, these policy methods often benefit from estimating a value or advantage function to reduce the variance in the policy search process, so we will also use techniques from Chapter 2.

The parametric policy will be denoted by $\pi_{\theta}(a|s)$, which is usually some form of neural network. For discrete actions, the final layer is usually passed through a softmax function and then into a categorical distribution. For continuous actions, we typically use a Gaussian output layer (potentially clipped to a suitable range, such as $[-1, 1]$), although it is also possible to use more expressive (multi-modal) distributions, such as diffusion models (which, when used as a policy, is known as a **diffusion policy** [Chi+23]).

There are many implementation details one needs to get right to get good performance when designing such neural networks. For example, [Fur+21] recommends using ELU instead of RELU activations, and using LayerNorm. (In [Gal+24] they recently proved that adding layer norm to the final layer of a DQN model is sufficient to guarantee that value learning is stable, even in the nonlinear setting.) However, we do not discuss these details in this manuscript.

For more details on policy gradient methods, see e.g., [Wen18b; Aga+21a; Leh24].

3.1 Policy gradient methods

In this section, we discuss how to compute the expected value of a policy, and the gradient of this expectation. This can be used, together with SGD, to learn a locally optimal policy. Our presentation is based in part on [KWW22, Ch.11].

3.1.1 Likelihood ratio estimate

We define the value of a policy as

$$J(\pi_{\theta}) = J(\theta) = \mathbb{E}_{p_{\theta}(\tau)}[R(\tau)] \quad (3.1)$$

where $R(\boldsymbol{\tau}) = \gamma^0 r_0 + \gamma^1 r_1 + \dots$ is the return along the trajectory, and $p_{\boldsymbol{\theta}}(\boldsymbol{\tau})$ is the distribution over trajectories induced by the policy (and world model):

$$p_{\boldsymbol{\theta}}(\boldsymbol{\tau}) = p(s_1) \prod_{k=1}^T \mathcal{T}(s_{k+1}|s_k, a_k) \pi_{\boldsymbol{\theta}}(a_k|s_k) \quad (3.2)$$

The gradient of the policy value is given by

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \int \nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(\boldsymbol{\tau}) R(\boldsymbol{\tau}) d\boldsymbol{\tau} = \int p_{\boldsymbol{\theta}}(\boldsymbol{\tau}) \frac{\nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(\boldsymbol{\tau})}{p_{\boldsymbol{\theta}}(\boldsymbol{\tau})} R(\boldsymbol{\tau}) d\boldsymbol{\tau} = \mathbb{E}_{\boldsymbol{\tau}} \left[\frac{\nabla_{\boldsymbol{\theta}} p_{\boldsymbol{\theta}}(\boldsymbol{\tau})}{p_{\boldsymbol{\theta}}(\boldsymbol{\tau})} R(\boldsymbol{\tau}) \right] \quad (3.3)$$

This is known as the **likelihood ratio estimator**.

Now consider the **log derivative trick**, which is the simple observation that $\nabla \log \pi = \frac{\nabla \pi}{\pi}$. Using this, we can rewrite the above expression as

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\tau}} [\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{\tau}) R(\boldsymbol{\tau})] \quad (3.4)$$

The expectations can be estimated using Monte Carlo sampling (rolling out the policy in the environment). The gradient can be computed from Equation (3.2) as follows:

$$\nabla_{\boldsymbol{\theta}} \log p_{\boldsymbol{\theta}}(\boldsymbol{\tau}) = \sum_{k=1}^T \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a_k|s_k) \quad (3.5)$$

Hence

$$\nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\tau}} \left[\left(\sum_{k=1}^T \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a_k|s_k) \right) R(\boldsymbol{\tau}) \right] \quad (3.6)$$

In statistics, the term $\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}|\mathbf{s})$ is called the (Fisher) **score function**¹, so sometimes Equation (3.6) is called the **score function estimator** or **SFE** [Fu15; Moh+20].

3.1.2 Variance reduction using reward-to-go

The likelihood ratio estimator can have high variance, since we are sampling entire trajectories. Fortunately we can reduce the variance using the temporal/causal structure of the problem. In particular, note that from Equation (3.6) we have

$$\nabla J(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\tau}} \left[\left(\sum_{k=1}^T \underbrace{\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a_k|s_k)}_{f_k} \right) \left(\sum_{k=1}^T r_k \gamma^{k-1} \right) \right] \quad (3.7)$$

$$= \mathbb{E}_{\boldsymbol{\tau}} [(f_1 + f_2 + \dots + f_T)(r_1 + r_2 \gamma + r_3 \gamma^2 + \dots + r_T \gamma^{T-1})] \quad (3.8)$$

Expanding the product inside the expectation we have

$$f_1 r_1 + f_1 r_2 \gamma + f_1 r_3 \gamma^2 + \dots + f_1 r_T \gamma^{T-1} \quad (3.9)$$

$$+ \cancel{f_2 r_1} + f_2 r_2 \gamma + f_2 r_3 \gamma^2 + \dots + f_2 r_T \gamma^{T-1} \quad (3.10)$$

$$+ \cancel{f_3 r_1} + \cancel{f_2 r_2 \gamma} + f_3 r_3 \gamma^2 + \dots + f_3 r_T \gamma^{T-1} \quad (3.11)$$

$$\vdots \quad (3.12)$$

$$+ \cancel{f_T r_1} + \cancel{f_T r_2 \gamma} + \cancel{f_T r_3 \gamma^2} + \dots + f_T r_T \gamma^{T-1} \quad (3.13)$$

¹This is distinct from the Stein score, which is the gradient wrt the argument of the log probability, $\nabla_{\mathbf{a}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}|\mathbf{s})$, as used in diffusion.

where we have canceled terms that are 0, due to the fact that the reward at step k cannot depend on actions at time steps in the future. Plugging in this simplified expression we get

$$\nabla J(\theta) = \mathbb{E}_{\tau} \left[\sum_{k=1}^T \nabla_{\theta} \log \pi_{\theta}(a_k | s_k) \left(\sum_{l=k}^T r^l \gamma^{l-1} \right) \right] \quad (3.14)$$

$$= \mathbb{E}_{\tau} \left[\sum_{k=1}^T \nabla_{\theta} \log \pi_{\theta}(a_k | s_k) \left(\gamma^{k-1} \sum_{l=k}^T r^l \gamma^{l-k} \right) \right] \quad (3.15)$$

$$= \mathbb{E}_{\tau} \left[\sum_{k=1}^T \nabla_{\theta} \log \pi_{\theta}(a_k | s_k) \gamma^{k-1} G_k \right] \quad (3.16)$$

where G_k is the **reward-to-go** or return

$$G_k \triangleq r_k + \gamma r_{k+1} + \gamma^2 r_{k+2} + \dots + \gamma^{T-k-1} r_{T-1} = \sum_{l=k}^{T-1} \gamma^{l-k} r_l \quad (3.17)$$

Note that the reward-to-go of a state-action pair (s, a) can be considered as a single sample approximation of the state-action value function $Q_{\theta}(s, a)$. Averaging over such samples gives

$$\nabla J(\theta) = \mathbb{E}_{\tau} \left[\sum_{k=1}^T \gamma^{k-1} Q_{\theta}(s_k, a_k) \nabla_{\theta} \log \pi_{\theta}(a_k | s_k) \right] \quad (3.18)$$

3.1.3 REINFORCE

In this section, we describe an algorithm that uses the above estimate of the gradient of the policy value, together with SGD, to fit a policy. That is, we use

$$\theta_{j+1} := \theta_j + \eta \sum_{k=1}^T \nabla_{\theta} \log \pi_{\theta_j}(a_k | s_k) \gamma^{k-1} G_k \quad (3.19)$$

where j is the SGD iteration number, and we draw a single trajectory at each step. This is called the **REINFORCE** algorithm [Wil92].²

The update equation in Equation (3.19) can be interpreted as follows: we compute the sum of discounted future rewards induced by a trajectory, and if this is positive, we increase θ so as to make this trajectory more likely, otherwise we decrease θ . Thus, we reinforce good behaviors, and reduce the chances of generating bad ones. See Algorithm 3 for the pseudocode.

3.1.4 The policy gradient theorem

We now turn to the discounted infinite horizon setting. We define the **discounted state visitation measure** as follows:

$$\rho_{\pi}^{\gamma}(s) \triangleq \gamma^0 P(s_0 = s | \pi) + \gamma P(s_1 = s | \pi) + \gamma^2 P(s_2 = s | \pi) + \dots \quad (3.20)$$

$$= \sum_{t=0}^{\infty} \gamma^t \underbrace{\sum_{s_0} p_0(s_0) p^{\pi}(s_0 \rightarrow s, t)}_{p_t^{\pi}(s)} \quad (3.21)$$

²The term “REINFORCE” is an acronym for “REward Increment = nonnegative Factor x Offset Reinforcement x Characteristic Eligibility”. The phrase “characteristic eligibility” refers to the $\nabla \log \pi_{\theta}(a_t | s_t)$ term; the phrase “offset reinforcement” refers to the $G_t - b(s_t)$ term, where b is a baseline to be defined later; and the phrase “nonnegative factor” refers to the learning rate η of SGD.

Algorithm 3: REINFORCE (episodic version)

```
1 Initialize policy parameters  $\theta$ 
2 repeat
3   Sample an episode  $\tau = (s_1, a_1, r_1, s_2, \dots, s_T)$  using  $\pi_\theta$ 
4   for  $k = 1, \dots, T$  do
5      $G_k = \sum_{l=k}^T \gamma^{l-k} R_l$ 
6      $\theta \leftarrow \theta + \eta \rho_\theta \gamma^{k-1} G_k \nabla_\theta \log \pi_\theta(a_k | s_k)$ 
7 until converged;
```

where $p^\pi(s_0 \rightarrow s, t)$ is the probability of going from s_0 to s in t steps, and $p_t^\pi(s)$ is the marginal probability of being in state s at time t (after each episodic reset). Note that ρ_π^γ is a measure of time spent in non-terminal states, but it is not a probability measure, since it is not normalized, i.e., $\sum_s \rho_\pi^\gamma(s) \neq 1$. However, we can define a normalized version of the measure ρ by noting that $\sum_{t=0}^{\infty} \gamma^t = \frac{1}{1-\gamma}$ for $\gamma < 1$. Hence the **normalized discounted state visitation distribution** is given by the following (note the change from ρ to p):

$$p_\pi^\gamma(s) = (1 - \gamma) \rho_\pi^\gamma(s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t p_t^\pi(s) \quad (3.22)$$

We can convert from the normalized distribution back to the measure using

$$\rho_\pi^\gamma(s) = \frac{1}{1 - \gamma} p_\pi^\gamma(s) \quad (3.23)$$

Using this notation, one can show (see [KL02]) that we can rewrite Equation (3.18) in terms of expectations over states rather than over trajectories:

$$\nabla_\theta J(\theta) = \mathbb{E}_{\rho_\pi^\gamma(s) \pi_\theta(a|s)} [Q^{\pi_\theta}(s, a) \nabla_\theta \log \pi_\theta(a|s)] \quad (3.24)$$

$$= \frac{1}{1 - \gamma} \mathbb{E}_{p_\pi^\gamma(s) \pi_\theta(a|s)} [Q^{\pi_\theta}(s, a) \nabla_\theta \log \pi_\theta(a|s)] \quad (3.25)$$

This is known as the **policy gradient theorem** [Sut+99].

3.1.5 Variance reduction using a baseline

In practice, estimating the policy gradient using Equation (3.6) can have a high variance. A **baseline** function $b(s)$ can be used for variance reduction to get

$$\nabla_\theta J(\theta) = \mathbb{E}_{\rho_\theta(s) \pi_\theta(a|s)} [(Q_{\pi_\theta}(s, a) - b(s)) \nabla_\theta \log \pi_\theta(a|s)] \quad (3.26)$$

Any function that satisfies $\mathbb{E}[\nabla_\theta b(s)] = 0$ is a valid baseline. This follows since

$$\sum_a \nabla_\theta \pi_\theta(a|s) (Q(s, a) - b(s)) = \nabla_\theta \sum_a \pi_\theta(a|s) Q(s, a) - \nabla_\theta [\sum_a \pi_\theta(a|s)] b(s) \quad (3.27)$$

$$= \nabla_\theta \sum_a \pi_\theta(a|s) Q(s, a) - 0 \quad (3.28)$$

A common choice for the baseline is $b(s) = V(s)$. This is a valid choice since $E[\nabla_\theta V(s)] = 0$ if we use an old (frozen) version of the policy that is independent of θ . This is a useful choice $V(s)$ and $Q(s, a)$ are correlated and have similar magnitudes, so the scaling factor in front of the gradient term will be small, ensuring the update steps are not too big.

Note that $Q(s, a) - V(s) = A(s, a)$ is the advantage function. In the finite horizon case we get

$$\nabla J(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\tau}} \left[\sum_{k=1}^T \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a_k | s_k) \gamma^{k-1} (Q_{\boldsymbol{\theta}}(s_k, a_k) - V_{\boldsymbol{\theta}}(s_k)) \right] = \mathbb{E}_{\boldsymbol{\tau}} \left[\sum_{k=1}^T \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a_k | s_k) \gamma^{k-1} A_{\boldsymbol{\theta}}(s_k, a_k) \right] \quad (3.29)$$

We can also apply a baseline to the reward-to-go formulation to get

$$\nabla J(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{\tau}} \left[\sum_{k=1}^T \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a_k | s_k) \gamma^{k-1} (G_k - b(s_k)) \right] \quad (3.30)$$

We can derive analogous baselines for the infinite horizon case, defined in terms of p_{π}^{γ} .

3.1.6 REINFORCE with baseline

We can recover the full REINFORCE algorithm by combining SGD with the score function estimate with a baseline, as follows:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta \sum_{k=1}^T \gamma^{k-1} (G_k - b(s_k)) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a_k | s_k) \quad (3.31)$$

See Algorithm 4 for the pseudocode, where we use the value function as a baseline, estimated using TD.

Algorithm 4: REINFORCE (episodic) with value function baseline

```

1 Initialize policy parameters  $\boldsymbol{\theta}$ , baseline parameters  $\boldsymbol{w}$ 
2 repeat
3   Sample an episode  $\boldsymbol{\tau} = (s_1, a_1, r_1, s_2, \dots, s_T)$  using  $\pi_{\boldsymbol{\theta}}$ 
4   for  $k = 1, \dots, T$  do
5      $G_k = \sum_{l=k}^T \gamma^{l-k} R_l$ 
6      $\delta_k = G_k - V_{\boldsymbol{w}}(s_k)$ 
7      $\boldsymbol{w} \leftarrow \boldsymbol{w} - \eta_{\boldsymbol{w}} \delta_k \nabla_{\boldsymbol{w}} V_{\boldsymbol{w}}(s_k)$ 
8      $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta_{\boldsymbol{\theta}} \gamma^{k-1} \delta_k \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a_k | s_k)$ 
9 until converged;
```

3.2 Actor-critic methods

An **actor-critic** method [BSA83] uses the policy gradient method, but where the expected return G_t is estimated using temporal difference learning of a value function instead of MC rollouts. (The term “actor” refers to the policy, and the term “critic” refers to the value function.) The use of bootstrapping in TD updates allows more efficient learning of the value function compared to MC, and further reduces the variance. In addition, it allows us to develop a fully online, incremental algorithm, that does not need to wait until the end of the trajectory before updating the parameters.

3.2.1 Advantage actor critic (A2C)

Concretely, consider the use of the one-step TD method to estimate the return in the episodic case, i.e., we replace G_t with $G_{t:t+1} = r_t + \gamma V_{\boldsymbol{w}}(s_{t+1})$. If we use $V_{\boldsymbol{w}}(s_t)$ as a baseline, the REINFORCE update in

Equation (3.19) becomes

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta \sum_{t=0}^{T-1} \gamma^t (G_{t:t+1} - V_{\mathbf{w}}(s_t)) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a_t|s_t) \quad (3.32)$$

$$= \boldsymbol{\theta} + \eta \sum_{t=0}^{T-1} \gamma^t (r_t + \gamma V_{\mathbf{w}}(s_{t+1}) - V_{\mathbf{w}}(s_t)) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a_t|s_t) \quad (3.33)$$

Note that $\delta_t = r_{t+1} + \gamma V_{\mathbf{w}}(s_{t+1}) - V_{\mathbf{w}}(s_t)$ is a single sample approximation to the advantage function $\text{Adv}(s_t, a_t) = Q(s_t, a_t) - V(s_t)$. This method is therefore called **advantage actor critic** or **A2C**. See Algorithm 5 for the pseudo-code.³ (Note that $V_{\mathbf{w}}(s_{t+1}) = 0$ if s_t is a done state, representing the end of an episode.) Note that this is an on-policy algorithm, where we update the value function $V_{\mathbf{w}}^{\pi}$ to reflect the value of the current policy π . See Section 3.2.3 for further discussion of this point.

Algorithm 5: Advantage actor critic (A2C) algorithm (episodic)

```

1 Initialize actor parameters  $\boldsymbol{\theta}$ , critic parameters  $\mathbf{w}$ 
2 repeat
3   Sample starting state  $s_0$  of a new episode
4   for  $t = 0, 1, 2, \dots$  do
5     Sample action  $a_t \sim \pi_{\boldsymbol{\theta}}(\cdot|s_t)$ 
6      $(s_{t+1}, r_t, \text{done}_t) = \text{env.step}(s_t, a_t)$ 
7      $y_t = r_t + \gamma(1 - \text{done}_t)V_{\mathbf{w}}(s_{t+1})$  // Target
8      $\delta_t = y_t - V_{\mathbf{w}}(s_t)$  // Advantage
9      $\mathbf{w} \leftarrow \mathbf{w} + \eta_{\mathbf{w}} \delta_t \nabla_{\mathbf{w}} V_{\mathbf{w}}(s_t)$  // Critic
10     $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta_{\boldsymbol{\theta}} \gamma^t \delta_t \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(a_t|s_t)$  // Actor
11    if  $\text{done}_t = 1$  then
12      break
13 until converged;
```

In practice, we should use a stop-gradient operator on the target value for the TD update, for reasons explained in Section 2.5.2.5. Furthermore, it is common to add an entropy term to the policy, to act as a regularizer (to ensure the policy remains stochastic, which smoothens the loss function — see Section 3.6.8). If we use a shared network with separate value and policy heads, we need to use a single loss function for training all the parameters ϕ . Thus we get the following loss, for each trajectory, where we want to minimize TD loss, maximize the policy gradient (expected reward) term, and maximize the entropy term.

$$\mathcal{L}(\phi; \tau) = \frac{1}{T} \sum_{t=1}^T [\lambda_{TD} \mathcal{L}_{TD}(s_t, a_t, r_t, s_{t+1}) - \lambda_{PG} J_{PG}(s_t, a_t, r_t, s_{t+1}) - \lambda_{ent} J_{ent}(s_t)] \quad (3.34)$$

$$y_t = r_t + \gamma(1 - \text{done}(s_t))V_{\phi}(s_{t+1}) \quad (3.35)$$

$$\mathcal{L}_{TD}(s_t, a_t, r_t, s_{t+1}) = (\text{sg}(y_t) - V_{\phi}(s_t))^2 \quad (3.36)$$

$$J_{PG}(s_t, a_t, r_t, s_{t+1}) = (\text{sg}(y_t - V_{\phi}(s_t)) \log \pi_{\phi}(a_t|s_t)) \quad (3.37)$$

$$J_{ent}(s_t) = - \sum_a \pi_{\phi}(a|s_t) \log \pi_{\phi}(a|s_t) \quad (3.38)$$

To handle the dynamically varying scales of the different loss functions, we can use the **PopArt** method of [Has+16; Hes+19] to allow for a fixed set of hyper-parameter values for λ_i . (PopArt stands for “Preserving Outputs Precisely, while Adaptively Rescaling Targets”.)

³In [Mni+16], they proposed a distributed version of A2C known as **A3C** which stands for “asynchronous advantage actor critic”.

3.2.2 Generalized advantage estimation (GAE)

In A2C, we replaced the high variance, but unbiased, MC return G_t with the low variance, but biased, one-step bootstrap return $G_{t:t+1} = r_t + \gamma V_{\mathbf{w}}(s_{t+1})$. More generally, we can compute the n -step estimate

$$G_{t:t+n} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots + \gamma^{n-1} r_{t+n-1} + \gamma^n V_{\mathbf{w}}(s_{t+n}) \quad (3.39)$$

and thus obtain the (truncated) n -step advantage estimate as follows:

$$A_{\mathbf{w}}^{(n)}(s_t, a_t) = G_{t:t+n} - V_{\mathbf{w}}(s_t) \quad (3.40)$$

Unrolling to infinity, we get

$$A_t^{(1)} = r_t + \gamma v_{t+1} - v_t \quad (3.41)$$

$$A_t^{(2)} = r_t + \gamma r_{t+1} + \gamma^2 v_{t+2} - v_t \quad (3.42)$$

$$\vdots \quad (3.43)$$

$$A_t^{(\infty)} = r_t + \gamma r_{t+1} + \gamma^2 r_{t+2} + \dots - v_t \quad (3.44)$$

$A_t^{(1)}$ is high bias but low variance, and $A_t^{(\infty)}$ is unbiased but high variance.

Instead of using a single value of n , we can take a weighted average. That is, we define

$$A_t = \frac{\sum_{n=1}^T w_n A_t^{(n)}}{\sum_{n=1}^T w_n} \quad (3.45)$$

If we set $w_n = \lambda^{n-1}$ we get the following simple recursive calculation:

$$\delta_t = r_t + \gamma v_{t+1} - v_t \quad (3.46)$$

$$A_t = \delta_t + \gamma \lambda \delta_{t+1} + \dots + (\gamma \lambda)^{T-(t+1)} \delta_{T-1} = \delta_t + \gamma \lambda A_{t+1} \quad (3.47)$$

Here $\lambda \in [0, 1]$ is a parameter that controls the bias-variance tradeoff: larger values decrease the bias but increase the variance. This is called **generalized advantage estimation (GAE)** [Sch+16b]. See Algorithm 6 for some pseudocode. Using this, we can define a general actor-critic method, as shown in Algorithm 7.

Algorithm 6: Generalized Advantage Estimation

```

1 def GAE( $r_{1:T}, v_{1:T}, \gamma, \lambda$ )
2    $A_T = 0$ 
3   for  $t = T - 1 : 1$  do
4      $\delta_t = r_t + \gamma v_{t+1} - v_t$ 
5      $A_t = \delta_t + \gamma \lambda A_{t+1}$  // advantage
6    $y_t = A_t + v_t$  // Return
7 Return ( $A_{1:T}, y_{1:T}$ )
```

We can generalize this approach even further, by using gradient estimators of the form

$$\nabla J(\boldsymbol{\theta}) = \mathbb{E} \left[\sum_{t=0}^{\infty} \Psi_t \nabla \log \pi_{\boldsymbol{\theta}}(a_t | s_t) \right] \quad (3.48)$$

Algorithm 7: Actor critic with GAE

```
1 Initialize parameters  $\phi$ , environment state  $s$ 
2 repeat
3    $(s_1, a_1, r_1, \dots, s_T) = \text{rollout}(s, \pi_\phi)$ 
4    $v_{1:T} = V_\phi(s_{1:T})$ 
5    $(A_{1:T}, y_{1:T}) = \text{sg}(\text{GAE}(r_{1:T}, v_{1:T}, \gamma, \lambda))$ 
6    $\mathcal{L}(\phi) = \frac{1}{T} \sum_{t=1}^T [\lambda_{TD}(V_\phi(s_t) - y_t)^2 - \lambda_{PG} A_t \log \pi_\phi(a_t|s_t) - \lambda_{ent} \mathbb{H}(\pi_\phi(\cdot|s_t))]$ 
7    $\phi := \phi - \eta \nabla \mathcal{L}(\phi)$ 
8 until converged;
```

where Ψ_t may be any of the following:

$$\Psi_t = \sum_{i=t}^{\infty} \gamma^i r_i \quad \text{Monte Carlo target} \quad (3.49)$$

$$\Psi_t = \sum_{i=t}^{\infty} \gamma^i r_i - V_{\mathbf{w}}(s_t) \quad \text{MC with baseline} \quad (3.50)$$

$$\Psi_t = A_{\mathbf{w}}(s_t, a_t) \quad \text{advantage function} \quad (3.51)$$

$$\Psi_t = Q_{\mathbf{w}}(s_t, a_t) \quad \text{Q function} \quad (3.52)$$

$$\Psi_t = r_t + \gamma V_{\mathbf{w}}(s_{t+1}) - V_{\mathbf{w}}(s_t) \quad \text{TD residual} \quad (3.53)$$

See [Sch+16b] for details.

3.2.3 Two-time scale actor critic algorithms

In standard AC, we update the actor and critic in parallel. However, it is better to let critic $V_{\mathbf{w}}$ learn using a faster learning rate (or more updates), so that it reflects the value of the current policy π_{θ} more accurately, in order to get better gradient estimates for the policy update. This is known as two timescale learning or **bilevel optimization** [Yu17; Zha+19; Hon+23; Zhe+22a; Lor24]. (See also Section 4.3.1, where we discuss RL from a game theoretic perspective.)

An alternative approach is to alternate between updating the policy and the value function, rather than updating them simultaneously. This is called **phasic policy gradient** [Cob+20].

3.2.4 Natural policy gradient methods

In this section, we discuss an improvement to policy gradient methods that uses preconditioning to speedup convergence. In particular, we replace gradient descent with **natural gradient descent (NGD)** [Ama98; Mar20], which we explain below. We then show how to combine it with actor-critic.

3.2.4.1 Natural gradient descent

NGD is a second order method for optimizing the parameters of (conditional) probability distributions, such as policies, $\pi_{\theta}(\mathbf{a}|\mathbf{s})$. It typically converges faster and more robustly than SGD, but is computationally more expensive.

Before we explain NGD, let us review standard SGD, which is an update of the following form

$$\theta_{k+1} = \theta_k - \eta_k \mathbf{g}_k \quad (3.54)$$

where $\mathbf{g}_k = \nabla_{\theta} \mathcal{L}(\theta_k)$ is the gradient of the loss at the previous parameter values, and η_k is the learning rate. It can be shown that the above update is equivalent to minimizing a locally linear approximation to the loss,



Figure 3.1: Changing the mean of a Gaussian by a fixed amount (from solid to dotted curve) can have more impact when the (shared) variance is small (as in a) compared to when the variance is large (as in b). Hence the impact (in terms of prediction accuracy) of a change to μ depends on where the optimizer is in (μ, σ) space. From Figure 3 of [Hon+10], reproduced from [Val00]. Used with kind permission of Antti Honkela.

$\hat{\mathcal{L}}_k$, subject to the constraint that the new parameters do not move too far (in Euclidean distance) from the previous parameters:

$$\theta_{k+1} = \underset{\theta}{\operatorname{argmin}} \hat{\mathcal{L}}_k(\theta) \quad \text{s.t.} \quad \|\theta - \theta_k\|_2^2 \leq \epsilon \quad (3.55)$$

$$\hat{\mathcal{L}}_k(\theta) = \mathcal{L}(\theta_k) + \mathbf{g}_k^\top (\theta - \theta_k) \quad (3.56)$$

where the step size η_k is proportional to ϵ . This is called a **proximal update** [PB+14].

One problem with the SGD update is that Euclidean distance in parameter space does not make sense for probabilistic models. For example, consider comparing two Gaussians, $p_\theta = p(y|\mu, \sigma)$ and $p_{\theta'} = p(y|\mu', \sigma')$. The (squared) Euclidean distance between the parameter vectors decomposes as $\|\theta - \theta'\|_2^2 = (\mu - \mu')^2 + (\sigma - \sigma')^2$. However, the predictive distribution has the form $\exp(-\frac{1}{2\sigma^2}(y - \mu)^2)$, so changes in μ need to be measured relative to σ . This is illustrated in Figure 3.1(a-b), which shows two univariate Gaussian distributions (dotted and solid lines) whose means differ by ϵ . In Figure 3.1(a), they share the same small variance σ^2 , whereas in Figure 3.1(b), they share the same large variance. It is clear that the difference in μ matters much more (in terms of the effect on the distribution) when the variance is small. Thus we see that the two parameters interact with each other, which the Euclidean distance cannot capture.

The key to NGD is to measure the notion of distance between two probability distributions in terms of the KL divergence. This can be approximated in terms of the **Fisher information matrix** (FIM). In particular, for any given input \mathbf{x} , we have

$$D_{\text{KL}}(p_\theta(\mathbf{y}|\mathbf{x}) \parallel p_{\theta+\delta}(\mathbf{y}|\mathbf{x})) \approx \frac{1}{2} \delta^\top \mathbf{F}_\mathbf{x} \delta \quad (3.57)$$

where $\mathbf{F}_\mathbf{x}$ is the FIM

$$\mathbf{F}_\mathbf{x}(\theta) = -\mathbb{E}_{p_\theta(\mathbf{y}|\mathbf{x})} [\nabla^2 \log p_\theta(\mathbf{y}|\mathbf{x})] = \mathbb{E}_{p_\theta(\mathbf{y}|\mathbf{x})} [(\nabla \log p_\theta(\mathbf{y}|\mathbf{x}))(\nabla \log p_\theta(\mathbf{y}|\mathbf{x}))^\top] \quad (3.58)$$

We now replace the Euclidean distance between the parameters, $d(\theta_k, \theta_{k+1}) = \|\delta\|_2^2$, with

$$d(\theta_k, \theta_{k+1}) = \delta^\top \mathbf{F}_k \delta \quad (3.59)$$

where $\delta = \theta_{k+1} - \theta_k$ and $\mathbf{F}_k = \mathbf{F}_\mathbf{x}(\theta_k)$ for a randomly chosen input \mathbf{x} . This gives rise to the following constrained optimization problem:

$$\delta_k = \underset{\delta}{\operatorname{argmin}} \hat{\mathcal{L}}_k(\theta_k + \delta) \quad \text{s.t.} \quad \delta^\top \mathbf{F}_k \delta \leq \epsilon \quad (3.60)$$

If we replace the constraint with a Lagrange multiplier, we get the unconstrained objective:

$$J_k(\delta) = \mathcal{L}(\theta_k) + \mathbf{g}_k^\top \delta + \eta_k \delta^\top \mathbf{F}_k \delta \quad (3.61)$$

Solving $J_k(\boldsymbol{\delta}) = 0$ gives the update

$$\boldsymbol{\delta} = -\eta_k \mathbf{F}_k^{-1} \mathbf{g}_k \quad (3.62)$$

The term $\mathbf{F}^{-1} \mathbf{g}$ is called the **natural gradient**. This is equivalent to a preconditioned gradient update, where we use the inverse FIM as a preconditioning matrix. We can compute the (adaptive) learning rate using

$$\eta_k = \sqrt{\frac{\epsilon}{\mathbf{g}_k^\top \mathbf{F}_k^{-1} \mathbf{g}_k}} \quad (3.63)$$

Computing the FIM can be hard. A simple approximation is to replace the model’s distribution with the empirical distribution. In particular, define $p_{\mathcal{D}}(\mathbf{x}, \mathbf{y}) = \frac{1}{N} \sum_{n=1}^N \delta_{\mathbf{x}_n}(\mathbf{x}) \delta_{\mathbf{y}_n}(\mathbf{y})$, $p_{\mathcal{D}}(\mathbf{x}) = \frac{1}{N} \sum_{n=1}^N \delta_{\mathbf{x}_n}(\mathbf{x})$ and $p_{\theta}(\mathbf{x}, \mathbf{y}) = p_{\mathcal{D}}(\mathbf{x}) p(\mathbf{y}|\mathbf{x}, \theta)$. Then we can compute the **empirical Fisher** [Mar16] as follows:

$$\mathbf{F}(\theta) = \mathbb{E}_{p_{\theta}(\mathbf{x}, \mathbf{y})} [\nabla \log p(\mathbf{y}|\mathbf{x}, \theta) \nabla \log p(\mathbf{y}|\mathbf{x}, \theta)^\top] \quad (3.64)$$

$$\approx \mathbb{E}_{p_{\mathcal{D}}(\mathbf{x}, \mathbf{y})} [\nabla \log p(\mathbf{y}|\mathbf{x}, \theta) \nabla \log p(\mathbf{y}|\mathbf{x}, \theta)^\top] \quad (3.65)$$

$$= \frac{1}{|\mathcal{D}|} \sum_{(\mathbf{x}, \mathbf{y}) \in \mathcal{D}} \nabla \log p(\mathbf{y}|\mathbf{x}, \theta) \nabla \log p(\mathbf{y}|\mathbf{x}, \theta)^\top \quad (3.66)$$

3.2.4.2 Natural actor critic

To apply NGD to RL, we can adapt the A2C algorithm in Algorithm 7. In particular, define

$$\mathbf{g}_{kt} = \nabla_{\theta_k} A_t \log \pi_{\theta}(\mathbf{a}_t | \mathbf{s}_t) \quad (3.67)$$

where A_t is the advantage function at step t of the random trajectory generated by the policy at iteration k . Now we compute

$$\mathbf{g}_k = \frac{1}{T} \sum_{t=1}^T \mathbf{g}_{kt}, \quad \mathbf{F}_k = \frac{1}{T} \sum_{t=1}^T \mathbf{g}_{kt} \mathbf{g}_{kt}^\top \quad (3.68)$$

and compute $\boldsymbol{\delta}_{k+1} = -\eta_k \mathbf{F}_k^{-1} \mathbf{g}_k$. This approach is called **natural policy gradient** [Kak01; Raj+17].

We can compute $\mathbf{F}_k^{-1} \mathbf{g}_k$ without having to invert \mathbf{F}_k by using the conjugate gradient method, where each CG step uses efficient methods for Hessian-vector products [Pea94]. This is called **Hessian free optimization** [Mar10]. Similarly, we can efficiently compute $\mathbf{g}_k^\top (\mathbf{F}_k^{-1} \mathbf{g}_k)$.

As a more accurate alternative to the empirical Fisher, [MG15] propose the **KFAC** method, which stands for “Kronecker factored approximate curvature”; this approximates the FIM of a DNN as a block diagonal matrix, where each block is a Kronecker product of two small matrices. This was applied to policy gradient learning in [Wu+17].

3.2.5 Architectural issues

It is common to use a single neural network for both the actor and critic, but using different output heads: a scalar output for the value function, and a vector output for the policy. For example, the **Amago** method [GFZ24; Gri+24] uses a transformer backbone. To train the shared model, they construct a unified objective $\mathcal{L} = E[\lambda_0 \mathcal{L}_{TD} + \lambda_1 \mathcal{L}_{PG}]$, where the TD and policy gradient losses are dynamically normalized using the **PopArt** method of [Has+16; Hes+19] to allow for a fixed set of hyper-parameter values for λ_i , even as the range of the losses change over time. (PopArt stands for “Preserving Outputs Precisely, while Adaptively Rescaling Targets”).

However, [Gar+25] argue that it can be better to use different networks for the actor and critic (at least when using MLPs/CNNs), since they need to extract different kinds of features.

3.2.6 Deterministic policy gradient methods

In this section, we consider an actor critic method that uses a deterministic policy, that predicts a unique action for each state, so $a_t = \mu_{\theta}(s_t)$, rather than $a_t \sim \pi_{\theta}(s_t)$. This is trained to match the optimal action from $Q_w(s, a)$. Thus we can think of the resulting method as a version of DQN designed for continuous actions. (We require that the actions are continuous, because we will take the Jacobian of the Q function wrt the actions.)

The benefit of using a deterministic policy, as opposed to a stochastic policy, is that we can modify the policy gradient method so that it can work off policy (e.g., using replay buffers) without needing importance sampling, in contrast to the methods discussed in Section 3.4. In addition, the feedback signal for learning is based on the vector-valued gradient of the value function, which is more informative than a scalar reward signal.

3.2.6.1 Deterministic policy gradient theorem

As before, we define the value of a policy as the expected discounted reward per state:

$$J(\mu_{\theta}) \triangleq \mathbb{E}_{\rho_{\mu_{\theta}}(s)} [R(s, \mu_{\theta}(s))] \quad (3.69)$$

The **deterministic policy gradient theorem** [Sil+14] tells us that the gradient of this expression is given by

$$\nabla_{\theta} J(\mu_{\theta}) = \mathbb{E}_{\rho_{\mu_{\theta}}(s)} [\nabla_{\theta} Q_{\mu_{\theta}}(s, \mu_{\theta}(s))] \quad (3.70)$$

$$= \mathbb{E}_{\rho_{\mu_{\theta}}(s)} [\nabla_{\theta} \mu_{\theta}(s) \nabla_a Q_{\mu_{\theta}}(s, a)|_{a=\mu_{\theta}(s)}] \quad (3.71)$$

where $\nabla_{\theta} \mu_{\theta}(s)$ is the $N_{\theta} \times N_A$ Jacobian matrix, and N_A and N_{θ} are the dimensions of \mathcal{A} and θ , respectively. The intuition for this equation is as follows: the change in the expected value due to changing the parameters, $\nabla_{\theta} J(\mu_{\theta}) \in \mathbb{R}^{N_{\theta}}$, is equal to the change in the policy output (i.e., the actions) due to changing the parameters, $\nabla_{\theta} \mu_{\theta}(s) \in \mathbb{R}^{N_{\theta} \times N_A}$ times the change in the expected value due to the change in the actions, $\nabla_a Q_{\mu_{\theta}}(s, a) \in \mathbb{R}^{N_A}$.

For stochastic policies of the form $\pi_{\theta}(a|s) = \mu_{\theta}(s) + \text{noise}$, the standard policy gradient theorem reduces to the above form as the noise level goes to zero.

Note that the gradient estimate in Equation (3.71) integrates over the states but not over the actions, which helps reduce the variance in gradient estimation from sampled trajectories. However, since the deterministic policy does not do any exploration, we need to use an off-policy method for training. This collects data from a stochastic behavior policy π_b , whose stationary state distribution is $p_{\pi_b}^{\gamma}$. The original objective, $J(\mu_{\theta})$, is approximated by the following:

$$J_b(\mu_{\theta}) \triangleq \mathbb{E}_{p_{\pi_b}^{\gamma}(s)} [V_{\mu_{\theta}}(s)] = \mathbb{E}_{p_{\pi_b}^{\gamma}(s)} [Q_{\mu_{\theta}}(s, \mu_{\theta}(s))] \quad (3.72)$$

with the off-policy deterministic policy gradient from [DWS12] is approximated by

$$\nabla_{\theta} J_b(\mu_{\theta}) \approx \mathbb{E}_{p_{\pi_b}^{\gamma}(s)} [\nabla_{\theta} [Q_{\mu_{\theta}}(s, \mu_{\theta}(s))]] = \mathbb{E}_{p_{\pi_b}^{\gamma}(s)} [\nabla_{\theta} \mu_{\theta}(s) \nabla_a Q_{\mu_{\theta}}(s, a)|_{a=\mu_{\theta}(s)}] \quad (3.73)$$

where we have a dropped a term that depends on $\nabla_{\theta} Q_{\mu_{\theta}}(s, a)$ and is hard to estimate [Sil+14].

To apply Equation (3.73), we may learn $Q_w \approx Q_{\mu_{\theta}}$ with TD, giving rise to the following updates:

$$\delta = r_t + \gamma Q_w(s_{t+1}, \mu_{\theta}(s_{t+1})) - Q_w(s_t, a_t) \quad (3.74)$$

$$w_{t+1} \leftarrow w_t + \eta_w \delta \nabla_w Q_w(s_t, a_t) \quad (3.75)$$

$$\theta_{t+1} \leftarrow \theta_t + \eta_{\theta} \nabla_{\theta} \mu_{\theta}(s_t) \nabla_a Q_w(s_t, a)|_{a=\mu_{\theta}(s_t)} \quad (3.76)$$

So we learn both a state-action critic Q_w and an actor μ_{θ} . This method avoids importance sampling in the actor update because of the deterministic policy gradient, and we avoids it in the critic update because of the use of Q-learning.

If $Q_{\mathbf{w}}$ is linear in \mathbf{w} , and uses features of the form $\phi(s, a) = \mathbf{a}^\top \nabla_{\boldsymbol{\theta}} \mu_{\boldsymbol{\theta}}(s)$, then we say the function approximator for the critic is **compatible** with the actor; in this case, one can show that the above approximation does not bias the overall gradient.

The basic off-policy DPG method has been extended in various ways, some of which we describe below.

3.2.6.2 DDPG

The **DDPG** algorithm of [Lil+16], which stands for “deep deterministic policy gradient”, uses the DQN method (Section 2.5.2.2) to learn the Q function, and then uses this to evaluate the policy. In more detail, the actor tries to minimize the output of the critic

$$\mathcal{L}_{\boldsymbol{\theta}}(s) = Q_{\mathbf{w}}(s, \mu_{\boldsymbol{\theta}}(s)) \quad (3.77)$$

where the loss is averaged over states s drawn from the replay buffer. The critic tries to minimize the 1-step TD loss, as in Q-learning:

$$\mathcal{L}_{\mathbf{w}}(s, a, r, s') = [Q_{\mathbf{w}}(s, a) - (r + \gamma Q_{\bar{\mathbf{w}}}(s', \mu_{\boldsymbol{\theta}}(s')))]^2 \quad (3.78)$$

where $Q_{\bar{\mathbf{w}}}$ is the target critic network, and the samples (s, a, r, s') are drawn from a replay buffer. (See Section 2.5.2.6 for a discussion of target networks.)

The **D4PG** algorithm [BM+18], which stands for “distributed distributional DDPG”, extends DDPG to handle distributed training, and to handle distributional RL (see Section 7.3).

3.2.6.3 Twin Delayed DDPG (TD3)

The **TD3** (“twin delayed deep deterministic”) method of [FHM18] extends DDPG in 3 main ways. First, it uses **target policy smoothing**, in which noise is added to the action, to encourage generalization:

$$\tilde{\mathbf{a}} = \mu_{\boldsymbol{\theta}}(s) + \text{noise} = \pi_{\boldsymbol{\theta}}(s) \quad (3.79)$$

Second it uses **clipped double Q learning**, which is an extension of the double Q-learning discussed in Section 2.5.3.1 to avoid over-estimation bias. In particular, the target values for TD learning are defined using

$$y(r, s'; \bar{\mathbf{w}}_{1:2}, \bar{\boldsymbol{\theta}}) = r + \gamma \min_{i=1,2} Q_{\bar{\mathbf{w}}_i}(s', \pi_{\bar{\boldsymbol{\theta}}}(s')) \quad (3.80)$$

Third, it uses **delayed policy updates**, in which it only updates the policy after the value function has stabilized. (See also Section 3.2.3.) See Algorithm 8 for the pseudocode.

3.2.6.4 Wasserstein Policy Optimization (WPO)

As we noted above, one advantage of DPG-based methods is that they can use the gradient of the value with respect to actions. In [Pfa+25] it was shown that by approximating **Wasserstein gradient flows** over the space of all parametric policies, we arrive at an update very similar to DPG, but for general stochastic policies.⁴ The derivation is somewhat complex, but the final algorithm is quite simple. In particular, we should update the policy using the following:⁵

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \mathcal{F}^{-1} \mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \pi_{\boldsymbol{\theta}}(\cdot|\mathbf{s})} [\nabla_{\boldsymbol{\theta}} (\nabla_{\mathbf{a}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}|\mathbf{s})^\top) \nabla_{\mathbf{a}} Q^{\pi_{\boldsymbol{\theta}}}(\mathbf{s}, \mathbf{a})] \quad (3.81)$$

where \mathcal{F} is the Fisher information matrix (FIM):

$$\mathcal{F} = \mathbb{E}_{\mathbf{s} \sim \mathcal{D}, \mathbf{a} \sim \pi_{\boldsymbol{\theta}}(\cdot|\mathbf{s})} [\nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}|\mathbf{s}) \nabla_{\boldsymbol{\theta}} \log \pi_{\boldsymbol{\theta}}(\mathbf{a}|\mathbf{s})^\top] \quad (3.82)$$

⁴Although the method of [Zha+18] and the **SVG(0)** method of [Hee+15] also support stochastic policies, they rely on the reparameterization trick, which is not always applicable (e.g., if the policy is a mixture of Gaussians). In addition, WPO makes use of natural gradients, whereas these are first-order methods.

⁵Note that $f(\mathbf{a}, \boldsymbol{\theta}) = \log \pi_{\boldsymbol{\theta}}(\mathbf{a}|\mathbf{s})$ is a scalar-valued function of $\boldsymbol{\theta}$ and \mathbf{a} (for a fixed \mathbf{s}); the notation $\nabla_{\boldsymbol{\theta}}(\nabla_{\mathbf{a}} f(\mathbf{a}, \boldsymbol{\theta})^\top)$ is another way of writing the Jacobian matrix $[\frac{\partial f}{\partial \theta_i \partial a_j}]_{ij}$.

Algorithm 8: TD3

```
1 Initialize environment state  $\mathbf{s}$ , policy parameters  $\boldsymbol{\theta}$ , target policy parameters  $\bar{\boldsymbol{\theta}}$ , critic parameters  $\mathbf{w}_i$ ,  
   target critic parameters  $\bar{\mathbf{w}}_i = \mathbf{w}_i$ , replay buffer  $\mathcal{D} = \emptyset$ , discount factor  $\gamma$ , EMA rate  $\rho$ , step size  $\eta_{\mathbf{w}}$ ,  
    $\eta_{\boldsymbol{\theta}}$ .  
2 repeat  
3    $\mathbf{a} = \mu_{\boldsymbol{\theta}}(\mathbf{s}) + \text{noise}$   
4    $(\mathbf{s}', r) = \text{step}(\mathbf{a}, \mathbf{s})$   
5    $\mathcal{D} := \mathcal{D} \cup \{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}')\}$   
6    $\mathbf{s} \leftarrow \mathbf{s}'$   
7   for  $G$  updates do  
8     Sample a minibatch  $\mathcal{B} = \{(\mathbf{s}_j, \mathbf{a}_j, r_j, \mathbf{s}'_j)\}$  from  $\mathcal{D}$   
9      $\mathbf{w} = \text{update-critics}(\boldsymbol{\theta}, \mathbf{w}, \mathcal{B})$   
10    Sample a minibatch  $\mathcal{B} = \{(\mathbf{s}_j, \mathbf{a}_j, r_j, \mathbf{s}'_j)\}$  from  $\mathcal{D}$   
11     $\boldsymbol{\theta} = \text{update-policy}(\boldsymbol{\theta}, \mathbf{w}, \mathcal{B})$   
12 until converged;  
13 .  
14 def  $\text{update-critics}(\boldsymbol{\theta}, \mathbf{w}, \mathcal{B})$ :  
15   Let  $(\mathbf{s}_j, \mathbf{a}_j, r_j, \mathbf{s}'_j)_{j=1}^B = \mathcal{B}$   
16   for  $j = 1 : B$  do  
17      $\tilde{\mathbf{a}}_j = \mu_{\bar{\boldsymbol{\theta}}}(\mathbf{s}'_j) + \text{clip}(\text{noise}, -c, c)$   
18      $y_j = r_j + \gamma \min_{i=1,2} Q_{\bar{\mathbf{w}}_i}(\mathbf{s}'_j, \tilde{\mathbf{a}}_j)$   
19   for  $i = 1 : 2$  do  
20      $\mathcal{L}(\mathbf{w}_i) = \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}')_j \in \mathcal{B}} (Q_{\mathbf{w}_i}(\mathbf{s}_j, \mathbf{a}_j) - \text{sg}(y_j))^2$   
21      $\mathbf{w}_i \leftarrow \mathbf{w}_i - \eta_{\mathbf{w}} \nabla \mathcal{L}(\mathbf{w}_i)$  // Descent  
22      $\bar{\mathbf{w}}_i := \rho \bar{\mathbf{w}}_i + (1 - \rho) \mathbf{w}_i$  // Update target networks with EMA  
23 Return  $\mathbf{w}_{1:N}, \bar{\mathbf{w}}_{1:N}$   
24 .  
25 def  $\text{update-actor}(\boldsymbol{\theta}, \mathbf{w}, \mathcal{B})$ :  
26    $J(\boldsymbol{\theta}) = \frac{1}{|\mathcal{B}|} \sum_{\mathbf{s} \in \mathcal{B}} (Q_{\mathbf{w}_1}(\mathbf{s}, \mu_{\boldsymbol{\theta}}(\mathbf{s})))^2$   
27    $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta_{\boldsymbol{\theta}} \nabla J(\boldsymbol{\theta})$  // Ascent  
28    $\bar{\boldsymbol{\theta}} := \rho \bar{\boldsymbol{\theta}} + (1 - \rho) \boldsymbol{\theta}$  // Update target policy network with EMA  
29 Return  $\boldsymbol{\theta}, \bar{\boldsymbol{\theta}}$ 
```

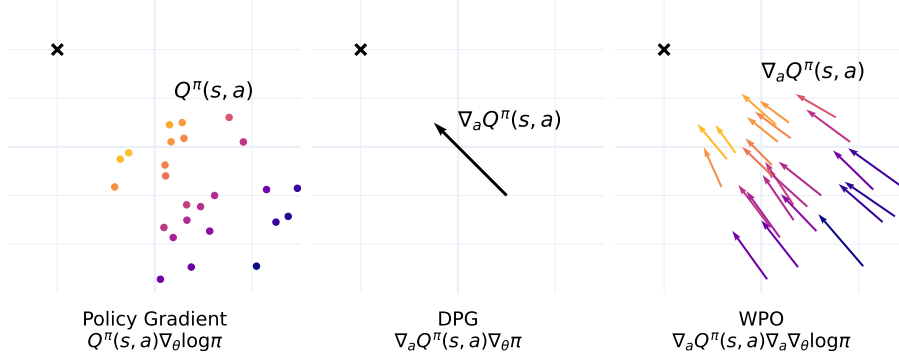


Figure 3.2: Conceptual illustration of how Wasserstein policy optimization (WPO) combines elements of stochastic and deterministic policy gradient methods for a 2d action space. Left: “classic” policy gradient. Samples are taken from a stochastic policy. Each sample contributes a scalar $Q^\pi(s, a)$ factor to the gradient. Middle: deterministic policy gradient (DPG). A deterministic action is chosen and the policy gradient depends on the gradient of $Q^\pi(s, a)$. Right: WPO. Samples are taken from a stochastic policy, as in classic policy gradient, but depend on the gradient of Q^π with respect to the action, as in DPG. *foo* From Figure 1 of [Pfa+25]. Used with kind permission of David Pfau.

(Note that the states are sampled from a replay buffer, so may be off-policy, but the actions are sampled from the current policy, so are on-policy.)

If we ignore the FIM preconditioner \mathcal{F}^{-1} , we see that the update is similar to the one used in the DPG theorem, except we replace the Jacobian $\nabla_{\theta} \mu_{\theta}(s)$ with $\nabla_{\theta} (\nabla_a \log \pi_{\theta}(a|s))^T \in \mathbb{R}^{N_{\theta} \times N_A}$. Intuitively this captures the change in **probability flow** over the action space due to a change in the parameters. See Figure 3.2 for an illustration.

However, the use of the FIM preconditioner keeps the update closer to the true gradient flow. (Indeed, in the case of a Gaussian policy and quadratic value function, WPO is exactly the Wasserstein gradient flow if you use the FIM, but is very different if you don’t.) Furthermore, this preconditioner can avoid numerical issues which can arise as the policy converges to a deterministic policy, leading to a blowing up of the gradient term $\nabla_a \log \pi_{\theta}(a|s)$.

In general, computing the FIM can be intractable. However, the authors assume the policy is a diagonal Gaussian, for which the FIM is diagonal:

$$\mathcal{F}(\mu, \sigma) = \begin{bmatrix} \text{diag} \left(\frac{1}{\sigma_1^2}, \frac{1}{\sigma_2^2}, \dots, \frac{1}{\sigma_d^2} \right) & 0 \\ 0 & \text{diag} \left(\frac{2}{\sigma_1^2}, \frac{2}{\sigma_2^2}, \dots, \frac{2}{\sigma_d^2} \right) \end{bmatrix} \quad (3.83)$$

This is fast to compute and invert.

After updating the policy with the above approach at each step, they updated the critic using a conventional DQN-like update, but one could use more sophisticated critic updates, such as TD3.

3.3 Policy improvement methods

In this section, we discuss methods that try to monotonically improve performance of the policy at each step, rather than just following the gradient, which can result in a high variance estimate where performance can increase or decrease at each step. These are called **policy improvement** methods. Our presentation is based on [QPC24].

3.3.1 Policy improvement lower bound

We start by stating a useful result from [Ach+17]. Let π_k be the current policy at step k , and let π be any other policy (e.g., a candidate new one). Let $p_{\pi_k}^{\gamma}$ be the normalized discounted state visitation distribution

for π_k , defined in Equation (3.22). Let $A^{\pi_k}(s, a) = Q^{\pi_k}(s, a) - V^{\pi_k}(s)$ be the advantage function. Finally, let the total variation distance between two distributions be given by

$$\text{TV}(p, q) \triangleq \frac{1}{2} \|\mathbf{p} - \mathbf{q}\|_1 = \frac{1}{2} \sum_s |p(s) - q(s)| \quad (3.84)$$

Then one can show [Ach+17] that

$$J(\pi) - J(\pi_k) \geq \frac{1}{1-\gamma} \underbrace{\mathbb{E}_{p_{\pi_k}^\gamma(s) \pi_k(a|s)} \left[\frac{\pi(a|s)}{\pi_k(a|s)} A^{\pi_k}(s, a) \right]}_{L(\pi, \pi_k)} - \frac{2\gamma C^{\pi, \pi_k}}{(1-\gamma)^2} \mathbb{E}_{p_{\pi_k}^\gamma(s)} [\text{TV}(\pi(\cdot|s), \pi_k(\cdot|s))] \quad (3.85)$$

where $C^{\pi, \pi_k} = \max_s |\mathbb{E}_{\pi(a|s)} [A^{\pi_k}(s, a)]|$. In the above, $L(\pi, \pi_k)$ is a surrogate objective, and the second term is a penalty term.

If we can optimize this lower bound (or a stochastic approximation, based on samples from the current policy π_k), we can guarantee monotonic policy improvement (in expectation) at each step. We will replace this objective with a trust-region update that is easier to optimize:

$$\pi_{k+1} = \underset{\pi}{\text{argmax}} L(\pi, \pi_k) \quad \text{s.t.} \quad \mathbb{E}_{p_{\pi_k}^\gamma(s)} [\text{TV}(\pi, \pi_k)(s)] \leq \epsilon \quad (3.86)$$

The constraint bounds the worst-case performance decline at each update. The overall procedure becomes an approximate policy improvement method. There are various ways of implementing the above method in practice, some of which we discuss below. (See also [GDWF22], who propose a framework called **mirror learning**, that justifies these “approximations” as in fact being the optimal thing to do for a different kind of objective; see also [Vas+21].)

3.3.2 Trust region policy optimization (TRPO)

In this section, we describe the **trust region policy optimization (TRPO)** method of [Sch+15b]. This implements an approximation to Equation (3.86). First, it leverages the fact that if

$$\mathbb{E}_{p_{\pi_k}^\gamma(s)} [D_{\text{KL}}(\pi_k \parallel \pi)(s)] \leq \delta \quad (3.87)$$

then π also satisfies the TV constraint with $\delta = \frac{\epsilon^2}{2}$. Next it considers a first-order expansion of the surrogate objective to get

$$L(\pi, \pi_k) = \mathbb{E}_{p_{\pi_k}^\gamma(s) \pi_k(a|s)} \left[\frac{\pi(a|s)}{\pi_k(a|s)} A^{\pi_k}(s, a) \right] \approx \mathbf{g}_k^\top (\boldsymbol{\theta} - \boldsymbol{\theta}_k) \quad (3.88)$$

where $\mathbf{g}_k = \nabla_{\boldsymbol{\theta}} L(\pi_{\boldsymbol{\theta}}, \pi_k)|_{\boldsymbol{\theta}_k}$. Finally it considers a second-order expansion of the KL term to get the approximate constraint

$$\mathbb{E}_{p_{\pi_k}^\gamma(s)} [D_{\text{KL}}(\pi_k \parallel \pi)(s)] \approx \frac{1}{2} (\boldsymbol{\theta} - \boldsymbol{\theta}_k)^\top \mathbf{F}_k (\boldsymbol{\theta} - \boldsymbol{\theta}_k) \quad (3.89)$$

where $\mathbf{F}_k = \mathbf{g}_k \mathbf{g}_k^\top$ is an approximation to the Fisher information matrix (see Equation (3.68)). We then use the update

$$\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \eta_k \mathbf{v}_k \quad (3.90)$$

where $\mathbf{v}_k = \mathbf{F}_k^{-1} \mathbf{g}_k$ is the natural gradient, and the step size is initialized to $\eta_k = \sqrt{\frac{2\delta}{\mathbf{v}_k^\top \mathbf{F}_k \mathbf{v}_k}}$. (In practice we compute \mathbf{v}_k by approximately solving the linear system $\mathbf{F}_k \mathbf{v} = \mathbf{g}_k$ using conjugate gradient methods, which just require matrix vector multiplies.) We then use a backtracking line search procedure to ensure the trust region is satisfied.

3.3.3 Proximal Policy Optimization (PPO)

In this section, we describe the **proximal policy optimization** or **PPO** method of [Sch+17], which is a simplification of TRPO.

We start by noting the following result:

$$\mathbb{E}_{p_{\pi_k}^\gamma(s)} [\text{TV}(\pi, \pi_k)(s)] = \frac{1}{2} \mathbb{E}_{(s,a) \sim p_{\pi_k}^\gamma} \left[\left| \frac{\pi(a|s)}{\pi_k(a|s)} - 1 \right| \right] \quad (3.91)$$

This holds provided the support of π is contained in the support of π_k at every state. We then use the following update:

$$\pi_{k+1} = \underset{\pi}{\operatorname{argmax}} \mathbb{E}_{(s,a) \sim p_{\pi_k}^\gamma} [\min(\rho_k(s,a)A^{\pi_k}(s,a), \tilde{\rho}_k(s,a)A^{\pi_k}(s,a))] \quad (3.92)$$

where $\rho_k(s,a) = \frac{\pi(a|s)}{\pi_k(a|s)}$ is the likelihood ratio, and $\tilde{\rho}_k(s,a) = \text{clip}(\rho_k(s,a), 1 - \epsilon, 1 + \epsilon)$, where $\text{clip}(x, l, u) = \min(\max(x, l), u)$. See [GDFW22] for a theoretical justification for these simplifications. Furthermore, this can be modified to ensure monotonic improvement as discussed in [WHT19], making it a true bound optimization method.

Some pseudocode for PPO (with GAE) is given in Algorithm 9. It is basically identical to the AC code in Algorithm 7, except the policy loss has the form $\min(\rho_t A_t, \tilde{\rho}_t A_t)$ instead of $A_t \log \pi_\phi(a_t|s_t)$, and we perform multiple policy updates per rollout, for increased sample efficiency. For all the implementation details, see <https://iclr-blog-track.github.io/2022/03/25/ppo-implementation-details/>.

Algorithm 9: PPO with GAE

```

1 Initialize parameters  $\phi$ , environment state  $s$ 
2 for iteration  $k = 1, 2, \dots$  do
3    $\phi_{\text{old}} \leftarrow \phi$ 
4    $(s_1, a_1, r_1, \dots, s_T) = \text{rollout}(s, \pi_{\phi_{\text{old}}})$ 
5    $v_t = V_\phi(s_t)$  for  $t = 1 : T$ 
6    $(A_{1:T}, y_{1:T}) = \text{GAE}(r_{1:T}, v_{1:T}, \gamma, \lambda)$ 
7   for  $m = 1 : M$  do
8      $\rho_t = \frac{\pi_\phi(a_t|s_t)}{\pi_{\phi_{\text{old}}}(a_t|s_t)}$  for  $t = 1 : T$ 
9      $\tilde{\rho}_t = \text{clip}(\rho_t)$  for  $t = 1 : T$ 
10     $\mathcal{L}(\phi) = \frac{1}{T} \sum_{t=1}^T [\lambda_{TD} \text{clip}((V_\phi(s_t) - y_t)^2) - \lambda_{PG} \min(\rho_t A_t, \tilde{\rho}_t A_t) - \lambda_{ent} \mathbb{H}(\pi_\phi(\cdot|s_t))]$ 
11     $\phi := \phi - \eta \nabla_\phi \mathcal{L}(\phi)$ 
```

3.3.3.1 Simplified form of the clipping term

In <https://spinningup.openai.com/en/latest/algorithms/ppo.html>, they propose a simplified form of the clipping term, which is easier to understand. First we rewrite the loss (inside the expectation of Equation (3.92)) in shorthand form as follows:

$$L(s, a, \theta_k, \theta) = \min(\rho A, \text{clip}(\rho)A) \quad (3.93)$$

where $\rho = \frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)}$ and $A = A^{\pi_{\theta_k}}(s, a)$. If the advantage is positive, $A > 0$, the action was better than expected, so we want to increase the probability ρ , but not too much (so the new policy stays close the previous one). The term $\text{clip}(\rho)A$ restricts ρA to be in $[(1 - \epsilon)A, (1 + \epsilon)A]$. The outer min selects the smaller of ρA and this clipped value. If $\rho < 1 - \epsilon$, then $\rho A < (1 - \epsilon)A$, so the min operator will naturally select ρA without needing the lower clip. So in this case we can use a loss of the form $\min(\rho A, (1 + \epsilon)A)$. If the

advantage is negative, $A < 0$, a similar argument shows we can use a loss of the form $\min(\rho A, (1 - \epsilon)A)$. We can combine these two cases by defining

$$L(s, a, \theta_k, \theta) = \min(\rho A, g(\epsilon, A)) \quad (3.94)$$

where

$$g(\epsilon, A) = \begin{cases} (1 + \epsilon)A & \text{if } A \geq 0 \\ (1 - \epsilon)A & \text{if } A < 0 \end{cases} \quad (3.95)$$

3.3.3.2 PPO for diffusion policies

The PPO algorithm relies on being able to compute the likelihood of a candidate action. This is difficult to do for diffusion policies. In [Ren+24], they propose **DPPO**, which is PPO for diffusion models. The key idea is to treat each step of the diffusion process as a step of an “inner” MDP, nested inside the outer (main) MDP, and then to apply PPO to this combined system. Unfortunately this results in very long horizon tasks. In [Yan+25], they propose noise-conditioned diffusion policy optimization, or **NCDPO**, which fixes the noise sequence for all the diffusion steps, and then deterministically backpropagates gradients through the entire denoising chain. The result is a single MDP, in which the policy happens to be implemented by a deep network, rather than a set of two nested MDPs. (See [Ueh+24] for a more general review of RL for diffusion models.)

3.3.3.3 Simple policy optimization

More recently, [Xie+25] propose **Simple Policy Optimization** or **SPO**, which improves upon ratio clipping, offering stronger theoretical properties and better constraining the probability ratio within the trust region.

3.3.4 Variational Maximum a Posteriori Policy Optimization (VMPO)

In this section, we discuss the **VMPO** algorithm of [FS+19], which is an on-policy extension of the earlier off-policy MPO (MAP policy optimization) algorithm that we discuss in Section 3.6.5. VMPO was originally explained in terms of the “control as inference” framework (see Section 3.6), but we can also view it as a constrained policy improvement method, based on Equation (3.86). In particular, VMPO leverages the fact that if

$$\mathbb{E}_{p_{\pi_k}^\gamma(s)} [D_{\text{KL}}(\pi \parallel \pi_k)(s)] \leq \delta \quad (3.96)$$

then π also satisfies the TV constraint with $\delta = \frac{\epsilon^2}{2}$.

Note that here the KL is reversed compared to TRPO in Section 3.3.2. This new version will encourage π to be mode-covering, so it will naturally have high entropy, which can result in improved robustness. Unfortunately, this kind of KL is harder to compute, since we are taking expectations wrt the unknown distribution π .

To solve this problem, VMPO adopts an EM-type approach. In the E step, we compute a non-parametric version of the state-action distribution given by the unknown new policy:

$$\psi(s, a) = \pi(a|s)p_{\pi_k}^\gamma(s) \quad (3.97)$$

The optimal new distribution is given by

$$\psi_{k+1} = \operatorname{argmax}_{\psi} \mathbb{E}_{\psi(s,a)} [A^{\pi_k}(s, a)] \quad \text{s.t.} \quad D_{\text{KL}}(\psi \parallel \psi_k) \leq \delta \quad (3.98)$$

where $\psi_k(s, a) = \pi_k(a|s)p_{\pi_k}^\gamma(s)$. The solution to this is

$$\psi_{k+1}(s, a) = p_{\pi_k}^\gamma(s)\pi_k(a|s)w(s, a) \quad (3.99)$$

$$w(s, a) = \frac{\exp(A^{\pi_k}(s, a)/\lambda^*)}{Z(\lambda^*)} \quad (3.100)$$

$$Z(\lambda) = \mathbb{E}_{(s,a) \sim p_{\pi_k}^\gamma} [\exp(A^{\pi_k}(s, a)/\lambda)] \quad (3.101)$$

$$\lambda^* = \operatorname{argmin}_{\lambda \geq 0} \lambda \delta + \lambda \log Z(\lambda) \quad (3.102)$$

In the M step, we project this target distribution back onto the space of parametric policies, while satisfying the KL trust region constraint:

$$\pi_{k+1} = \operatorname{argmax}_{\pi} \mathbb{E}_{(s,a) \sim p_{\pi_k}^\gamma} [w(s, a) \log \pi(a|s)] \quad \text{s.t.} \quad \mathbb{E}_{p_{\pi_k}^\gamma} [D_{\text{KL}}(\psi_k \parallel \psi)(s)] \leq \delta \quad (3.103)$$

3.4 Off-policy methods

In many cases, it is useful to train a policy using data collected from a distinct **behavior policy** $\pi_b(a|s)$ that is not the same as the **target policy** $\pi(a|s)$ that is being learned. For example, this could be data collected from earlier trials or parallel workers (with different parameters θ') and stored in a **replay buffer**, or it could be **demonstration data** from human experts. This is known as **off-policy RL**, and can be much more sample efficient than the on-policy methods we have discussed so far, since these methods can use data from multiple sources. However, off-policy methods are more complicated, as we will explain below.

The basic difficulty is that the target policy that we want to learn may want to try an action in a state that has not been experienced before in the existing data, so there is no way to predict the outcome of this new (s, a) pair. In this section, we tackle this problem by assuming that the target policy is not too different from the behavior policy, so that the ratio $\pi(a|s)/\pi_b(a|s)$ is bounded, which allows us to use methods based on importance sampling. In the online learning setting, we can ensure this property by using conservative incremental updates to the policy. Alternatively we can use policy gradient methods with various regularization methods, as we discuss below.

In Section 7.7, we discuss offline RL, which is an extreme instance of off-policy RL where we have a fixed behavioral dataset, possibly generated from an unknown behavior policy, and we can never collect any new data. Some of the solution methods for offline RL build on concepts discussed here.

3.4.1 Policy evaluation using importance sampling

Assume we have a dataset of the form $\mathcal{D} = \{\tau^{(i)}\}_{1 \leq i \leq n}$, where each trajectory is a sequence $\tau^{(i)} = (s_0^{(i)}, a_0^{(i)}, r_0^{(i)}, s_1^{(i)} \dots)$, where the actions are sampled according to a behavior policy π_b , and the reward and next states are sampled according to the reward and transition models. We want to use this offline dataset to evaluate the performance of some target policy π ; this is called **off-policy policy evaluation** or **OPE**. If the trajectories $\tau^{(i)}$ were sampled from π , we could use the standard Monte Carlo estimate:

$$\hat{J}(\pi) \triangleq \frac{1}{n} \sum_{i=1}^n \sum_{t=0}^{T-1} \gamma^t r_t^{(i)} \quad (3.104)$$

However, since the trajectories are sampled from π_b , we use **importance sampling** (IS) to correct for the distributional mismatch, as first proposed in [PSS00]. This gives

$$\hat{J}_{\text{IS}}(\pi) \triangleq \frac{1}{n} \sum_{i=1}^n \frac{p(\tau^{(i)}|\pi)}{p(\tau^{(i)}|\pi_b)} \sum_{t=0}^{T-1} \gamma^t r_t^{(i)} \quad (3.105)$$

It can be verified that $\mathbb{E}_{\pi_b} [\hat{J}_{\text{IS}}(\pi)] = J(\pi)$, that is, $\hat{J}_{\text{IS}}(\pi)$ is **unbiased**, provided that $p(\tau|\pi_b) > 0$ whenever $p(\tau|\pi) > 0$. The **importance ratio**, $\frac{p(\tau^{(i)}|\pi)}{p(\tau^{(i)}|\pi_b)}$, is used to compensate for the fact that the data is sampled from π_b and not π . It can be simplified as follows:

$$\frac{p(\tau|\pi)}{p(\tau|\pi_b)} = \frac{p(s_0) \prod_{t=0}^{T-1} \pi(a_t|s_t) p_S(s_{t+1}|s_t, a_t) p_R(r_t|s_t, a_t, s_{t+1})}{p(s_0) \prod_{t=0}^{T-1} \pi_b(a_t|s_t) p_S(s_{t+1}|s_t, a_t) p_R(r_t|s_t, a_t, s_{t+1})} = \prod_{t=0}^{T-1} \frac{\pi(a_t|s_t)}{\pi_b(a_t|s_t)} \quad (3.106)$$

This simplification makes it easy to apply IS, as long as the target and behavior policies are known. (If the behavior policy is unknown, we can estimate it from \mathcal{D} , and replace π_b by its estimate $\hat{\pi}_b$. For convenience, define the **per-step importance ratio** at time t by

$$\rho_t(\tau) \triangleq \pi(a_t|s_t) / \pi_b(a_t|s_t) \quad (3.107)$$

We can reduce the variance of the estimator by noting that the reward r_t is independent of the trajectory beyond time t . This leads to a **per-decision importance sampling** variant:

$$\hat{J}_{\text{PDIS}}(\pi) \triangleq \frac{1}{n} \sum_{i=1}^n \sum_{t=0}^{T-1} \prod_{t' \leq t} \rho_{t'}(\tau^{(i)}) \gamma^t r_t^{(i)} \quad (3.108)$$

3.4.2 Off-policy actor critic methods

In this section, we discuss how to extend actor-critic methods to work with off-policy data.

3.4.2.1 Learning the critic using V-trace

In this section we build on Section 3.4.1 to develop a practical method, known as **V-trace** [Esp+18], to estimate the value function for a target policy using off-policy data. (This is an extension of the earlier **Retrace** algorithm [Mun+16], which estimates the Q function using off-policy data.)

First consider the n -step target value for $V(s_i)$ in the on-policy case:

$$V_i = V(s_i) + \sum_{t=i}^{i+n-1} \gamma^{t-i} r_t + \gamma^n V(s_{i+n}) \quad (3.109)$$

$$= V(s_i) + \sum_{t=i}^{i+n-1} \gamma^{t-i} \underbrace{(r_t + \gamma V(s_{t+1}) - V(s_t))}_{\delta_t} \quad (3.110)$$

where we define $\delta_t = (r_t + \gamma V(s_{t+1}) - V(s_t))$ as the TD error at time t . To extend this to the off-policy case, we use the per-step importance ratio trick. However, to bound the variance of the estimator, we truncate the IS weights. In particular, we define

$$c_t = \min \left(\bar{c}, \frac{\pi(a_t|s_t)}{\pi_b(a_t|s_t)} \right), \quad \rho_t = \min \left(\bar{\rho}, \frac{\pi(a_t|s_t)}{\pi_b(a_t|s_t)} \right) \quad (3.111)$$

where \bar{c} and $\bar{\rho}$ are hyperparameters. We then define the V-trace target value for $V(s_i)$ as

$$v_i = V(s_i) + \sum_{t=i}^{i+n-1} \gamma^{t-i} \left(\prod_{t'=i}^{t-1} c_{t'} \right) \rho_t \delta_t \quad (3.112)$$

Note that we can compute these targets recursively using

$$v_i = V(s_i) + \rho_i \delta_i + \gamma c_i (v_{i+1} - V(s_{i+1})) \quad (3.113)$$

The product of the weights $c_i \dots c_{t-1}$ (known as the “trace”) measures how much a temporal difference δ_t at time t impacts the update of the value function at earlier time i . If the policies are very different, the variance of this product will be large. So the truncation parameter \bar{c} is used to reduce the variance. In [Esp+18], they find $\bar{c} = 1$ works best.

The use of the target $\rho_t \delta_t$ rather than δ_t means we are evaluating the value function for a policy that is somewhere between π_b and π . For $\bar{\rho} = \infty$ (i.e., no truncation), we converge to the value function V^π , and for $\bar{\rho} \rightarrow 0$, we converge to the value function V^{π_b} . In [Esp+18], they find $\bar{\rho} = 1$ works best. (An alternative to clipping the importance weights is to use a resampling technique, and then use unweighted samples to estimate the value function [Sch+19].)

In [SHS20], they analyse the variance of the V-trace estimator, used to compute ρ_t in Equation (3.111). They show that to keep this bounded, it is necessary to mix some off-policy data (from the replay buffer) with some fresh online data from the current policy.

Note that if $\bar{c} = \bar{\rho}$, then $c_i = \rho_i$. This gives rise to the simplified form

$$v_t = V(s_t) + \sum_{j=0}^{n-1} \gamma^j \left(\prod_{m=0}^j c_{t+m} \right) \delta_{t+j} \quad (3.114)$$

We can use the above V-trace targets to learn an approximate value function by minimizing the usual ℓ_2 loss:

$$\mathcal{L}(\mathbf{w}) = \mathbb{E}_{t \sim \mathcal{D}} [(v_t - V_{\mathbf{w}}(s_t))^2] \quad (3.115)$$

3.4.2.2 Learning the actor

We now discuss how to update the actor using an off-policy estimate of the policy gradient. We start by defining the objective to be the expected value of the new policy, where the states are drawn from the behavior policy’s state distribution, but the actions are drawn from the target policy:

$$J_{\pi_b}(\pi_\theta) = \sum_s p_{\pi_b}^\gamma(s) V_\pi(s) = \sum_s p_{\pi_b}^\gamma(s) \sum_a \pi_\theta(a|s) Q_\pi(s, a) \quad (3.116)$$

Differentiating this and ignoring the term $\nabla_\theta Q_\pi(s, a)$, as suggested by [DWS12], gives a way to (approximately) estimate the **off-policy policy-gradient** using a one-step IS correction ratio:

$$\nabla_\theta J_{\pi_b}(\pi_\theta) \approx \sum_s \sum_a p_{\pi_b}^\gamma(s) \nabla_\theta \pi_\theta(a|s) Q_\pi(s, a) \quad (3.117)$$

$$= \mathbb{E}_{p_{\pi_b}^\gamma(s), \pi_b(a|s)} \left[\frac{\pi_\theta(a|s)}{\pi_b(a|s)} \nabla_\theta \log \pi_\theta(a|s) Q_\pi(s, a) \right] \quad (3.118)$$

In practice, we can approximate $Q_\pi(s_t, a_t)$ by $q_t = r_t + \gamma v_{t+1}$, where v_{t+1} is the V-trace estimate for state s_{t+1} . If we use $V(s_t)$ as a baseline, to reduce the variance, we get the following gradient estimate for the policy:

$$\nabla J(\theta) = \mathbb{E}_{t \sim \mathcal{D}} [\rho_t \nabla_\theta \log \pi_\theta(a_t|s_t) (r_t + \gamma v_t - V_{\mathbf{w}}(s_t))] \quad (3.119)$$

We can also replace the 1-step IS-weighted TD error $\rho_t(r_t + \gamma v_t - V_{\mathbf{w}}(s_t))$ with an IS-weighted GAE value by modifying the generalized advantage estimation method in Section 3.2.2 to replace A_t with $\rho_t A_t$. See Algorithm 10 for some pseudocode.

Algorithm 10: Actor-Critic (off-policy)

Data: Learning rates $\alpha_\theta, \alpha_\phi$, discount γ , GAE λ , tapering func. $\tau(\cdot)$

```
1 Initialize: Actor  $\pi_\theta$ , Critic  $V_\phi$ , Target  $V_{\phi'} \leftarrow V_\phi$ , Replay buffer  $\mathcal{D}$ 
2 for  $episode = 1$  to  $M$  do
3   Sample initial state  $s_0$ 
4   Initialize empty episode buffer  $\mathcal{E} \leftarrow []$ 
5   for  $t = 0$  to  $T - 1$  do
6     Get action  $a_t \sim \pi_b(\cdot|s_t)$ 
7     Execute  $a_t$ , observe  $s_{t+1}$  and  $r_t$ 
8     Store  $(s_t, a_t, r_t, s_{t+1}, \pi_b(a_t|s_t))$  in  $\mathcal{E}$ 
9      $s_t \leftarrow s_{t+1}$ 
10  end
11  Store trajectory  $\mathcal{E}$  in the replay buffer  $\mathcal{D}$ 
12  if enough data in  $\mathcal{D}$  then
13    Sample a batch of trajectories  $\{\mathcal{E}_i\}$  from  $\mathcal{D}$ 
14    for each trajectory  $\mathcal{E}_i$  in the batch do
15      // Calculate Advantage Estimates using GAE
16      for  $t = T - 1, \dots, 0$  do
17         $\delta_t \leftarrow r_t + \gamma V_{\phi'}(s_{t+1}) - V_\phi(s_t)$  // TD Error
18         $\hat{A}_t \leftarrow \delta_t + \gamma \lambda \hat{A}_{t+1}$  // Advantage
19      end
20      // Update Critic
21       $L(\phi) \leftarrow \frac{1}{T} \sum_{t=0}^{T-1} (r_t + \gamma V_{\phi'}(s_{t+1}) - V_\phi(s_t))^2$ 
22       $\phi \leftarrow \phi - \alpha_\phi \nabla_\phi L(\phi)$ 
23      // Update Actor
24       $J(\theta) \leftarrow \frac{1}{T} \sum_{t=0}^{T-1} \rho_t \log \pi_\theta(a_t|s_t) \hat{A}_t$ 
25      where  $\rho_t = \min\left(\bar{\rho}, \frac{\pi_\theta(a_t|s_t)}{\pi_b(a_t|s_t)}\right)$ 
26       $\theta \leftarrow \theta + \alpha_\theta \nabla_\theta J(\theta)$ 
27    end
28  Update target network:  $\phi' \leftarrow (1 - \alpha_{\text{polyak}})\phi' + \alpha_{\text{polyak}}\phi$ 
29 end
```

3.4.2.3 Example: IMPALA

As an example of an off-policy AC method, we consider **IMPALA**, which stands for “Importance Weighted Actor-Learning Architecture”. [Esp+18]. This uses shared parameters for the policy and value function (with different output heads), and adds an entropy bonus to ensure the policy remains stochastic. Thus we end up with the following objective, which is very similar to on-policy actor-critic shown in Algorithm 7:

$$\mathcal{L}(\phi) = \mathbb{E}_{t \sim \mathcal{D}} [\lambda_{TD} (V_\phi(s_t) - v_t)^2 - \lambda_{PG} \rho_t A_t \log \pi_\phi(a_t|s_t) - \lambda_{ent} \mathbb{H}(\pi_\phi(\cdot|s_t))] \quad (3.120)$$

The only difference from standard A2C is that we need to store the probabilities of each action, $\pi_b(a_t|s_t)$, in addition to (s_t, a_t, r_t, s_{t+1}) in the dataset \mathcal{D} , which can be used to compute the importance ratio ρ_t in Equation (3.111). [Esp+18] was able to use this method to train a single agent (using a shared CNN and LSTM for both value and policy) to play all 57 games at a high level. Furthermore, they showed that their method — thanks to its off-policy corrections — outperformed the A3C method (a parallel version of A2C) in Section 3.2.1.

3.4.2.4 Off-policy learning with deterministic policies

In Section 3.2.6, we discuss actor-critic methods that learn deterministic policies for continuous actions, based on the deterministic policy gradient (DPG) method. These methods can work with off-policy without the need to do importance sampling correction.

3.4.2.5 PGQL: Combining off-policy Q-learning with policy gradient

In [O'D+16], they present a method known as **PGQL**, which uses Q learning to learn from off-policy data in a replay buffer, and uses this to regularize the on-policy actor-critic learner.

3.4.3 Off-policy policy improvement methods

So far we have focused on actor-critic methods. However, policy improvement methods, such as PPO, are often preferred to AC methods, since they monotonically improve the objective. In [QPC21] they propose one way to extend PPO to the off-policy case. This method was generalized in [QPC24] to cover a variety of policy improvement algorithms, including TRPO and VMPO. We give a brief summary below.

The key insight is to realize that we can generalize the lower bound in Equation (3.85) to any reference policy

$$J(\pi) - J(\pi_k) \geq \frac{1}{1-\gamma} \mathbb{E}_{p_{\pi_{\text{ref}}}(s)\pi_k(a|s)} \left[\frac{\pi(a|s)}{\pi_{\text{ref}}(a|s)} A^{\pi_k}(s, a) \right] - \frac{2\gamma C^{\pi, \pi_k}}{(1-\gamma)^2} \mathbb{E}_{p_{\pi_{\text{ref}}}(s)} [\text{TV}(\pi(\cdot|s), \pi_{\text{ref}}(\cdot|s))] \quad (3.121)$$

The reference policy can be any previous policy, or a convex combination of them. In particular, if π_k is the current policy, we can consider the reference policy to be $\pi_{\text{ref}} = \sum_{i=1}^M \nu_i \pi_{k-i}$, where $0 \leq \nu_i \leq 1$ and $\sum_i \nu_i = 1$ are mixture weights. We can approximate the expectation by sampling from the replay buffer, which contains samples from older policies. That is, $(s, a) \sim p_{\pi_{\text{ref}}}^\gamma$ can be implemented by $i \sim \nu$ and $(s, a) \sim p_{\pi_{k-i}}^\gamma$.

To compute the advantage function A^{π_k} from off policy data, we can adapt the V-trace method of Equation (3.114) to get

$$A_{\text{trace}}^{\pi_k}(s_t, a_t) = \delta_t + \sum_{j=0}^{n-1} \gamma^j \left(\prod_{m=1}^j c_{t+m} \right) \delta_{t+j} \quad (3.122)$$

where $\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t)$, and $c_t = \min\left(\bar{c}, \frac{\pi_k(a_t|s_t)}{\pi_{k-i}(a_t|s_t)}\right)$ is the truncated importance sampling ratio.

To compute the TV penalty term from off policy data, we need to choose between the PPO (Section 3.3.3), VMPO (Section 3.3.4) and TRPO (Section 3.3.2) approach.

We can derive an off-policy version of PPO using an update of the following form (known as **Generalized PPO**):

$$\pi_{k+1} = \underset{\pi}{\operatorname{argmax}} \mathbb{E}_{i \sim \nu} \left[\mathbb{E}_{(s,a) \sim p_{\pi_{k-i}}^\gamma} [\min(\rho_{k-i}(s, a) A^{\pi_k}(s, a), \tilde{\rho}_{k-i}(s, a) A^{\pi_k}(s, a))] \right] \quad (3.123)$$

where $\rho_{k-i}(s, a) = \frac{\pi(a|s)}{\pi_{k-i}(a|s)}$ and $\tilde{\rho}_{k-i}(s, a) = \text{clip}(\frac{\pi(a|s)}{\pi_{k-i}(a|s)}, l, u)$, where $l = \frac{\pi_k(a|s)}{\pi_{k-i}(a|s)} - \epsilon$ and $u = \frac{\pi_k(a|s)}{\pi_{k-i}(a|s)} + \epsilon$. (For other off-policy variants of PPO, see e.g., [Men+23; LMW24].)

For details on the off-policy version of TRPO, see [QPC24].

For an off-policy version of VMPO, see the discussion of MPO in Section 3.6.5.

3.5 Gradient-free policy optimization

So far, we have focused on fitting parametric policies, represented by differentiable functions $\pi_\theta(a|s)$, using methods based on the policy gradient theorem. Unfortunately, such gradient-based methods can get stuck in poor local optima. In addition, gradient descent cannot be applied to non-differentiable policies, such as programs, or functions with discrete latent variables (e.g., if-then branches). We can therefore consider other

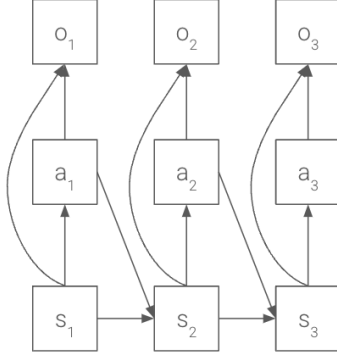


Figure 3.3: A graphical model for optimal control.

kinds of methods for policy learning, based on **blackbox optimization**, aka **derivative-free optimization**. This includes techniques such as **cross-entropy method** and **evolutionary strategies**. For details on such algorithms, see e.g. [Mur23, Sec 7.7]. For some applications of ES to RL, see e.g. [MGR18] (who obtain good results by training linear policies with random search), [Sal+17] (who use evolutionary strategies to optimize the policy of a robotic controller), and [Sar+25] (who discuss a method called “Evolution Guided General Optimization via Low-rank Learning” (EGGROLL), which provides a way to scale to very large models).

3.6 RL as inference

In this section, we discuss an approach to policy optimization that reduces it to probabilistic inference. This is called **control as inference**, or **RL as inference**, and has been discussed in numerous works (see e.g., [Att03; TS06; Tou09; ZABD10; RTV12; BT12; KGO12; HR17; Lev18; Fur+21; Zha+24d]). The primary advantage of this approach is that it enables policy learning using off-policy data, while avoiding the need to use (potentially high variance) importance sampling corrections. (This is because the inference approach takes expectations wrt $d_q(s)$ instead of $d_\pi(s)$, where q is an auxiliary distribution, π is the policy which is being optimized, and d is the state visitation measure.) A secondary advantage is that it enables us to use the large toolkit of methods for probabilistic modeling and inference to solve RL problems.⁶ The resulting framework forms the foundation of the the MPO discussed in Section 3.6.5, the SAC method discussed in Section 3.6.8, as well as the SMC planning method discussed in Section 4.2.3, and some kinds of LLM test-time inference, as discussed in Section 6.2.4.11.

The core of these methods is based on the probabilistic model shown in Figure 3.3. This shows an MDP augmented with new variables, \mathcal{O}_t . These are called **optimality variables**, and indicating whether the action at time t is optimal or not. We assume these have the following probability distribution:

$$p(\mathcal{O}_t = 1 | s_t, a_t) \propto \exp(\eta^{-1} G(s_t, a_t)) \quad (3.124)$$

where $\eta > 0$ is a temperature parameter, and $G(s, a)$ is some quality function, such as $G(s, a) = R(s, a)$, or $G(s, a) = Q(s, a)$ or $G(s, a) = A(s, a)$. For brevity, we will just write $p(\mathcal{O} = 1 | \cdot)$ to denote the probability of the event that $\mathcal{O}_t = 1$ for all time steps. (Note that the specific value of 1 is arbitrary; this likelihood function is really just a non-negative weighting term that biases the action trajectory, as we show below.)

⁶Note, however, that we do not tackle the problem of epistemic uncertainty (exploration). Solving this in the context of RL-as-inference requires additional machinery, as discussed in [TLO23].

3.6.1 Deterministic case (planning/control as inference)

Our goal is to find trajectories that are optimal. That is, we would like to find the mode (or posterior samples) from the following distribution:

$$p(\boldsymbol{\tau}|\mathcal{O} = 1, \pi) \propto p(\boldsymbol{\tau}, \mathcal{O} = 1|\pi) \propto \left[p(s_1) \prod_{t=1}^{T-1} \pi(a_t|s_t) p(s_{t+1}|s_t, a_t) \right] \left[\prod_{t=1}^T p(\mathcal{O}_t = 1|s_t, a_t) \right] \quad (3.125)$$

where π is the policy.

Let us start by considering the deterministic case, where $p(s_{t+1}|s_t, a_t)$ is either 1 or 0, depending on whether the transition is feasible or not. In this case, rather than learning a policy π that maps states to actions we just need to learn a plan (a specific sequence of action $\mathbf{a}_{1:T}$) for each starting state s_1 . This is equivalent to a **shortest path** problem, i.e., we want to maximize

$$p(\boldsymbol{\tau}|\mathcal{O} = 1, \mathbf{a}_{1:T}) \propto p(s_1) \left[\prod_{t=1}^{T-1} p(s_{t+1}|s_t, a_t) \right] \left[\exp\left(\sum_{t=1}^T R(s_t, a_t)\right) \right] \quad (3.126)$$

(Typically the initial state s_1 is known, in which case $p(s_1)$ is a delta function.)

The MAP sequence of actions, which we denote by $\hat{\mathbf{a}}_{1:T}(s_1)$, is the optimal **open loop plan**. (It is called “open loop” since the agent does not need to observe the state, since s_t is uniquely determined by s_1 and $\mathbf{a}_{1:t}$, both of which are known.) Computing this trajectory is known as the **control as inference** problem [Wat+21]. Such open loop planning problems can be solved using model predictive control methods, discussed in Section 4.2.4.

3.6.2 Stochastic case (policy learning as variational inference)

In the stochastic case, we want to learn a policy π which maps states to actions, and which generates a distribution over trajectories which are optimal. Thus we define the objective as

$$\log p(\mathcal{O} = 1|\pi) = \log \int p_\pi(\boldsymbol{\tau}) p(\mathcal{O} = 1|\boldsymbol{\tau}) d\boldsymbol{\tau} \quad (3.127)$$

where we define

$$p_\pi(\boldsymbol{\tau}) = p(s_1) \prod_t p(s_{t+1}|s_t, a_t) \pi(a_t|s_t) \quad (3.128)$$

Since marginalizing over trajectories is difficult, we introduce a variational distribution $q(\boldsymbol{\tau})$ to simplify the computations. We assume q factors in the same way:

$$q(\boldsymbol{\tau}) = p(s_1) \prod_t p(s_{t+1}|s_t, a_t) \pi_q(a_t|s_t) \quad (3.129)$$

Note that we use the true dynamics model $p(s_{t+1}|s_t, a_t)$ when defining q , and only introduce the variational distribution for the actions, $\pi_q(a_t|s_t)$. This is one way to avoid the **optimism bias** that can arise if we sample from an unconstrained $q(\boldsymbol{\tau}|\mathcal{O} = 1)$. To see this, suppose $\mathcal{O} = 1$ is the event that we win the lottery. We do not want conditioning on this outcome to influence our belief in the the probability of chance events, which is governed by $p(s_{t+1}|s_t, a_t)$ and not $p(s_{t+1}|s_t, a_t, \mathcal{O} = 1)$. See [Lev18] for further discuss of this point.

Now note the following identity

$$D_{\text{KL}}(q(\boldsymbol{\tau}) \parallel p_\pi(\boldsymbol{\tau}|\mathcal{O} = 1)) = \mathbb{E}_q \left[\log q(\boldsymbol{\tau}) - \log \frac{p_\pi(\mathcal{O} = 1|\boldsymbol{\tau}) p_\pi(\boldsymbol{\tau})}{p_\pi(\mathcal{O} = 1)} \right] \quad (3.130)$$

$$= \mathbb{E}_q [\log q(\boldsymbol{\tau}) - \log p_\pi(\mathcal{O} = 1|\boldsymbol{\tau}) - \log p_\pi(\boldsymbol{\tau})] + \log p_\pi(\mathcal{O} = 1) \quad (3.131)$$

Hence

$$\log p_\pi(\mathcal{O} = 1) = \mathbb{E}_q \left[\log p(\mathcal{O} = 1 | \tau) - \log \frac{q(\tau)}{p(\tau)} + \log \frac{q(\tau)}{p(\tau | \mathcal{O} = 1)} \right] \quad (3.132)$$

$$= J(p_\pi, q) + D_{\text{KL}}(q(\tau) \parallel p_\pi(\tau | \mathcal{O} = 1)) \quad (3.133)$$

where J is defined by

$$J(p_\pi, \pi_q) = \mathbb{E}_q [\log p_\pi(\mathcal{O} = 1 | \tau)] - D_{\text{KL}}(q(\tau) \parallel p_\pi(\tau)) \quad (3.134)$$

$$= \sum_{t=1}^T \mathbb{E}_q [\eta^{-1} G(s_t, a_t) - D_{\text{KL}}(\pi_q(\cdot | s_t) \parallel \pi_p(\cdot | s_t))] \quad (3.135)$$

Since $D_{\text{KL}}(q(\tau) \parallel p_\pi(\tau | \mathcal{O} = 1)) \geq 0$, we see that $\log p(\mathcal{O} = 1 | \pi) \geq J(p_\pi, q)$; hence J is called the **evidence lower bound** or **ELBO**. We can define the policy learning task as maximizing the ELBO, subject to the constraints that π_p and π_q are distributions that integrate to 1 across actions for all states.

To extend to the infinite time discounted case, we define $d_\pi(s)$ as the unnormalized discounted distribution over states

$$d_\pi(s) = \sum_{t=1}^{\infty} \gamma^t p(s_t = s | \pi) \quad (3.136)$$

We now replace the $\sum_t E_q(s)$ with $E_{d_q(s)}$ to get the constrained objective

$$\max_{\pi_p, \pi_q} J(\pi_p, \pi_q) \quad \text{s.t.} \quad \int d_q(s) \int \pi_p(a | s) da ds = 1, \quad \int d_q(s) \int \pi_q(a | s) da ds = 1 \quad (3.137)$$

There are two main ways to solve this optimization problem, which we call “EM control” and “KL control”, following [Fur+21]. We describe these below.

3.6.3 EM control

In this section, we discuss ways to optimize Equation (3.137) using the **Expectation Maximization** or **EM** algorithm, which is a widely used **bound optimization** method, also called a **MM** (majorize / maximize) method, that monotonically increases a lower bound on its objective (see [HL04] for a tutorial). In the E step, we maximize J wrt a non-parametric representation of the variational posterior π_q , while holding the parametric prior $\pi_p = \pi_{\theta_p}^{k-1}$ fixed at the value from the previous ($k-1$ ’th) iteration, to get π_q^k . In the M step, we then maximize J wrt π_p , holding the variational posterior fixed at π_q^k , to get the updated policy $\pi_{\theta_p}^k$.

In more detail, in the E step we maximize the following wrt π_q :

$$\begin{aligned} J(\pi_{\theta_p}^{k-1}, \pi_q) &= \int d_q(s) \int \pi_q(s | a) \eta^{-1} G(s, a) da ds \\ &\quad - \int d_q(s) \int \pi_q(a | s) \log \frac{\pi_q(a | s)}{\pi_{\theta_p}^{k-1}(a | s)} da ds + \lambda \left(1 - \int d_{\pi(s)} \int \pi_q(a | s) da ds \right) \end{aligned} \quad (3.138)$$

where λ is a Lagrange multiplier. The optimal (non-parametric) solution to this is

$$\pi_q^k(a | s) = Z(s)^{-1} \pi_{\theta_p}^{k-1}(a | s) \exp(\eta^{-1} G(s, a)) \quad (3.139)$$

where Z is the partition function

$$Z(s) = \int \pi_{\theta_p}^{k-1}(a | s) \exp(\eta^{-1} G(s, a)) da \quad (3.140)$$

In the M step, we maximize the following wrt π_{θ_p} :

$$J(\pi_q^k, \pi_p) = \mathbb{E}_{d_q(s) \pi_q^k(a | s)} [\log \pi_{\theta_p}(a | s)] \quad (3.141)$$

which we recognize as a weighted maximum likelihood problem.

3.6.4 KL control (maximum entropy RL)

In KL control, we only optimize the variational posterior π_q , holding the prior π_p fixed. Thus we only have an E step. In addition, we represent π_q parameterically, as π_{θ_q} , instead of the non-parametric approach used by EM. If the prior π_p is uniform, and we use $G(s_t, a_t) = R(s_t, a_t)$, then Equation (3.135) becomes

$$\eta J(\pi_p, \pi_q) = \sum_{t=1}^T \mathbb{E}_q [R(s_t, a_t) - \eta H(\pi_q(\cdot|s_t))] \quad (3.142)$$

where $-H(q) = D_{\text{KL}}(q \parallel \text{unif}) = \sum_a q(a) \log \frac{q(a)}{c}$ is the negative entropy function and c is a constant. This is called the **maximum entropy RL** objective [ZABD10; Haa+18a; Haa+18b]. This differs from the standard objective used in RL training (namely a lower bound on sum of expected rewards) by virtue of the addition of the entropy regularizer on the policy. See Section 3.6.8 for further discussion.

3.6.5 Maximum a Posteriori Policy Optimization (MPO)

In this section, we discuss the **MPO** method of [Abd+18]. This is an instance of EM control, where $G(s, a) = Q(s, a)$, which is estimated using the retrace algorithm (see Section 3.4.2.1) or a single-step Bellman update.

It implements the E step using Equation (3.139), where we approximate $Z(s)$ with Monte Carlo:

$$q^k(a|s) = \frac{1}{\hat{Z}(s)} \pi_{\theta_p^{k-1}}(a|s) \exp(\eta^{-1} G(s, a)) \quad (3.143)$$

$$Z(s) \approx \frac{1}{M} \sum_{j=1}^M \exp(\eta^{-1} G(s, a_j)), \quad a_j \sim \pi_{\theta_p^{k-1}}(\cdot|s) \quad (3.144)$$

In addition, the (inverse) temperature parameter η is solved for by minimizing the dual of Equation (3.139), which is given by

$$g(\eta) = \eta \epsilon + \eta \log \mathbb{E}_{d_q(s) \pi_{\theta_p^{k-1}}(a|s)} \left[\exp \left(\eta^{-1} Q^{\pi_{\theta_p^{k-1}}}(s, a) \right) \right] \quad (3.145)$$

In the M step, MPO augments the objective in Equation (3.141) with a log prior at the k 'th step of the form $\log p_k(\theta_p)$ to create a MAP estimate. That is, it optimizes the following wrt θ_p :

$$J(q^k, \pi_{\theta_p}) = \mathbb{E}_{d_q(s) q(a|s)} [\log \pi_{\theta_p}(a|s)] + \log p_k(\theta_p) \quad (3.146)$$

We can think of this step as projecting the non-parametric policy q back to the space of parameterizable policies Π_{θ} .

We assume the prior is a Gaussian centered at the previous parameters,

$$p_k(\theta) = \mathcal{N}(\theta | \theta_k, \lambda \mathbf{F}_k) = c \exp(-\lambda(\theta - \theta_k)^\top \mathbf{F}_k^{-1}(\theta - \theta_k)) \quad (3.147)$$

where \mathbf{F}_k is the Fisher information matrix. If we view this as a second order approximation to the KL, we can rewrite the objective as

$$\max_{\theta_p} E_{d_q(s)} [E_{q(a|s)} \log \pi(a|s, \theta_p) - \lambda D_{\text{KL}}(\pi(a|s, \theta_k) \parallel \pi(a|s, \theta_p))] \quad (3.148)$$

We can approximate the expectation wrt $d_q(s)$ by sampling states from a replay buffer, and the expectation wrt $q(a|s)$ by sampling from the policy. The KL term can be computed analytically for Gaussian policies. We can then optimize this objective using SGD.

Note that we can also rewrite this as a constrained optimization problem

$$\max_{\theta_p} E_{d_q(s)} [E_{q(a|s)} \log \pi(a|s, \theta_p)] \quad \text{s.t.} \quad E_{d_q(s)} [D_{\text{KL}}(\pi(a|s, \theta_k) \parallel \pi(a|s, \theta_p))] \leq \epsilon_m \quad (3.149)$$

This can be optimized using a trust region method.

3.6.6 Sequential Monte Carlo Policy Optimisation (SMC-PO)

In this section, we discuss **SMC-PO** method of [Mac+24]. This is a model-based version of MPO, which uses Sequential Monte Carlo (SMC) to perform approximate inference in the E step. In particular, it samples from a distribution over optimal future trajectories starting from the current state, s_t , and using the current policy π_{θ_p} and dynamics model $\mathcal{T}(s'|s, a)$. From this it derives a non-parametric distribution over the optimal actions to take at the next step, $q(a_t|s_t)$. (see Section 4.2.3 for details). This becomes a target for the parametric policy update in the M step, which is the same weighted maximum likelihood method used by MPO.

3.6.7 AWR and AWAC

The **Advantage Weighted Regression** or **AWR** method of [Pen+19] and the **Advantage Weighted Actor Critic** or **AWAC** method of [Nai+20] are both EM control methods. AWR uses $G(s, a) = A(s, a)$, where the advantage function is estimated using GAE. The value function $V(s)$ is estimated using TD(λ), and is the value for the average of all previous policies, $\tilde{\pi}_{p^k} = \frac{1}{k} \sum_{j=0}^{k-1} \pi_{\theta_p^j}$. In contrast, AWAC uses $G(s, a) = Q(s, a)$, which is estimated by TD(0).

The (non-parametric) E step is closed form, as in other EM control methods, where the temperate η is treated as a hyper-parameter. The (parametric) M step is a weighted maximum likelihood step that is solved with SGD.

3.6.8 Soft Actor Critic (SAC)

The **soft actor-critic (SAC)** algorithm [Haa+18a; Haa+18b] is an off-policy actor-critic method based on the maximum entropy RL method we discussed in Section 3.6.4. This is an instance of the KL control scheme where the variational posterior policy $\pi_q = \pi_{\theta_q}$ is parameterized, but the prior policy π_p is fixed to the uniform distribution. (Thus SAC only has an E step (implemented with SGD), but no M step.) SAC uses $G(s, a) = Q^{\text{soft}}(s, a)$, where the soft-Q function is defined below.

Crucially, even though SAC is off-policy and utilizes a replay buffer to sample past experiences, the policy update is done using the actor's own probability distribution, eliminating the need to use importance sampling to correct for discrepancies between the behavior policy (used to collect data) and the target policy (used for updating), as we will see below.

3.6.8.1 SAC objective

We can write the maxent RL objective for the E step by using Equation (3.142) with slightly modified notation:

$$J^{\text{SAC}}(\theta) \triangleq \mathbb{E}_{p_{\pi_{\theta}}(s) \pi_{\theta}(a|s)} [R(s, a) + \alpha \mathbb{H}(\pi_{\theta}(\cdot|s))] \quad (3.150)$$

Note that the entropy term makes the objective easier to optimize, and encourages exploration. To optimize this, we can perform a policy evaluation step, and then a policy improvement step.

3.6.8.2 Policy evaluation: tabular case

We can perform policy evaluation by repeatedly applying a modified Bellman backup operator \mathcal{T}^{π} defined as

$$\mathcal{T}^{\pi}Q(s_t, \mathbf{a}_t) = r(s_t, \mathbf{a}_t) + \gamma \mathbb{E}_{\mathbf{s}_{t+1} \sim p} [V(\mathbf{s}_{t+1})] \quad (3.151)$$

where

$$V(s_t) = \mathbb{E}_{\mathbf{a}_t \sim \pi} [Q(s_t, \mathbf{a}_t) - \alpha \log \pi(\mathbf{a}_t|s_t)] \quad (3.152)$$

is the **soft value function**. If we iterate $Q^{k+1} = \mathcal{T}^{\pi}Q^k$, this will converge to the soft Q function for π .

In the tabular case, we can derive the optimal soft value function as follows. First, by definition, we have

$$V^*(s) := \max_{\pi} \sum_a \pi(a | s) [Q^*(s, a) - \alpha \log \pi(a | s)]. \quad (3.153)$$

This is a constrained optimization problem, where $\pi(\cdot | s)$ is a probability distribution. We introduce a Lagrange multiplier λ to enforce the normalization constraint:

$$\mathcal{L}(\pi, \lambda) = \sum_a \pi(a | s) [Q^*(s, a) - \alpha \log \pi(a | s)] + \lambda \left(1 - \sum_a \pi(a | s) \right). \quad (3.154)$$

Taking the derivative of \mathcal{L} with respect to $\pi(a | s)$ and setting it to zero:

$$\frac{\partial \mathcal{L}}{\partial \pi(a | s)} = Q^*(s, a) - \alpha(1 + \log \pi(a | s)) - \lambda = 0. \quad (3.155)$$

Solving for $\pi(a | s)$:

$$\log \pi(a | s) = \frac{Q^*(s, a) - \lambda - \alpha}{\alpha} \Rightarrow \pi(a | s) \propto \exp\left(\frac{Q^*(s, a)}{\alpha}\right). \quad (3.156)$$

The optimal policy is therefore the softmax over Q-values:

$$\pi^*(a | s) = \frac{\exp\left(\frac{Q^*(s, a)}{\alpha}\right)}{\sum_{a'} \exp\left(\frac{Q^*(s, a')}{\alpha}\right)}. \quad (3.157)$$

Plugging this back into the soft value function:

$$V^*(s) = \sum_a \pi^*(a | s) [Q^*(s, a) - \alpha \log \pi^*(a | s)]. \quad (3.158)$$

Since

$$\log \pi^*(a | s) = \frac{Q^*(s, a)}{\alpha} - \log \sum_{a'} \exp\left(\frac{Q^*(s, a')}{\alpha}\right), \quad (3.159)$$

we have

$$Q^*(s, a) - \alpha \log \pi^*(a | s) = \alpha \log \sum_{a'} \exp\left(\frac{Q^*(s, a')}{\alpha}\right). \quad (3.160)$$

Therefore, the optimal soft value function is given by

$$V^*(s) = \sum_a \pi^*(a | s) \cdot \alpha \log \sum_{a'} \exp\left(\frac{Q^*(s, a')}{\alpha}\right) = \alpha \log \sum_a \exp\left(\frac{Q^*(s, a)}{\alpha}\right). \quad (3.161)$$

3.6.8.3 Policy evaluation: general case

We now generalize this to the non-tabular case. We hold the policy parameters π fixed and optimize the parameters \mathbf{w} of the Q function by minimizing

$$J_Q(\mathbf{w}) = \mathbb{E}_{(\mathbf{s}_t, \mathbf{a}_t, r_{t+1}, \mathbf{s}_{t+1}) \sim \mathcal{D}} \left[\frac{1}{2} (Q_{\mathbf{w}}(\mathbf{s}_t, \mathbf{a}_t) - y(r_{t+1}, \mathbf{s}_{t+1}))^2 \right] \quad (3.162)$$

where \mathcal{D} is a replay buffer,

$$y(r_{t+1}, \mathbf{s}_{t+1}) = r_{t+1} + \gamma V_{\bar{\mathbf{w}}}(\mathbf{s}_{t+1}) \quad (3.163)$$

is the frozen target value, and $V_{\bar{w}}(s)$ is a frozen version of the soft value function from Equation (3.152):

$$V_{\bar{w}}(s_t) = \mathbb{E}_{\pi(a_t|s_t)} [Q_{\bar{w}}(s_t, a_t) - \alpha \log \pi(a_t|s_t)] \quad (3.164)$$

where \bar{w} is the EMA version of w . (The use of a frozen target is to avoid bootstrapping instabilities discussed in Section 2.5.2.5.)

To avoid the positive overestimation bias that can occur with actor-critic methods, [Haa+18a], suggest fitting two soft Q functions, by optimizing $J_Q(w_i)$, for $i = 1, 2$, independently. Inspired by clipped double Q learning, used in TD3 (Section 3.2.6.3), the targets are defined as

$$y(r_{t+1}, s_{t+1}; \bar{w}_{1:2}, \theta) = r_{t+1} + \gamma \left[\min_{i=1,2} Q_{\bar{w}_i}(s_{t+1}, \tilde{a}_{t+1}) - \alpha \log \pi_{\theta}(\tilde{a}_{t+1}|s_{t+1}) \right] \quad (3.165)$$

where $\tilde{a}_{t+1} \sim \pi_{\theta}(s_{t+1})$ is a sampled next action. In [Che+20], they propose the REDQ method (Section 2.5.3.3) which uses a random ensemble of $N \geq 2$ networks instead of just 2.

3.6.8.4 Policy improvement

In the policy improvement step, we derive the new policy based on the soft Q function by softmaxing over the possible actions for each state. We then project the update back on to the policy class Π :

$$\pi_{\text{new}} = \arg \min_{\pi' \in \Pi} D_{\text{KL}} \left(\pi'(\cdot|s_t) \parallel \frac{\exp(\frac{1}{\alpha} Q^{\pi_{\text{old}}}(s_t, \cdot))}{Z^{\pi_{\text{old}}}(s_t)} \right) \quad (3.166)$$

(The partition function $Z^{\pi_{\text{old}}}(s_t)$ may be intractable to compute for a continuous action space, but it cancels out when we take the derivative of the objective, so this is not a problem, as we show below.) After solving the above optimization problem, we are guaranteed to satisfy the soft policy improvement theorem, i.e., $Q^{\pi_{\text{new}}}(s_t, a_t) \geq Q^{\pi_{\text{old}}}(s_t, a_t)$ for all s_t and a_t .

We now generalize this to the non-tabular case. For policy improvement, we hold the value function parameters w fixed and optimize the parameters θ of the policy by minimizing the objective below, which is derived from the KL term by multiplying by α and dropping the constant Z term:

$$J_{\pi}(\theta) = \mathbb{E}_{s_t \sim \mathcal{D}} [\mathbb{E}_{a_t \sim \pi_{\theta}} [\alpha \log \pi_{\theta}(a_t|s_t) - Q_w(s_t, a_t)]] \quad (3.167)$$

Since we are taking gradients wrt θ , which affects the inner expectation term, we need to either use the REINFORCE estimator from Equation (3.26) or the **reparameterization trick** (see e.g., [Moh+20]). The latter is much lower variance, so is preferable.

To explain this in more detail, let us assume the policy distribution has the form $\pi_{\theta}(a_t|s_t) = \mathcal{N}(\mu_{\theta}(s_t), \sigma^2 \mathbf{I})$. We can write the random action as $a_t = f_{\theta}(s_t, \epsilon_t)$, where f is a deterministic function of the state and a noise variable ϵ_t , since $a_t = \mu(s_t) + \sigma^2 \epsilon_t$, where $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. The objective now becomes

$$J_{\pi}(\theta) = \mathbb{E}_{s_t \sim \mathcal{D}, \epsilon_t \sim \mathcal{N}} [\alpha \log \pi_{\theta}(f_{\theta}(s_t, \epsilon_t)|s_t) - Q_w(s_t, f_{\theta}(s_t, \epsilon_t))] \quad (3.168)$$

where we have replaced the expectation of a_t wrt π_{θ} with an expectation of ϵ_t wrt its noise distribution \mathcal{N} . Hence we can now safely take stochastic gradients. See Algorithm 11 for the pseudocode.

For discrete actions, we can replace the Gaussian reparameterization with the gumbel-softmax reparameterization [JGP16; MMT17]. Alternatively, we can eschew sampling and compute the expectations over the actions explicitly, to derive lower variance versions of the equations; this is known as **SAC-Discrete** [Chr19].

3.6.8.5 Adjusting the temperature

In [Haa+18b] they propose to automatically adjust the temperature parameter α by optimizing

$$J(\alpha) = \mathbb{E}_{s_t \sim \mathcal{D}, a_t \sim \pi_{\theta}} [-\alpha(\log \pi_{\theta}(a_t|s_t) + \bar{H})]$$

where \bar{H} is the target entropy (a hyper-parameter). This objective is approximated by sampling actions from the replay buffer.

Algorithm 11: SAC

```
1 Initialize environment state  $\mathbf{s}$ , policy parameters  $\boldsymbol{\theta}$ ,  $N$  critic parameters  $\mathbf{w}_i$ , target parameters  
    $\bar{\mathbf{w}}_i = \mathbf{w}_i$ , replay buffer  $\mathcal{D} = \emptyset$ , discount factor  $\gamma$ , EMA rate  $\rho$ , step size  $\eta_w$ ,  $\eta_\pi$ .  
2 repeat  
3   Take action  $\mathbf{a} \sim \pi_\theta(\cdot|\mathbf{s})$   
4    $(\mathbf{s}', r) = \text{step}(\mathbf{a}, \mathbf{s})$   
5    $\mathcal{D} := \mathcal{D} \cup \{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}')\}$   
6    $\mathbf{s} \leftarrow \mathbf{s}'$   
7   for  $G$  updates do  
8     Sample a minibatch  $\mathcal{B} = \{(\mathbf{s}_j, \mathbf{a}_j, r_j, \mathbf{s}'_j)\}$  from  $\mathcal{D}$   
9      $\mathbf{w} = \text{update-critics}(\boldsymbol{\theta}, \mathbf{w}, \mathcal{B})$   
10    Sample a minibatch  $\mathcal{B} = \{(\mathbf{s}_j, \mathbf{a}_j, r_j, \mathbf{s}'_j)\}$  from  $\mathcal{D}$   
11     $\boldsymbol{\theta} = \text{update-policy}(\boldsymbol{\theta}, \mathbf{w}, \mathcal{B})$   
12 until converged;  
13 .  
14 def update-critics( $\boldsymbol{\theta}, \mathbf{w}, \mathcal{B}$ ):  
15   Let  $(\mathbf{s}_j, \mathbf{a}_j, r_j, \mathbf{s}'_j)_{j=1}^B = \mathcal{B}$   
16    $y_j = y(r_j, \mathbf{s}'_j; \bar{\mathbf{w}}_{1:N}, \boldsymbol{\theta})$  for  $j = 1 : B$   
17   for  $i = 1 : N$  do  
18      $\mathcal{L}(\mathbf{w}_i) = \frac{1}{|\mathcal{B}|} \sum_{(\mathbf{s}, \mathbf{a}, r, \mathbf{s}')_j \in \mathcal{B}} (Q_{\mathbf{w}_i}(\mathbf{s}_j, \mathbf{a}_j) - \text{sg}(y_j))^2$   
19      $\mathbf{w}_i \leftarrow \mathbf{w}_i - \eta_w \nabla \mathcal{L}(\mathbf{w}_i)$  // Descent  
20      $\bar{\mathbf{w}}_i := \rho \bar{\mathbf{w}}_i + (1 - \rho) \mathbf{w}_i$  // Update target networks  
21   Return  $\mathbf{w}_{1:N}, \bar{\mathbf{w}}_{1:N}$   
22 .  
23 def update-actor( $\boldsymbol{\theta}, \mathbf{w}, \mathcal{B}$ ):  
24    $\hat{Q}(s, a) \triangleq \frac{1}{N} \sum_{i=1}^N Q_{\mathbf{w}_i}(s, a)$  // Average critic  
25    $J(\boldsymbol{\theta}) = \frac{1}{|\mathcal{B}|} \sum_{\mathbf{s} \in \mathcal{B}} \left( \hat{Q}(\mathbf{s}, \tilde{\mathbf{a}}_\theta(\mathbf{s})) - \alpha \log \pi_\theta(\tilde{\mathbf{a}}_\theta(\mathbf{s})|\mathbf{s}) \right)$ ,  $\tilde{\mathbf{a}}_\theta(\mathbf{s}) \sim \pi_\theta(\cdot|\mathbf{s})$   
26    $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} + \eta_\theta \nabla J(\boldsymbol{\theta})$  // Ascent  
27   Return  $\boldsymbol{\theta}$ 
```

3.6.9 Active inference

Control as inference is closely related to a technique known as **active inference**, as we explain below. For more details on the connection, see [Mil+20; WIP20; LÖW21; Saj+21; Tsc+20].

The active inference technique was developed in the neuroscience community, that has its own vocabulary for standard ML concepts. We start with the **free energy principle** [Fri09; Buc+17; SKM18; Ger19; Maz+22]. The FEP is equivalent to using variational inference to perform state estimation (perception) and parameter estimation (learning) in a latent variable model. In particular, consider an LVM $p(\mathbf{z}, \mathbf{o}|\boldsymbol{\theta})$ with hidden states \mathbf{z} , observations \mathbf{o} , and parameters $\boldsymbol{\theta}$. We define the variational free energy to be

$$\mathcal{F}(\mathbf{o}|\boldsymbol{\theta}) = \min_{q \in \mathcal{Q}} \mathcal{F}(\mathbf{o}|q, \boldsymbol{\theta}) \quad (3.169)$$

where

$$\mathcal{F}(\mathbf{o}|q, \boldsymbol{\theta}) = D_{\text{KL}}(q(\mathbf{z}|\mathbf{o}, \boldsymbol{\theta}) \parallel p(\mathbf{z}|\mathbf{o}, \boldsymbol{\theta})) - \log p(\mathbf{o}|\boldsymbol{\theta}) = \mathbb{E}_{q(\mathbf{z}|\mathbf{o}, \boldsymbol{\theta})} [\log q(\mathbf{z}|\mathbf{o}, \boldsymbol{\theta}) - \log p(\mathbf{o}, \mathbf{z}|\boldsymbol{\theta})] \geq -\log p(\mathbf{o}|\boldsymbol{\theta}) \quad (3.170)$$

which is the KL between the approximate variational posterior q and the true posterior p , minus a normalization constant, $\log p(\mathbf{o}|\boldsymbol{\theta})$, which is known as the free energy. State estimation (perception) corresponds to solving $\min_{q(\mathbf{z}|\mathbf{o}, \boldsymbol{\theta})} \mathcal{F}(\mathbf{o}|q, \boldsymbol{\theta})$, and parameter estimation (model fitting) corresponds to solving $\min_{\boldsymbol{\theta}} \mathcal{F}(\mathbf{o}|\boldsymbol{\theta})$, just as in the EM (expectation maximization) algorithm. (We can also be Bayesian about $\boldsymbol{\theta}$, as in variational Bayes EM, instead of just computing a point estimate.) This EM procedure will minimize the VFE, which is an upper bound on the negative log marginal likelihood of the data. In other words, it adjusts the model (belief state and parameters) so that it better predicts the observations, so the agent is less surprised (minimizes prediction errors). cf. [Ber+21].

To extend the above FEP to decision making problems, we define the **expected free energy** as follows

$$\mathcal{G}(\mathbf{a}) = \mathbb{E}_{q(\mathbf{o}|\mathbf{a})} [\mathcal{F}(\mathbf{o})] \quad (3.171)$$

$$= \underbrace{\mathbb{E}_{q(\mathbf{o}|\mathbf{a})} [D_{\text{KL}}(q(\mathbf{z}|\mathbf{o}) \parallel p(\mathbf{z}|\mathbf{o}))]}_{\mathcal{G}_{\text{epistemic}}(\mathbf{a})} - \underbrace{\mathbb{E}_{q(\mathbf{o}|\mathbf{a})} [\log p(\mathbf{o}|\boldsymbol{\theta})]}_{\mathcal{G}_{\text{extrinsic}}(\mathbf{a})} \quad (3.172)$$

where $q(\mathbf{o}|\mathbf{a})$ is the posterior predictive distribution over future observations given action sequence \mathbf{a} . (We should also condition on any observed history / agent state \mathbf{h} , and the model parameters $\boldsymbol{\theta}$, but we omit this from the notation for brevity.)

We see that we can decompose the EFE into two terms. First there is the **intrinsic value**, known as the **epistemic drive**. Minimizing this will encourage the agent to choose actions which maximize the mutual information between the observations \mathbf{o} and the hidden states \mathbf{z} , thus reducing uncertainty about the hidden states. (This is called **epistemic foraging**.) Second there is the **extrinsic value**, known as the **exploitation term**. Maximizing this will encourage the agent to choose actions that result in observations that match its prior. For example, if the agent predicts that the world will look brighter when it flips a light switch, it can take the action of flipping the switch to fulfill this prediction. This prior can be related to a reward function by defining as $p(\mathbf{o}) \propto e^{R(\mathbf{o})}$, encouraging goal directed behavior, exactly as in control-as-inference (cf. [Vri+25]). However, the active inference approach provides a way of choosing actions without needing to specify a reward.

Since solving to the optimal action at each step can be slow, it is possible to amortize this cost by training a policy network to compute $\pi(\mathbf{a}|\mathbf{h}) = \text{argmin}_{\mathbf{a}} \mathcal{G}(\mathbf{a}|\mathbf{h})$, where \mathbf{h} is the observation history (or current state), as shown in [Mil20; HL20]; this is called **“deep active inference”**.

Overall, we see that this framework provides a unified theory of both perception and action, both of which try to minimize some form of free energy. In particular, minimizing the expected free energy will cause the agent to pick actions to reduce its uncertainty about its hidden states, which can then be used to improve its predictive model $p_{\boldsymbol{\theta}}$ of observations; this in turn will help minimize the VFE of future observations, by updating the internal belief state $q(\mathbf{z}|\mathbf{o}, \boldsymbol{\theta})$ to explain the observations. In other words, the agent acts so it can learn so it becomes less surprised by what it sees. This ensures the agent is in **homeostasis** with its environment.

Note that active inference is often discussed in the context of **predictive coding**. This is equivalent to a special case of FEP where two assumptions are made: (1) the generative model $p(\mathbf{z}, \mathbf{o}|\boldsymbol{\theta})$ is a nonlinear hierarchical Gaussian model (similar to a VAE decoder), and (2) the variational posterior approximation uses a diagonal Laplace approximation, $q(\mathbf{z}|\mathbf{o}, \boldsymbol{\theta}) = \mathcal{N}(\mathbf{z}|\hat{\mathbf{z}}, \mathbf{H})$ with the mode $\hat{\mathbf{z}}$ being computed using gradient descent, and \mathbf{H} being the Hessian at the mode. This can be considered a non-amortized version of a VAE, where inference (E step) is done with iterated gradient descent, and parameter estimation (M step) is also done with gradient descent. (A more efficient incremental EM version of predictive coding, which updates $\{\hat{\mathbf{z}}_n : n = 1 : N\}$ and $\boldsymbol{\theta}$ in parallel, was recently presented in [Sal+24], and an amortized version in [Tsc+23].) For more details on predictive coding, see [RB99; Fri03; Spr17; Buc+17; HM20; MSB21; Mar21; OK22; Sal+23; Sal+24].

Chapter 4

Model-based RL

4.1 Introduction

Model-free approaches to RL typically need a lot of interactions with the environment to achieve good performance. For example, state of the art methods for the Atari benchmark, such as rainbow (Section 2.5.2.2), use millions of frames, equivalent to many days of playing at the standard frame rate. By contrast, humans can achieve the same performance in minutes [Tsi+17]. Similarly, OpenAI’s robot hand controller [And+20] needs 100 years of simulated data to learn to manipulate a rubiks cube.

One promising approach to greater sample efficiency is **model-based RL (MBRL)**. In the simplest approach to MBRL, we first learn the state transition or dynamics model $p_S(s'|s, a)$ — also called a **world model** (see e.g., [Zhu+24; Din+25; Sha+25]) — and the reward function $R(s, a)$, using some offline trajectory data, and then we use these models to compute a policy (e.g., using dynamic programming, as discussed in Section 2.2, or using some model-free policy learning method on simulated data, as discussed in Chapter 3). It can be shown that the sample complexity of learning the dynamics is less than the sample complexity of learning the policy [ZHR24].¹

The above two-stage approach — where we first learn the model, and then plan with it — can suffer from the usual problems encountered in offline RL (Section 7.7), i.e., the policy may query the model at a state for which no data has been collected, so predictions can be unreliable, causing the policy to learn the wrong thing. To get better results, we have to interleave the model learning and policy learning, so that one helps the other (since the policy determines what data is collected).

There are two main ways to perform MBRL. In the first approach, known as **decision-time planning** or **model predictive control**, we use the model to choose the next action by searching over possible future trajectories. We then score each trajectory, pick the action corresponding to the best one, take a step in the environment, and repeat. (We can also optionally update the model based on the rollouts.) This is discussed in Section 4.2.

The second approach is to use the current model and policy to rollout imaginary trajectories, and to use this data (optionally in addition to empirical data) to improve the policy using model-free RL; this is called **background planning**, and is discussed in Section 4.3.

The advantage of decision-time planning is that it allows us to train a world model on reward-free data, and then use that model to optimize any reward function. This can be particularly useful if the reward contains changing constraints, or if it is an intrinsic reward (Section 7.4) that frequently changes based on the knowledge state of the agent. The downside of decision-time planning is that it is much slower. However, it is possible to combine the two methods, as we discuss below. For an empirical comparison of background planning and decision-time planning, see [AP24].

Some generic pseudo-code for an MBRL agent is given in Algorithm 12. (The **rollout** function is defined in Algorithm 13; some simple code for model learning is shown in Algorithm 14, although we discuss

¹It can also be shown that learning an agent that can learn a policy to achieve any multi-step goal-conditioned task must be *implicitly* learning a world model [Ric+25]. This is potentially easier in practice, since it may be more robust to modeling error.

other loss functions in Section 4.4; finally, the code for the policy learning is given in other parts of this manuscript.) For more details on general MBRL, see e.g., [Wan+19; Moe+23; PKP21; Luo+22]. See also <https://github.com/alexzhang13/world-models-papers>.

Algorithm 12: MBRL agent

```

1 def MBRL-agent( $M_{\text{env}}; T, H, N$ ):
2   Initialize state  $s \sim M_{\text{env}}$ 
3   Initialize data buffer  $\mathcal{D} = \emptyset$ , model  $\hat{M}$ 
4   Initialize value function  $V$ , policy proposal  $\pi$ 
5   repeat
6     // Collect data from environment
7      $\tau_{\text{env}} = \text{rollout}(s, \pi, T, M_{\text{env}})$ ;
8      $s = \tau_{\text{env}}[-1]$ ;
9      $\mathcal{D} = \mathcal{D} \cup \tau_{\text{env}}$ 
10    // Update model
11    if Update model online then
12       $\hat{M} = \text{update-model}(\hat{M}, \tau_{\text{env}})$ 
13    if Update model using replay then
14       $\tau_{\text{replay}}^n = \text{sample-trajectory}(\mathcal{D}), n = 1 : N$ 
15       $\hat{M} = \text{update-model}(\hat{M}, \tau_{\text{replay}}^{1:N})$ 
16    // Update policy
17    if Update on-policy with real then
18       $(\pi, V) = \text{update-on-policy}(\pi, V, \tau_{\text{env}})$ 
19    if Update on-policy with imagination then
20       $\tau_{\text{imag}}^n = \text{rollout}(\text{sample-init-state}(\mathcal{D}), \pi, T, \hat{M}), n = 1 : N$ 
21       $(\pi, V) = \text{update-on-policy}(\pi, V, \tau_{\text{imag}}^{1:N})$ 
22    if Update off-policy with real then
23       $\tau_{\text{replay}}^n = \text{sample-trajectory}(\mathcal{D}), n = 1 : N$ 
24       $(\pi, V) = \text{update-off-policy}(\pi, V, \tau_{\text{replay}}^{1:N})$ 
25    if Update off-policy with imagination then
26       $\tau_{\text{imag}}^n = \text{rollout}(\text{sample-state}(\mathcal{D}), \pi, T, \hat{M}), n = 1 : N$ 
27       $(\pi, V) = \text{update-off-policy}(\pi, V, \tau_{\text{imag}}^{1:N})$ 
28  until until converged;
```

4.2 Decision-time (online) planning

In this section, we discuss how to choose the best action at each step based on planning forward from the current state using a known (or learned) world model. This is called **decision time planning** or “**planning in the now**” [KLP11], and is in contrast to methods that try to learn a policy which can be applied to all possible situations. In this section, we summarize some approaches to this problem. Our presentation is based in part on [KWW22, Ch. 9].

4.2.1 Receding horizon control

In **receding horizon control** or **RHC**, we plan from the current state s_t to a maximum fixed depth (horizon into the future) of d . We then take the first action a_t based on this future planning, observe the

Algorithm 13: Rollout

```
1 def rollout( $s_1, \pi, T, M$ )
2    $\tau = [s_1]$ 
3   for  $t = 1 : T$  do
4      $a_t = \pi(s_t)$ 
5      $(s_{t+1}, r_{t+1}) \sim M(s_t, a_t)$ 
6      $\tau \leftarrow [a_t, r_{t+1}, s_{t+1}]$ 
7   Return  $\tau$ 
```

Algorithm 14: Model learning

```
1 def update-model( $M, \tau^{1:N}$ ) :
2    $\ell(M) = -\frac{1}{NT} \sum_{n=1}^N \sum_{(s_t, a_t, r_{t+1}, s_{t+1}) \in \tau^n} \log M(s_{t+1}, r_{t+1} | s_t, a_t)$  // NLL
3    $M = M - \eta_M \nabla_M \ell(M)$ 
4   Return  $M$ 
```

new state s_{t+1} , and then replan. This approach can be quite slow, since it needs to perform a search or optimization procedure at each step. However, it can give good results, since it can choose an action that is tailored to the current state (and likely future), rather than relying on the generalization properties of a policy that was learned offline. In the sections below, we discuss various ways to implement this procedure.

4.2.1.1 Forward search

In **forward search**, we examine all possible transitions up to depth d by starting from the current state, and then considering all possible actions, and then considering all possible next states, etc. An example of the resulting **search tree** is given in Figure 4.1. We can compute the reward associated with each edge in the tree. At the leaves of the tree, we compute the remaining reward-to-go based on a utility or value function, $V(s)$, which can be learned offline using value-based methods. We then find the path with the highest score, and return the first action on this path. This process takes $O((|\mathcal{S}| \times |\mathcal{A}|)^d)$ time.

4.2.1.2 Branch and bound

In **branch and bound**, we try to avoid the exponential complexity of forward search by pruning paths that we determine are suboptimal. To do this, we need to know a lower bound on the value function, $\underline{V}(s)$, and an upper bound on the action value function, $\overline{Q}(s, a)$. At each state node s , we examine the actions in decreasing order of their upper bound. If we find an action a where $\overline{Q}(s, a)$ is less than the current best lower bound, we prune this branch of the tree, otherwise we expand it, and explore below. We continue this

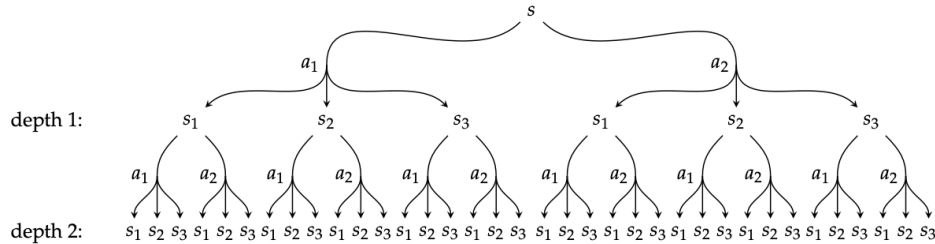


Figure 4.1: Illustration of forward search applied to a problem with 3 discrete states and 2 discrete actions. From Figure 9.1 of [KWW22]. Used with kind permission of Mykel Kochenderfer.

process until we hit a leaf node s (at the maximum depth), in which case we return the lower bound $\underline{V}(s)$. Depending on the tightness of the bounds, this approach can be significantly faster than forward search.

4.2.1.3 Sparse sampling

A simple way to speed up forward search (and branch and bound) is to sample a subset of m possible next states for each action. This is called **sparse sampling** [KM99]. The resulting complexity is $O((m \times |\mathcal{A}|)^d)$, which is independent of $|\mathcal{S}|$.

4.2.1.4 Heuristic search

In **heuristic search**, we start with a heuristic function $\bar{V}(s)$, which we use to initialize the value function $V(s)$. We then perform m Monte Carlo rollouts starting from the root node s . At each state node, we pick the greedy action wrt the current V , i.e., we choose $\arg\max_a R(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s')$. We then update $V(s) = \max_a R(s, a) + \gamma \sum_{s'} p(s'|s, a) V(s')$, and sample a next state $s' \sim p(s|s, a)$. We repeat this process until we hit the max depth. Finally we return the greedy action wrt V applied to the root node.

If the heuristic function is an upper bound on the optimal value function, then it is called an **admissible heuristic**. In this case, heuristic search is guaranteed to converge to the optimal value. The efficiency depends on the tightness of the upper bound, but in the worst case it is $O(m \times d \times |\mathcal{S}| \times |\mathcal{A}|)$.

4.2.2 Monte Carlo tree search (MCTS)

Monte Carlo tree search or **MCTS** is a receding horizon control procedure that works as follows (see e.g., [Mun14] for more details). Given the root node s_t , we perform m Monte Carlo rollouts to estimate $Q(s_t, a)$, and then we return the best action $\arg\max_a Q(s_t, a)$ or the action distribution $\text{softmax}(Q(s_t, a))$. To perform a rollout from a state s , we proceed as follows.

- **Action selection:** If we have not visited s before, we initialize the node by setting $N(s, a) = 0$ and $Q(s, a) = 0$ and returning $U(s)$ as the value, where U is some estimated value function. Otherwise we pick the next action to explore from state s . To explore actions, we first try each action once, and we then use the **Upper Confidence Tree** or **UCT** heuristic (based on UCB from Section 7.2.3) to select subsequent actions, i.e. we use

$$a = \arg\max_{a \in \mathcal{A}(s)} Q(s, a) + c \sqrt{\frac{\log N(s)}{N(s, a)}} \quad (4.1)$$

where $N(s) = \sum_a N(s, a)$ is the total visit count to s , and c is an exploration bonus scaling term. (Various other expressions are used in the literature, see [Bro+12] for a discussion.) If we have a predictor or prior over over actions, $P(s, a)$, we can instead use

$$a = \arg\max_{a \in \mathcal{A}(s)} Q(s, a) + c \left(P(s, a) \frac{\sqrt{N(s)}}{1 + N(s, a)} \right) \quad (4.2)$$

- **Expansion:** After choosing action a , we sample the next state $s' \sim p(s'|s, a)$.
- **Rollout:** we recursively estimate $u = U(s')$ using MCTS from that node. At some depth, we stop and use the value function to return $u = r + \gamma v(s')$.
- **Backup:** Finally we update the Q function for the root node using a running average:

$$Q(s, a) \leftarrow Q(s, a) + \frac{1}{N(s, a)} (u - Q(s, a)) \quad (4.3)$$

where the learning rate is given by $\frac{1}{N(s, a)}$. We also increment $N(s, a)$ by 1.

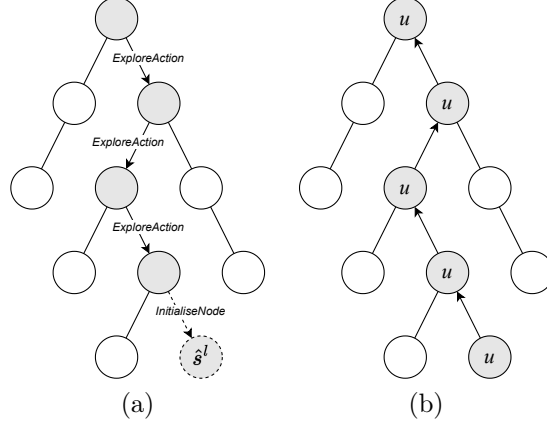


Figure 4.2: Illustration of MCTS. (a) Expanding nodes until we hit a new (previously unexplored) leaf node. (b) Propagating leaf value u back up the tree. From Figure 9.25 of [ACS24].

When we return from the recursive call, we are effectively backpropagating the value u from the leaves up the tree, as illustrated in Figure 4.2(b). A sketch of a non-recursive version of the algorithm (Algorithm 25 of [ACS24]) is shown in Algorithm 15.

Algorithm 15: Monte Carlo Tree Search (MCTS)

```

1 for  $t = 0, 1, 2, 3, \dots$  do
2   Observe current state  $s_t$ ;
3   for  $k$  simulations do
4      $\tau \leftarrow t$ ;
5      $\hat{s}_\tau \leftarrow s_t$ 
6     while  $\hat{s}_\tau$  is non-terminal and node  $\hat{s}_\tau$  exists in tree do
7        $\hat{a}_\tau \leftarrow \text{ExploreAction}(\hat{s}_\tau)$ ;
8        $\hat{s}_{\tau+1} \sim \mathcal{T}(\cdot \mid \hat{s}_\tau, \hat{a}_\tau)$ ;
9        $\hat{r}_\tau \leftarrow \mathcal{R}(\hat{s}_\tau, \hat{a}_\tau, \hat{s}_{\tau+1})$ ;
10       $\tau \leftarrow \tau + 1$ ;
11     if node  $\hat{s}_\tau$  does not exist in tree then
12       InitializeNode( $\hat{s}_\tau$ )
13     while  $\tau > t$  do
14       // Backpropagate
15        $\tau \leftarrow \tau - 1$ ;
16       Update( $Q, \hat{s}_\tau, \hat{a}_\tau$ );
17   // Select action for state  $s_t$ 
18    $\pi_t \leftarrow \text{BestAction}(s_t)$ ;
19    $a_t \sim \pi_t$ ;

```

4.2.2.1 MCTS for 2p0s games: AlphaGo, AlphaGoZero, and AlphaZero

MCTS can be applied to any kind of MDP, but some of its most famous applications are to games. We discuss general stochastic games in Chapter 5, but here focus on the special case of two-player, zero-sum symmetric games. In this case, the agent can model the opponent using its own policy, but with the roles reversed (this is known as **self-play**, see Section 5.3.5 for details). This lets the main player treat its opponent as part of the environment, thus creating a (non-stationary) single-agent problem.

In addition to choosing the next best action a_t (as in RHC), MCTS can be used to return a distribution

over good actions for the current state s ; we denote this by $\pi_s^{\text{MCTS}}(a) = [N(s, a) / (\sum_b N(s, b))]^{1/\tau}$, where τ is a temperature. This can be used as a target for policy improvement.

This method was used in the **AlphaGo** system of [Sil+16], which was the first AI system to beat a human grandmaster at the board game Go. AlphaGo was followed up by **AlphaGoZero** [Sil+17a], which had a much simpler design, and did not train on any human data, i.e., it was trained entirely using RL and self play. It significantly outperformed the original AlphaGo. This was generalized to **AlphaZero** [Sil+18], which can play expert-level Go, chess, and shogi (Japanese chess), without using any domain knowledge (except in the design of the neural network used to guide MCTS). (An open source version of AlphaZero, known as **Leela**, is available at <https://lczero.org/>.)

In more detail, AlphaZero used MCTS (with self-play), combined with a neural network which computes $(v^s, \pi^s) = f(s; \theta)$, where v^s is the expected outcome of the game from state s (either +1 for a win, -1 for a loss, or 0 for a draw), and π^s is the policy (distribution over actions) for state s . The policy is used internally by MCTS whenever a new node is initialized to give an additional exploration bonus to the most promising / likely actions. This controls the breadth of the search tree. In addition, the learned value function $v^s = f(s; \theta)_v$ is used to provide the value for leaf nodes in cases where we cannot afford to rollout to termination. This controls the depth of the search tree.

The policy/value network f is trained by optimizing the actor-critic loss

$$\mathcal{L}(\theta) = \mathbb{E}_{(s, \pi_s^{\text{MCTS}}, V^{\text{MCTS}}(s)) \sim \mathcal{D}} \left[(V^{\text{MCTS}}(s) - V_\theta(s))^2 - \sum_a \pi_s^{\text{MCTS}}(a) \log \pi_\theta(a|s) \right] \quad (4.4)$$

where $\mathcal{D} = \{(s, \pi_s^{\text{MCTS}}, V_s^{\text{MCTS}})\}$ is a dataset collected from MCTS rollouts starting at state s . These rollouts generate a distribution over actions at the root node s using $\pi_s^{\text{MCTS}}(a) = [N(s, a) / (\sum_b N(s, b))]^{1/\tau}$, where τ is a temperature. The rollouts also provide an estimate of $Q(s, a)$ for each visited (state, action) pair. From this we can estimate the non-parametric state-value function $V^{\text{MCTS}}(s) = \max_a Q^{\text{MCTS}}(s, a)$.

The above self-play approach trains an agent against the current version of itself, which can result in overfitting. To combat this, we can store multiple past versions of the policy, and then select any of these policies as a proxy for the opponent’s policy. This increases robustness of the main agent. Nevertheless, the resulting learned agent is susceptible to **adversarial attacks**; indeed, [Wan+23] showed that using this method, humans can beat “super-human” Go-playing AI agents.

4.2.2.2 MCTS with learned world model: MuZero and EfficientZero

AlphaZero and related methods assume the world model is known. The **MuZero** method of [Sch+20] learns a world model, by training a latent representation (embedding function) of the observations, $\mathbf{z}_t = e_\phi(\mathbf{o}_t)$, and a corresponding latent dynamics (and reward) model $(\mathbf{z}_t, r_t) = M_w(\mathbf{z}_t, a_t)$. The world model is trained to predict the immediate reward, the future reward (i.e, the value), and the optimal policy, where the optimal policy is computed using MCTS.

In more detail, we use MCTS to select action a_t , take a step, and add $(\mathbf{o}_t, a_t, r_t, \mathbf{o}_{t+1}, \pi_t^{\text{MCTS}}, V_t^{\text{MCTS}})$ to the replay buffer. To train the model, we augment the loss in Equation (4.4) by adding a term that measures how well the learned model predicts the observed rewards. Also, we now optimize this wrt the policy/value parameters θ as well as the model parameters w and embedding parameters ϕ :

$$\mathcal{L}(\theta, w, \phi) = \mathbb{E}_{(\mathbf{o}, a_t, r, \mathbf{o}', \pi_z^{\text{MCTS}}, V_z^{\text{MCTS}}) \sim \mathcal{D}} \left\{ (V^{\text{MCTS}}(z) - V_\theta(e_\phi(\mathbf{o}))^2 - \sum_a \pi_z^{\text{MCTS}}(a) \log \pi_\theta(a|e_\phi(\mathbf{o})) \right. \quad (4.5)$$

$$\left. + (r - M_w^r(e_\phi(\mathbf{o}), a_t))^2 \right\} \quad (4.6)$$

MuZero was applied to 3 perfect information board games (Go, Chess, and Shogi), as well as to Atari. The **Stochastic MuZero** method of [Ant+22] extends MuZero to allow for stochastic environments, such as the games 2048 and Backgammon. The **Sampled MuZero** method of [Hub+21] extends MuZero to allow for large and/or continuous action spaces. The **Gumbel MuZero** method of [Dan+22] proposes a better policy improvement algorithm, based on sampling actions without replacement, that improves sample

efficiency. The **MuZero Unplugged** method of [Sch+21b] proposes the **Reanalyse** algorithm, that applies an MCTS-based policy and value improvement operator, together with a learned world model, to offline trajectories.²

The **Efficient Zero** paper [Ye+21] extends MuZero by adding an additional self-prediction loss of the form $(z_{t+1} - M_w^z(z_t, a_t))^2$ to Equation (4.6) to help train the world model. (See Section 4.4.2.6 for further discussion of such losses.) It also makes several other changes. In particular, it replaces the empirical sum of instantaneous rewards, $\sum_{i=0}^{n-1} \gamma^i r_{t+i}$, used in computing V_t^{MCTS} , with an LSTM model that predicts the sum of rewards for a trajectory starting at z_t ; they call this the value prefix. In addition, it replaces the stored value at the leaf nodes of trajectories in the replay buffer with new values, by rerunning MCTS using the current model applied to the leaves, using the Reanalyse algorithm. They show that all three changes help, but the biggest gain is from the self-prediction loss. The recent **Efficient Zero V2** [Wan+24b] extends this to also work with continuous actions, by replacing tree search with sampling-based Gumbel search, amongst other changes.

4.2.2.3 MCTS in belief space

In [Mos+24], they present **BetaZero**, which performs MCTS in belief space. The current state is represented by a belief state, b_t , which is passed to the network to generate an initial policy proposal $\pi_\theta(a|b)$ and value function $v_\theta(b)$. (Instead of passing the belief state to the network, they actually pass features derived from the belief state, namely the mean and variance of the states.³)

To rollout out trajectories from the current root node b , they proceed as follows:

- Select an action $a \sim \pi_\theta(\cdot|b)$ using UCT heuristic.
- Expand the node as follows: sample the current hidden state $s \sim b$, sample the next hidden state $s' \sim T(s'|s, a)$ sample the observation $o \sim O(s')$, sample the reward $r \sim R(s, a, s')$; and finally derive the new belief state $b' = \text{Update}(b, a, o)$ using e.g., a particle filter (see e.g., [Lim+23]).
- Simulate future returns using rollouts to get $u = r + \gamma V_\theta(b')$ (assuming single step for notational simplicity);
- Backup the values using $Q(b, a) + = u$.

At the end of tree search, they derive the tree policy $\pi_t = \pi^{\text{MCTS}}(b_t)$ from the root, and compute the empirical reward-to-go $g_t = \sum_{i=t}^T \gamma^{i-t} r_i$ based on all rewards observed below root node b_t ; this is added to a dataset $\mathcal{D} = \{(b_t, \pi_t, g_t)\}$ which is used to update the policy and value network.

4.2.3 Sequential Monte Carlo (SMC) for online planning

Although MCTS is powerful, it is inherently serial, and can be complicated to apply to continuous action spaces. In this section, we discuss a more general method known as **SPO**, which stands for Sequential Monte Carlo Policy Optimisation [Mac+24].

SPO is based on the “RL as inference” framework, and is discussed in more detail in Section 3.6.4. In brief, the goal is to sample trajectories (sequences of states and actions) that are likely to be high scoring. That is, we want to sample from the following distribution

$$q(\tau) \propto d_q(s_0) \prod_{t \geq 0} \mathcal{T}(s_{t+1}|s_t, a_t) \pi(a_t|s_t, \theta_t) \exp\left(\frac{A(s_t, a_t)}{\eta}\right) \quad (4.7)$$

²Pseudocode for all these algorithm variants can be found at <https://www.julian.ac/about/>.

³The **POMCP** algorithm of [SV10] (Partially Observable Monte Carlo Planning) is related to BetaZero, but passes observation-action histories as input to the policy/value network, instead of features derived from the belief state. The **POMCPow** algorithm of [SK18] (POMCP with observation widening) extends this to continuous domains by sampling observations and actions.

where s_0 is the current state, A is the advantage function

$$A(s_t, a_t) = Q(s_t, a_t) - V(s_t) \approx r_t + V(s_{t+1}) - V(s_t) \quad (4.8)$$

and η is a temperature parameter obtained by maximizing Equation (3.145). (The state-value function V can be learned via TD(0).)

Let the resulting empirical distribution over trajectories be denoted by

$$\hat{q}_i(\tau) = \sum_{n=1}^N \bar{w}^n \delta(\tau - \tau^n) \quad (4.9)$$

where τ^n is the n 'th sample, and \bar{w}^n is its (normalized) weight. We can derive the distribution over next best action as follows:

$$\hat{q}(a|s_0) = \sum_n \bar{w}^n \delta(a - a_0^n) \quad (4.10)$$

One way to sample trajectories from such a distribution is to use **SMC** (Sequential Monte Carlo), which is a generalization of **particle filtering**. This is an approach to approximate inference in state space models based on sequential importance sampling with resampling (see e.g., [NLS19]). At each step, we use a proposal distribution $\beta(\tau_t|\tau_{1:t-1})$, which extends the previous sampled trajectory with a new value of $x_t = (s_t, a_t)$. We then compute the weight of this proposed extension by comparing it to the target $q_i(\tau_t|\tau_{1:t})$ to get

$$w(\tau_{1:t}) \propto w(\tau_{1:t-1}) \frac{q_i(\tau_{1:t})}{\beta(\tau_{1:t})} \quad (4.11)$$

Suppose we use the following proposal

$$\beta_i(\tau_t|\tau_{1:t-1}) \propto \hat{\mathcal{T}}(s_{t+1}|s_t, a_t) \pi(a_t|s_t, \theta_i) \quad (4.12)$$

Then the weight is given by

$$w(\tau_{1:t}) \propto w(\tau_{1:t-1}) \frac{\mathcal{T}(s_t|s_{t-1}, a_{t-1})}{\hat{\mathcal{T}}(s_t|s_{t-1}, a_{t-1})} \cdot \frac{\exp(A_i(s_t, a_t)/\eta_i^*) \pi(a_t|s_t, \theta_i)}{\pi(a_t|s_t, \theta_i)} \quad (4.13)$$

If we assume the learned model $\hat{\mathcal{T}}$ is accurate, this simplifies to

$$w(\tau_{1:t}) \propto w(\tau_{1:t-1}) \cdot \exp(A(s_t, a_t)/\eta) \quad (4.14)$$

In SMC, at each step we propose a new particle according to β , and then weight it according to the above equation. We can then optionally resample the particles every few steps, or when the effective sample size becomes too small; after a resampling step, we reset the weights to 1, since we now have a weighted sample. At the end, we return an empirical distribution over actions that correspond to high scoring trajectories, from which we can estimate the next best action (e.g., by taking the mean or mode of this distribution). See Algorithm 16 for details. (See also [Pic+19; Lio+22] for related methods.)

Note that the above framework is a special case of **twisted SMC** [NLS19; Law+22; Zha+24e], where the advantage function plays the role of a “twist” function, summarizing expected future rewards from the current state.

4.2.4 Model predictive control (MPC), aka open loop planning

In this section, we describe a method known as **model predictive control** (MPC), which is an **open loop** version of receding horizon control [MM90; CA13; RMD22]. (Connections with RL are discussed in e.g., [Ber24; Law+25].) In particular, at each step, it solves for the sequence of subsequent actions that is most likely to achieve high expected reward:

$$\mathbf{a}_{t:t+d}^* = \operatorname{argmax}_{\mathbf{a}_{t:t+d}} \mathbb{E}_{s_{t+1:t+d} \sim \mathcal{T}(\cdot|s_t, \mathbf{a}_{t:t+d})} \left[\sum_{h=0}^d R(s_{t+h}, a_{t+h}) + \hat{V}(s_{t+d+1}) \right] \quad (4.15)$$

Algorithm 16: SMC-RHC (Sequential Monte Carlo for Receding Horizon Control)

```

1 def SMC-RHC( $s_t, \pi_i, V_i$ ):
2   Initialize particles:  $\{\mathbf{s}_t^n = \mathbf{s}_t\}_{n=1}^N$ 
3   Initialize weights:  $\{w_t^n = 1\}_{n=1}^N$ 
4   for  $j = t + 1 : t + d$  do
5      $\{a_j^n \sim \pi_i(\cdot | s_j^n)\}_{n=1}^N$ 
6      $\{s_{j+1}^n \sim \hat{\mathcal{T}}(s_j^n, a_j^n)\}_{n=1}^N$ 
7      $\{r_j^n \sim \hat{\mathcal{R}}(s_j^n, a_j^n)\}_{n=1}^N$ 
8      $\{x_j^n = (s_j^n, a_j^n, r_j^n)\}_{n=1}^N$ 
9      $\{A_j^n = r_j^n + V_i(s_{j+1}^n) - V_i(s_j^n)\}_{n=1}^N$ 
10     $\{w_j^n = w_{j-1}^n \exp(A_j^n / \eta_i^*)\}_{n=1}^N$ 
11    if Resample then
12       $\{\mathbf{x}_{t,j}^n\} \sim \text{Multinom}(n; w_i^1, \dots, w_i^N)$ 
13       $\{w_j^n = 1\}_{n=1}^N$ 
14     $\{\bar{w}^n = \frac{w^n}{\sum_{n'} w^{n'}}\}_{n=1}^N$ 
15    Let  $\{a_t^n\}_{n=1}^N$  be the set of sampled actions at the start of  $\{x_{t:t+d}^n\}_{n=1}^N$ 
16    Return  $\hat{q}(a | s_t) = \sum_n \bar{w}^n \delta(a - a_t^n)$ 

```

where \mathcal{T} is the dynamics model. It then returns a_t^* as the best action, takes a step, and replans.

Crucially, the future actions are chosen without knowing what the future states are; this is what is meant by “open loop”. This can be much faster than interleaving the search for actions and future states. However, it can also lead to suboptimal decisions, as we discuss below. Nevertheless, the fact that we replan at each step can reduce the harms of this approximation, making the method quite popular for some problems, especially ones where the dynamics are deterministic, and the actions are continuous (so that Equation (4.15) becomes a standard optimization problem over the real valued sequence of vectors $\mathbf{a}_{t:t+d-1}$).

4.2.4.1 Suboptimality of open-loop planning for stochastic environments

Consider the example in Figure 4.3, where there are 9 states, 2 actions (going up or down), and the planning horizon is $d = 2$. All transitions are deterministic, except that going up from s_1 can either end up in s_2 wp 0.5 or in s_3 wp 0.5.

There are 4 open-loop plans, with the following expected utilities:

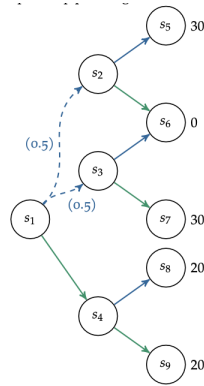


Figure 4.3: Illustration of the suboptimality of open-loop planning. From Figure 9.9 of [KWW22]. Used with kind permission of Mykel Kochenderfer.

- $U(\text{up}, \text{up}) = 0.5 \times 30 + 0.5 \times 0 = 15$
- $U(\text{up}, \text{down}) = 0.5 \times 0 + 0.5 \times 30 = 15$
- $U(\text{down}, \text{up}) = 20$
- $U(\text{down}, \text{down}) = 20$

Thus the best open-loop action is to choose down, with an expected reward of 20. However, closed-loop planning can reason that, after taking the first action, the agent can sense the resulting state. If it initially chooses to go up from s_1 , then it can decide to next go up or down, depending on whether it is in s_2 or s_3 , thereby guaranteeing a reward of 30.

4.2.4.2 Trajectory optimization

If the dynamics is deterministic, the problem becomes one of solving

$$\max_{a_{1:d}, s_{2:d}} \sum_{t=1}^d \gamma^t R(s_t, a_t) \quad (4.16)$$

$$\text{s.t. } s_{t+1} = \mathcal{T}(s_t, a_t) \quad (4.17)$$

where \mathcal{T} is the transition function. This is called a **trajectory optimization** problem. We discuss various ways to solve this below.

4.2.4.3 LQR

If the system dynamics are linear and the reward function is quadratic, then the optimal action sequence can be computed exactly using a method similar to Kalman filtering. This is known as the **linear quadratic regulator (LQR)**. For details, see e.g., [AM89; HR17; Pet08].

If the model is nonlinear, we can use **differential dynamic programming (DDP)** [JM70; TL05] to approximately solve the problem. In each iteration, DDP starts with a reference trajectory, and linearizes the system dynamics around states on the trajectory to form a locally quadratic approximation of the reward function. This system can be solved using LQG, whose optimal solution results in a new trajectory. The algorithm then moves to the next iteration, with the new trajectory as the reference trajectory.

4.2.4.4 Random shooting

For general nonlinear models (such as neural networks), a simple approach is to pick a sequence of random actions to try (from some proposal policy), evaluate the reward for each trajectory, and pick the best. This is called **random shooting** [Die+07; Rao10].

4.2.4.5 CEM

As an improvement upon random shooting, it is common to use black-box (gradient-free) optimization methods like the **cross-entropy method** or **CEM** in order to find the best action sequence. The CEM method is a simple derivative-free optimization method for continuous black-box functions $f : \mathbb{R}^D \rightarrow \mathbb{R}$. We start with a multivariate Gaussian, $\mathcal{N}(\mu_0, \Sigma_0)$, representing a distribution over possible action \mathbf{a} . We sample from this, evaluate all the proposals, pick the top K , then refit the Gaussian to these top K , and repeat, until we find a sample with sufficiently good score (or we perform moment matching on the top K scores). For details, see [Rub97; RK04; Boe+05]. In Section 4.2.4.6, we discuss the MPPI method, which is a common instantiation of CEM method. In [BXS20] they discuss how to combine CEM with gradient-based planning.

4.2.4.6 MPPI

The **model predictive path integral** or **MPPI** approach [WAT17] is a version of CEM. Originally MPPI was limited to models with linear dynamics, but it was extended to general nonlinear models in [Wil+17]. The basic idea is that the initial mean of the Gaussian at step t , namely $\mu_t = \mathbf{a}_{t:t+H}$, is computed based on shifting $\hat{\mu}_{t-1}$ forward by one step. (Here μ_t is known as a reference trajectory.)

In [Wag+19], they apply this method for robot control. They consider a state vector of the form $\mathbf{s}_t = (\mathbf{q}_t, \dot{\mathbf{q}}_t)$, where \mathbf{q}_t is the configuration of the robot. The deterministic dynamics has the form

$$\mathbf{s}_{t+1} = F(\mathbf{s}_t, \mathbf{a}_t) = \begin{pmatrix} \mathbf{q}_t + \dot{\mathbf{q}}_t \Delta t \\ \dot{\mathbf{q}}_t + f(\mathbf{s}_t, \mathbf{a}_t) \Delta t \end{pmatrix} \quad (4.18)$$

where f is a 2 layer MLP. This is trained using the **Dagger** method of [RGB11], which alternates between fitting the model (using supervised learning) on the current replay buffer (initialized with expert data), and then deploying the model inside the MPPI framework to collect new data (see Section 7.7.4.2 for details).

A similar method was used in the TD-MPC paper [HSW22; HSW24], which learns a non-generative world model in latent space, and then uses MPPI to implement MPC (see Section 4.4.2.12 for details). They initialize the population of K sampled action trajectories by applying the policy prior to generate $J < K$ samples, and then generate the remaining $K - J$ samples using the diagonal Gaussian prior from the previous time step.

4.2.4.7 GP-MPC

[KD18] proposed **GP-MPC**, which combines a Gaussian process dynamics model with model predictive control. They compute a Gaussian approximation to the future state trajectory given a candidate action trajectory, $p(\mathbf{s}_{t+1:t+H} | \mathbf{a}_{t:t+H-1}, \mathbf{s}_t)$, by moment matching, and use this to deterministically compute the expected reward and its gradient wrt $\mathbf{a}_{t:t+H-1}$. Using this, they can solve Equation (4.15) to find $\mathbf{a}_{t:t+H-1}^*$; finally, they execute the first step of this plan, \mathbf{a}_t^* , and repeat the whole process.

The key observation is that moment matching is a deterministic operator that maps $p(\mathbf{s}_t | \mathbf{a}_{1:t-1})$ to $p(\mathbf{s}_{t+1} | \mathbf{a}_{1:t})$, so the problem becomes one of deterministic optimal control, for which many solution methods exist. Indeed the whole approach can be seen as a generalization of the **LQG** method from classical control, which assumes a (locally) linear dynamics model, a quadratic cost function, and a Gaussian distribution over states [Rec19]. In GP-MPC, the moment matching plays the role of local linearization.

The advantage of GP-MPC over the earlier method known as **PILCO** (“probabilistic inference for learning control”), which learns a policy by maximizing the expected reward from rollouts (see [DR11; DFR15] for details), is that GP-MPC can handle constraints more easily, and it can be more data efficient, since it continually updates the GP model after every step (instead of at the end of an trajectory).

4.3 Background (offline) planning

In Section 4.2, we discussed how to use models to perform decision-time planning. However, this can be slow. Fortunately, we can amortize the planning process into a reactive policy. To do this, we can use the model to generate synthetic trajectories “in the background” (while executing the current policy), and use this imaginary data to train the policy; this is called “**background planning**”. We discuss a game theoretic formulation of this setup in Section 4.3.1. Then in Section 4.3.2, we discuss ways to combine model-based and model-free learning. Finally, in Section 4.4.4, we discuss ways to deal with model errors, that might lead the policy astray.

4.3.1 A game-theoretic perspective on MBRL

In MBRL, we optimize the policy wrt a possibly incorrect model, and then update the model given new data, in an iterative fashion. In this section, we discuss a game-theoretic framework for MBRL, as proposed in

[RMK20], which justifies why this approach should converge. (See Chapter 5 for a more general discussion of game theory and RL.)

We denote the true world model by M_{env} . To simplify the notation, we assume an MDP setup with a known reward function, so all that needs to be learned is the world model, \hat{M} , representing $p(s'|s, a)$. (It is trivial to also learn the reward function.) We define the value of a policy π when rolled out in some model M' as the (discounted) sum of expected rewards:

$$J(\pi, M') = \mathbb{E}_{M', \pi} \left[\sum_{t=0}^{\infty} \gamma^t R(s_t) \right]$$

We define the loss of a model \hat{M} given a distribution $\mu(s, a)$ over states and actions as

$$\ell(\hat{M}, \mu) = \mathbb{E}_{(s,a) \sim \mu} \left[D_{\text{KL}} \left(M_{\text{env}}(\cdot|s, a) \parallel \hat{M}(\cdot|s, a) \right) \right]$$

We now define MBRL as a two-player general-sum game (see Chapter 5 for details):

$$\overbrace{\max_{\pi} J(\pi, \hat{M})}^{\text{policy player}}, \overbrace{\min_{\hat{M}} \ell(\hat{M}, \mu_{M_{\text{env}}}^{\pi})}^{\text{model player}}$$

where $\mu_{M_{\text{env}}}^{\pi} = \frac{1}{T} \sum_{t=0}^T M_{\text{env}}(s_t = s, a_t = a)$ is the induced state visitation distribution when policy π is applied in the real world M_{env} , so that minimizing $\ell(\hat{M}, \mu_{M_{\text{env}}}^{\pi})$ gives the **maximum likelihood estimate** for \hat{M} .

Now consider a **Nash equilibrium** (Section 5.2.4) of this game, that is a pair (π, \hat{M}) that satisfies $\ell(\hat{M}, \mu_{M_{\text{env}}}^{\pi}) \leq \epsilon_{M_{\text{env}}}$ and $J(\pi, \hat{M}) \geq J(\pi', \hat{M}) - \epsilon_{\pi}$ for all π' . (That is, the model is accurate when predicting the rollouts from π , and π cannot be improved when evaluated in \hat{M} .) In [RMK20] they prove that the Nash equilibrium policy π is near optimal wrt the real world, in the sense that $J(\pi^*, M_{\text{env}}) - J(\pi, M_{\text{env}})$ is bounded by a constant, where π^* is an optimal policy for the real world M_{env} . (The constant depends on the ϵ parameters, and the TV distance between $\mu_{M_{\text{env}}}^{\pi^*}$ and $\mu_{\hat{M}}^{\pi^*}$.)

A natural approach to trying to find such a Nash equilibrium is to use **gradient descent ascent** or **GDA**, in which each player updates its parameters simultaneously, using

$$\begin{aligned} \pi_{k+1} &= \pi_k + \eta_{\pi} \nabla_{\pi} J(\pi_k, \hat{M}_k) \\ \hat{M}_{k+1} &= \hat{M}_k - \eta_M \nabla_{\hat{M}} \ell(\hat{M}_k, \mu_{M_{\text{env}}}^{\pi_k}) \end{aligned}$$

Unfortunately, GDA is often an unstable algorithm, and often needs very small learning rates η . In addition, to increase sample efficiency in the real world, it is better to make multiple policy improvement steps using synthetic data from the model \hat{M}_k at each step.

Rather than taking small steps in parallel, the **best response** strategy fully optimizes each player given the previous value of the other player, in parallel:

$$\begin{aligned} \pi_{k+1} &= \arg\max_{\pi} J(\pi, \hat{M}_k) \\ \hat{M}_{k+1} &= \arg\min_{\hat{M}} \ell(\hat{M}, \mu_{M_{\text{env}}}^{\pi_k}) \end{aligned}$$

Unfortunately, making such large updates in parallel can often result in a very unstable algorithm.

To avoid the above problems, [RMK20] propose to replace the min-max game with a **Stackelberg game**, which is a generalization of min-max games where we impose a specific player ordering (Section 5.2.13). In particular, let the players be A and B , let their parameters be θ_A and θ_B , and let their losses be $\mathcal{L}_A(\theta_A, \theta_B)$ and $\mathcal{L}_B(\theta_A, \theta_B)$. If player A is the leader, the Stackelberg game corresponds to the following **nested optimization problem**, also called a **bilevel optimization problem**:

$$\min_{\theta_A} \mathcal{L}_A(\theta_A, \theta_B^*(\theta_A)) \quad \text{s.t.} \quad \theta_B^*(\theta_A) = \arg\min_{\theta_B} \mathcal{L}_B(\theta_A, \theta_B)$$

Since the follower B chooses the best response to the leader A , the follower’s parameters are a function of the leader’s. The leader is aware of this, and can utilize this when updating its own parameters.

The main advantage of the Stackelberg approach is that one can derive gradient-based algorithms that will provably converge to a local optimum [CMS07; ZS22; GP23]. In particular, suppose we choose the **policy as leader (PAL)**. We then just have to solve the following optimization problem:

$$\begin{aligned}\hat{M}_{k+1} &= \underset{\hat{M}}{\operatorname{argmin}} \ell(\hat{M}, \mu_{M_{\text{env}}}^{\pi_k}) \\ \pi_{k+1} &= \pi_k + \eta_{\pi} \nabla_{\pi} J(\pi_k, \hat{M}_{k+1})\end{aligned}$$

We can solve the first step by executing π_k in the environment to collect data \mathcal{D}_k and then fitting a local (policy-specific) dynamics model by solving $\hat{M}_{k+1} = \underset{\hat{M}}{\operatorname{argmin}} \ell(\hat{M}, \mathcal{D}_k)$. (For example, this could be a locally linear model, such as those used in trajectory optimization methods discussed in Section 4.2.4.6.) We then (slightly) improve the policy to get π_{k+1} using a conservative update algorithm, such as natural actor-critic (Section 3.2.4) or TRPO (Section 3.3.2), on “imaginary” model rollouts from \hat{M}_{k+1} .

Alternatively, suppose we choose the **model as leader (MAL)**. We now have to solve

$$\begin{aligned}\pi_{k+1} &= \underset{\pi}{\operatorname{argmax}} J(\pi, \hat{M}_k) \\ \hat{M}_{k+1} &= \hat{M}_k - \eta_M \nabla_{\hat{M}} \ell(\hat{M}, \mu_{M_{\text{env}}}^{\pi_{k+1}})\end{aligned}$$

We can solve the first step by using any RL algorithm on “imaginary” model rollouts from \hat{M}_k to get π_{k+1} . We then apply this in the real world to get data \mathcal{D}_{k+1} , which we use to slightly improve the model to get \hat{M}_{k+1} by using a conservative model update applied to \mathcal{D}_{k+1} . (In practice we can implement a conservative model update by mixing \mathcal{D}_{k+1} with data generated from earlier models, an approach known as **data aggregation** [RB12].) Compared to PAL, the resulting model will be a more global model, since it is trained on data from a mixture of policies (including very suboptimal ones at the beginning of learning).

4.3.2 Dyna

The **Dyna** paper [Sut90] proposed an approach to MBRL that is related to the approach discussed in Section 4.3.1, in the sense that it trains a policy and world model in parallel, but it differs in one crucial way: the policy is also trained on real data, not just imaginary data. That is, we define $\pi_{k+1} = \pi_k + \eta_{\pi} \nabla_{\pi} J(\pi_k, \hat{\mathcal{D}}_k \cup \mathcal{D}_k)$, where \mathcal{D}_k is data from the real environment and $\hat{\mathcal{D}}_k = \text{rollout}(\pi_k, \hat{M}_k)$ is imaginary data from the model. This makes Dyna a hybrid model-free and model-based RL method, rather than a “pure” MBRL method.

In more detail, at each step of Dyna, the agent collects new data from the environment and adds it to a real replay buffer. This is then used to do an off-policy update. It also updates its world model given the real data. Then it simulates imaginary data, starting from a previously visited state (see **sample-init-state** function in Algorithm 12), and rolling out the current policy in the learned model. The imaginary data is then added to the imaginary replay buffer and used by an on-policy learning algorithm. This process continues until the agent runs out of time and must take the next step in the environment. See Figure 4.4 for an illustration.

4.3.2.1 Tabular Dyna

The original Dyna paper was developed under the assumption that the world model $s' = M(s, a)$ is deterministic and tabular, and the Q function is also tabular. See Algorithm 17 for the simplified pseudocode for this case. Since we assume a deterministic world model of the form $s' = M(s, a)$, then sampling a single step from this starting at a previously visited state is equivalent to experience replay (Section 2.5.2.3). Thus we can think of ER as a kind of non-parametric world model [HHA19].

4.3.2.2 Dyna with function approximation

It is easy to extend Dyna to work with function approximation and policy gradient methods. The code is identical to the MBRL code in Algorithm 12, where now we train the policy on real as well as imaginary data.

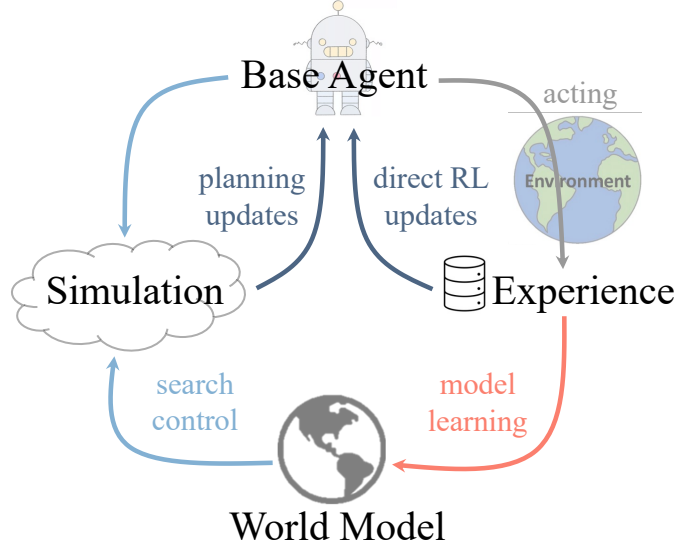


Figure 4.4: Illustration of Dyna. From Figure 2 of [Liu+24]. Used with kind permission of Zichen Liu.

Algorithm 17: Tabular Dyna-Q

```

1 def dyna-Q-agent( $s, M_{\text{env}}; \epsilon, \eta, \gamma$ ):
2   Initialize data buffer  $\mathcal{D} = \emptyset$ ,  $Q(s, a) = 0$  and  $\hat{M}(s, a) = 0$ 
3   repeat
4     // Collect real data from environment
5      $a = \text{eps-greedy}(Q, \epsilon)$ 
6      $(s', r) = \text{env.step}(s, a)$ 
7      $\mathcal{D} = \mathcal{D} \cup \{(s, a, r, s')\}$ 
8     // Update policy on real data
9      $Q(s, a) := Q(s, a) + \eta[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$ 
10    // Update model on real data
11     $\hat{M}(s, a) = (s', r)$ 
12     $s := s'$ 
13    // Update policy on imaginary data
14    for  $n=1:N$  do
15      Select  $(s, a)$  from  $\mathcal{D}$ 
16       $(s', r) = \hat{M}(s, a)$ 
17       $Q(s, a) := Q(s, a) + \eta[r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$ 
18  until until converged;

```

([Lai+21] argues that we should gradually increase the fraction of real data that is used to train the policy, to avoid suboptimal performance due to model limitations.) If we use real data from the replay buffer, we have to use an off-policy learner, since the replay buffer contains trajectories that may have been generated from old policies. (The most recent real trajectory, and all imaginary trajectories, are always from the current policy.)

We now mention some examples of this “generalized Dyna” framework. In [Sut+08] they extended Dyna to the case where the Q function is linear, and in [HTB18] they extended it to the DQN case. In [Jan+19a], they present the **MBPO** (model based policy optimization) algorithm, which uses Dyna with the off-policy SAC method. Their world model is an **ensemble of DNNs**, which generates diverse predictions (an approach which was originally proposed in the **PETS** (probabilistic ensembles with trajectory sampling) paper of [Chu+18]).⁴ In [Kur+19], they combine Dyna with TRPO (Section 3.3.2) and ensemble world models, and in [Wu+23] they combine Dyna with PPO and GP world models. (Technically speaking, these on-policy approaches are not valid with Dyna, but they can work if the replay buffer used for policy training is not too stale.)

4.4 World models

	Background planning	Online planning	Exploration
Observation prediction	Dyna, DreamerV3, IRIS, Delta-IRIS, Diamond	CEM: PlaNet Rnd shooting: TDM	SPR
Value + self prediction	DreamingV2, AIS	MCTS: MuZero, EfficientZero CEM: TD-MPC	BYOL-Explore

Table 4.1: Comparison of different world model-based methods.

In this section, we discuss various kinds of world models that have been proposed in the literature. These models can be trained to predict future observations (generative WMs) or just future rewards/values and/or future latent embeddings (non-generative / non-reconstructive WMs). Once trained, the models can be used for decision-time planning, background planning, or just as an auxiliary signal to aid in things like intrinsic curiosity. See Table 4.1 for a summary.

4.4.1 World models which are trained to predict observation targets

In this section, we discuss different kinds of world model $\mathcal{T}(\mathbf{s}'|\mathbf{s}, a)$. We can use this to generate imaginary trajectories by sampling from the following joint distribution:

$$p(\mathbf{s}_{t+1:T}, \mathbf{r}_{t+1:T}, \mathbf{a}_{t:T-1} | \mathbf{s}_t) = \prod_{i=t}^{T-1} \pi(\mathbf{a}_i | \mathbf{s}_i) \mathcal{T}(\mathbf{s}_{i+1} | \mathbf{s}_i, \mathbf{a}_i) R(r_{i+1} | \mathbf{s}_i, \mathbf{a}_i) \quad (4.19)$$

The model may be augmented with latent variables, as we discuss in Section 4.4.1.2.

If the state space is high dimensional (e.g., images), then we denote the observable data by \mathbf{o} . We can then learn $\mathcal{T}(\mathbf{o}'|\mathbf{o}, a)$ using standard techniques for conditional image generation such as diffusion (see e.g., the **Diamond** method of [Alo+24], the **Genie2** method of [al24], the **GAIA-1** model of [Hu+23a], etc. This kind of world model is equivalent to an action-conditional version of a **video generative model**, such as **Sora**, **Veo-3**, **seedance** [Gao+25], etc.

⁴In [Zhe+22b] they argue that the main benefit of an ensemble is that it limits the Lipschitz constant of the value function. They show that more direct methods for regularizing this can work just as well, and are much faster.

Note that these methods are trained to predict the entire observation vector, even if we use latent variables. (This is what we mean by “generative world model”.) One big disadvantage of this approach is that the observations may contain irrelevant or distractor variables, that are not necessary for task performance. In addition, such models are often slow to use slow, and there may be a distribution shift in the observation process between train and test time. Both of these factors can adversely affect the performance of generative WMs (see e.g., [Tom+23]). We discuss some non-generative approaches to WMs in Section 4.4.2.

4.4.1.1 Generative world models without latent variables

The simplest approach is to define $\mathcal{T}(\mathbf{o}'|\mathbf{o}, a)$ as a conditional generative model over states. If the observed states are low-dimensional vectors, such as proprioceptive states, we can use transformers (see e.g., the **Transformer Dynamics Model** of [Sch+23a]).

In some cases, the dimensions of the state vector \mathbf{s} represent distinct variables, and the joint Markov transition matrix $p(\mathbf{s}'|\mathbf{s}, a)$ has conditional independence properties which can be represented as a sparse graphical model. This is called a **factored MDP** [BDG00].

4.4.1.2 Generative world models with latent variables

In this section, we describe some methods that use latent variables as part of their world model. This can improve the speed of generating imaginary futures, and can provide a compact latent space as input to a policy.

We let \mathbf{z}_t denote the latent (or hidden) state at time t ; this can be a discrete or continuous variable (or vector of variables). The generative model has the form of a controlled HMM:

$$p(\mathbf{o}_{t+1:T}, \mathbf{z}_{t+1:T}, \mathbf{r}_{t+1:T}, \mathbf{a}_{t:T-1}|\mathbf{z}_t) = \prod_{i=t}^{T-1} \pi(\mathbf{a}_i|\mathbf{z}_i) \mathcal{M}(\mathbf{z}_{i+1}|\mathbf{z}_i, \mathbf{a}_i) R(r_i|\mathbf{z}_{i+1}, \mathbf{a}_i) D(\mathbf{o}_i|\mathbf{z}_i) \quad (4.20)$$

where $p(\mathbf{o}_t|\mathbf{z}_t) = D(\mathbf{o}_t|\mathbf{z}_t)$ is the decoder or likelihood function, $\mathcal{M}(\mathbf{z}'|\mathbf{z}, \mathbf{a})$ is the dynamics in latent space. $\pi(\mathbf{a}_t|\mathbf{z}_t)$ is the policy in latent space.

The world model is usually trained by maximizing the marginal likelihood of the observed outputs given an action sequence. (We discuss non-likelihood based loss functions in Section 4.4.2.) Computing the marginal likelihood requires marginalizing over the hidden variables $\mathbf{z}_{t+1:T}$. To make this computationally tractable, it is common to use amortized variational inference, in which we train an encoder network, $p(\mathbf{z}_t|\mathbf{o}_t) = E(\mathbf{z}_t|\mathbf{o}_t)$, to approximate the posterior over the latents. Many papers have followed this basic approach, such as the “world models” paper [HS18], and the methods we discuss below.

4.4.1.3 Example: Dreamer

In this section, we summarize the approach used in **Dreamer** paper [Haf+20] and its recent extensions, such as DreamerV2 [Haf+21], DreamerV3 [Haf+25] and DreamerV4 [HYL25]. These are all based on the background planning approach, in which the policy is trained on imaginary trajectories generated by a latent variable world model. (Note that Dreamer is based on an earlier approach called **PlaNet** [Haf+19], which used MPC instead of background planning.)

In Dreamer, the stochastic dynamic latent variables in Equation (4.20) are replaced by deterministic dynamic latent variables \mathbf{h}_t , since this makes the model easier to train. (We will see that \mathbf{h}_t acts like the posterior over the hidden state at time $t - 1$; this is also the prior predictive belief state before we see \mathbf{o}_t .) A static stochastic variable \mathbf{z}_t is now generated for each time step, and acts like a “random effect” in order to help generate the observations, without relying on \mathbf{h}_t to store all of the necessary information. (This simplifies the recurrent latent state.) In more detail, Dreamer defines the following functions:⁵

- A hidden dynamics (sequence) model: $\mathbf{h}_{t+1} = \mathcal{U}(\mathbf{h}_t, \mathbf{a}_t, \mathbf{z}_t)$

⁵We can map from our notation to the notation in the paper as follows: $\mathbf{o}_t \rightarrow x_t$, $U \rightarrow f_\phi$ (sequence model), $P_0 \rightarrow p_\phi(\hat{z}_t|h_t)$ (dynamics predictor), $D \rightarrow p_\phi(\hat{x}_t|h_t, \hat{z}_t)$ (decoder), $E \rightarrow q_\phi(\mathbf{z}_t|h_t, x_t)$ (encoder).

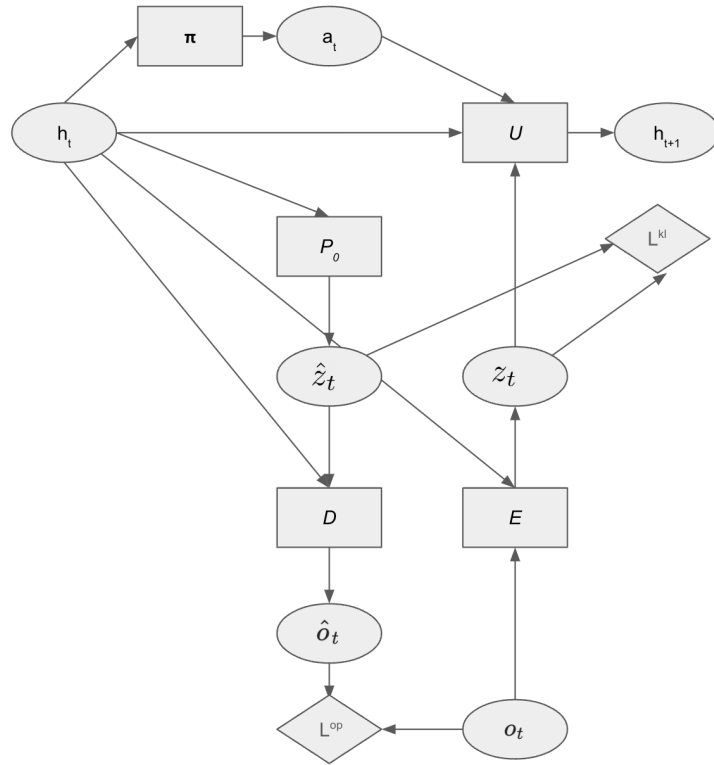


Figure 4.5: Illustration of Dreamer world model as a factor graph (squares are learnable functions, circles are variables, diamonds are fixed cost functions). We have unrolled the forwards prediction for only 1 step. Also, we have omitted the reward prediction loss.

- A latent state conditional prior: $\hat{\mathbf{z}}_t \sim P(\hat{\mathbf{z}}_t|\mathbf{h}_t)$
- A latent state decoder (observation predictor): $\hat{\mathbf{o}}_t \sim D(\hat{\mathbf{o}}_t|\mathbf{h}_t, \hat{\mathbf{z}}_t)$.
- A reward predictor: $\hat{r}_t \sim R(\hat{r}_t|\mathbf{h}_t, \hat{\mathbf{z}}_t)$
- A latent state encoder: $\mathbf{z}_t \sim E(\mathbf{z}_t|\mathbf{h}_t, \mathbf{o}_t)$.
- A policy function: $\mathbf{a}_t \sim \pi(\mathbf{a}_t|\mathbf{h}_t)$

See Figure 4.5 for an illustration of the system.

We now give a simplified explanation of how the world model is trained. The loss has the form

$$\mathcal{L}^{\text{WM}} = \mathbb{E}_{q(\mathbf{z}_{1:T})} \left[\sum_{t=1}^T \beta_o \mathcal{L}^o(\mathbf{o}_t, \hat{\mathbf{o}}_t) + \beta_z \mathcal{L}^z(\mathbf{z}_t, \hat{\mathbf{z}}_t) \right] \quad (4.21)$$

where the β terms are different weights for each loss, and q is the posterior over the latents, given by

$$q(\mathbf{z}_{1:T}|\mathbf{h}_0, \mathbf{o}_{1:T}, \mathbf{a}_{1:T}) = \prod_{t=1}^T E(\mathbf{z}_t|\mathbf{h}_t, \mathbf{o}_t) \delta(\mathbf{h}_t - \mathcal{U}(\mathbf{h}_{t-1}, \mathbf{a}_{t-1}, \mathbf{z}_{t-1})) \quad (4.22)$$

The loss terms correspond to the observation prediction cross entropy, and the posterior to prior KL penalty:

$$\mathcal{L}^o = -\ln D(\mathbf{o}_t|\mathbf{z}_t, \mathbf{h}_t) \quad (4.23)$$

$$\mathcal{L}^z = D_{\text{KL}}(E(\mathbf{z}_t|\mathbf{h}_t, \mathbf{o}_t) \parallel P(\mathbf{z}_t|\mathbf{h}_t)) \quad (4.24)$$

where we abuse notation somewhat, since \mathcal{L}^z is a function of two distributions, not of the variables \mathbf{z}_t and $\hat{\mathbf{z}}_t$.

In addition to the world model loss, we have the following actor-critic losses

$$\mathcal{L}^{\text{critic}} = \sum_{t=1}^T (V(\mathbf{h}_t) - \text{sg}(G_t^\lambda))^2 \quad (4.25)$$

$$\mathcal{L}^{\text{actor}} = - \sum_{t=1}^T \text{sg}((G_t^\lambda - V(\mathbf{h}_t))) \log \pi(\mathbf{a}_t|\mathbf{h}_t) \quad (4.26)$$

where G_t^λ is the GAE estimate of the reward to go:

$$G_t^\lambda = r_t + \gamma ((1 - \lambda)V(\mathbf{h}_t) + \lambda G_{t+1}^\lambda) \quad (4.27)$$

There have been several extensions to the original Dreamer paper. **DreamerV2** [Haf+21] adds categorical (discrete) latents and KL balancing between prior and posterior estimates. This was the first imagination-based agent to outperform humans in Atari games. **DayDreamer** [Wu+22] applies DreamerV2 to real robots. **DreamerV3** [Haf+25] builds upon DreamerV2 using various tricks — such as symlog encodings⁶ for the reward, critic, and decoder — to enable more stable optimization and domain independent choice of hyper-parameters. It was the first method to create diamonds in the Minecraft game without requiring human demonstration data. (However, reaching this goal took 17 days of training.)

Many other variants of Dreamer have been explored. For example, TransDreamer [Che+21a] and STORM [Zha+23b] replace the RNN world model with transformers, and the S4WM method of [DPA23] uses S4 (Structured State Space Sequence) models. The DreamingV2 paper of [OT22] replaces the generative loss with a non-generative self-prediction loss (see Section 4.4.2.6), and [RHH23; RHH25] use the VicReg non-generative representation learning method (see Section 4.4.2.10).

⁶The symlog function is defined as $\text{symlog}(x) = \text{sign}(x) \ln(|x| + 1)$, and its inverse is $\text{symexp}(x) = \text{sign}(x)(\exp(|x|) - 1)$. The symlog function squashes large positive and negative values, while preserving small values.

The most recent variant is **DreamerV4** [HYL25]. Here the world model is a standard conditional latent video diffusion model, $p(z_t|z_{t-c:t-1}, a_t)$, for context length c , based on a transformer backbone, combined with an autoencoder between pixel and latent space using $q(z_t|y_t)$ and $p(y_t|z_t)$. The WM is trained offline on 2500 hours of reward-free but action-labeled videos (keyboard-mouse actions) of people playing minecraft, derived from OpenAI’s VPT dataset [Bak+22]. After training the WM, a policy is initialized with behavior cloning, and then fine tuned in imagination using actor-critic methods, with minecraft-specific reward functions (learned from the VPT metadata).

The key difference from prior Dreamer models is that the world model is more powerful, and it is trained offline on diverse human-collected data. The downside is that the world model is slower to sample from.⁷ The upside it is that the model is sufficiently expressive, and the training data is sufficiently large and diverse, that the world model is accurate enough (in terms of dynamics and visual details) that it can be used to train policies purely in imagination, without mixing in any real world data. The resulting DreamerV4 recipe is similar in spirit to LLMs, that are large autoregressive transformer models that are pretrained on large diverse datasets, and then used to train policies on synthetic rollouts using RL. (In the LLM case, the policy and WM are the same model.)

4.4.1.4 Example: IRIS

The **IRIS** method (“Imagination with auto-Regression over an Inner Speech”) of [MAF22] follows the MBRL paradigm, in which it alternates between (1) learning a world model using real data D_r and then generate imaginary rollouts D_i using the WM, and (2) learning the policy given D_i and collecting new data D'_r for learning. In the model learning stage, Iris learns a discrete latent encoding using the VQ-VAE method, and then fits a transformer dynamics model to the latent codes. In the policy learning stage, it uses actor critic methods. The **Delta-IRIS** method of [MAF24] extends this by training the model to only predict the delta between neighboring frames. Note that, in both cases, the policy has the form $a_t = \pi(o_t)$, where o_t is an image, so the rollouts need to ground to pixel space, and cannot only be done in latent space.

4.4.1.5 Code world models

Recently it has become popular to represent the world model $p(s'|s, a)$ using code, such as Python. This is called a **code world model**. It is possible to learn such models from trajectory data using LLMs. See Section 6.3.3.2 for details.

4.4.1.6 Partial observation prediction

Predicting all the pixels in image may waste capacity and may distract the agent from the important bits. A natural alternative is to just predict some function of the observations, rather than the entire observation vector. This is known as a **partial world model** (see e.g., [AP23; TS11]). One way to implement this is to impose an information bottleneck between the latent state and the observed variables, to prevent the agent focusing on irrelevant observational details (see e.g., the **denoised MDP** method of [Wan+22]). We can also use code synthesis methods (see Section 6.3.3.3).

4.4.2 World models which are trained to predict other targets

In this section, we discuss training world models that are not necessarily able to predict all the future observations. These are often still (conditional) generative models (in that they return a distribution over potentially high dimensional outputs), but they are *lossy* models, because they do not capture all the details of the data.

⁷The reason is that diffusion is an iterative procedure. To speed it up, DreamerV4 uses the “shortcut” trick from [Fra+24a] to reduce the number of diffusion sampling steps to 1-2. (This is a special case of the self-consistency training methods described in [BAVE25].)

Loss	Policy	Usage	Examples
OP	Observables	Dyna	Diamond [Alo+24], Delta-Iris [MAF24]
OP	Observables	MCTS	TDM [Sch+23a]
OP	Latents	Dyna	Dreamer [Haf+25]
RP, VP, PP	Latents	MCTS	MuZero [Sch+20]
RP, VP, PP, ZP	Latents	MCTS	EfficientZero [Ye+21]
RP, VP, ZP	Latents	MPC-CEM	TD-MPC [HSW22; HSW24]
VP, ZP	Latents	Aux.	Minimalist [Ni+24]
VP, ZP	Latents	Dyna	DreamingV2 [OT22]
VP, ZP, OP	Latents	Dyna	AIS [Sub+22]
POP	Latents	Dyna	Denoised MDP [Wan+22]

Table 4.2: Summary of some world-modeling methods. The “loss” column refers to the loss used to train the latent encoder (if present) and the dynamics model (OP = observation prediction, ZP = latent state prediction, RP = reward prediction, VP = value prediction, PP = policy prediction, POP = partial observation prediction). The “policy” column refers to the input that is passed to the policy. (For MCTS methods, the policy is just used as a proposal over action sequences to initialize the search/ optimization process.) The “usage” column refers to how the world model is used: for background planning (which we call “Dyna”), or for decision-time planning (which we call “MCTS”), or just as an auxiliary loss on top of standard policy/value learning (which we call “Aux”). Thus Aux methods are single-stage (“end-to-end”), whereas the other methods alternate are two-phase, and alternate between improving the world model and then using it for improving the policy (or searching for the optimal action).

4.4.2.1 The objective mismatch problem

In Section 4.3.1, we argued that, if we can learn a sufficiently accurate world model, then solving for the optimal policy in simulation will give a policy that is close to optimal in the real world. However, a simple agent may not be able to capture the full complexity of the true environment; this is called the “**small agent, big world**” problem [DVRZ22; Lu+23; Aru+24a; Kum+24].

Consider what happens when the agent’s model is misspecified (i.e., it cannot represent the true world model), which is nearly always the case. The agent will train its model to reduce state (or observation) prediction error, by minimizing $\ell(\hat{M}, \mu_M^\pi)$. However, not all features of the state are useful for planning. For example, if the states are images, a dynamics model with limited representational capacity may choose to focus on predicting the background pixels rather than more control-relevant features, like small moving objects, since predicting the background reliably reduces the MSE more. This is due to “**objective mismatch**” [Lam+20; Wei+24], which refers to the discrepancy between the way a model is usually trained (to predict the observations) vs the way its representation is used for control. To tackle this problem, in this section we discuss methods for learning representations and models that don’t rely on predicting all the observations. Our presentation is based in part on [Ni+24] (which in turn builds on [Sub+22]). See Table 4.2 for a summary of some of the methods we will discuss.

4.4.2.2 Observation prediction

We consider a modeling paradigm where we learn an encoder, $\mathbf{z} = \phi(\mathbf{o})^8$; a dynamics model in latent space, $\mathbf{z}' = \mathcal{M}(\mathbf{z}, \mathbf{a})$, for future prediction; and an update model in latent space, $\mathbf{z}' = \mathcal{U}(\mathbf{z}, \mathbf{a}, \mathbf{o})$, for belief state tracking.

⁸Note that in general, the encoder may depend on the entire history of previous observations, denoted $\mathbf{z} = \phi(\mathcal{D})$.

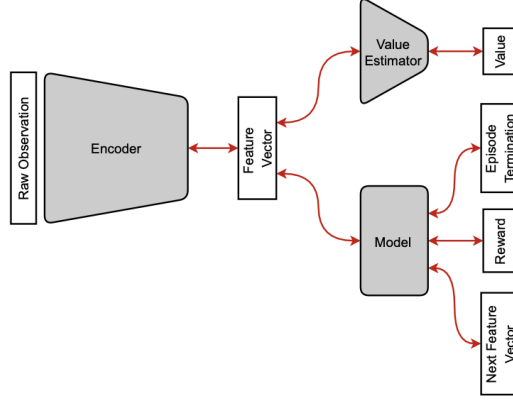


Figure 4.6: Illustration of an encoder $\mathbf{z}_t = E(\mathbf{o}_t)$, which is passed to a value estimator $v_t = V(\mathbf{z}_t)$, and a world model, which predicts the next latent state $\hat{\mathbf{z}}_{t+1} = \mathcal{M}(\mathbf{z}_t, a_t)$, the reward $r_t = R(\mathbf{z}_t, a_t)$, and the termination (done) flag, $d_t = \text{done}(\mathbf{z}_t)$. From Figure C.2 of [AP23]. Used with kind permission of Doina Precup.

A natural target to use for learning the encoder and dynamics is the next observation, using a one-step version of Equation (4.19). Indeed, [Ni+24] say that a representation ϕ satisfies the **OP** (observation prediction) criterion if it satisfies the following condition:

$$\exists D \text{ s.t. } p^*(\mathbf{o}'|\mathbf{h}, a) = D(\mathbf{o}'|\phi(\mathbf{h}), a) \quad \forall \mathbf{h}, a \quad (4.28)$$

where D is the decoder. In order to repeatedly apply this, we need to be able to update the encoding $\mathbf{z} = \phi(\mathbf{h})$ in a recursive or online way. Thus we must also satisfy the following recurrent encoder condition, which [Ni+24] call **Rec**:

$$\exists U \text{ s.t. } \phi(\mathbf{h}') = \mathcal{U}(\phi(\mathbf{h}), a, \mathbf{o}') \quad \forall \mathbf{h}, a, \mathbf{o}' \quad (4.29)$$

where \mathcal{U} is the update operator. Note that belief state updates (as in a POMDP) satisfy this property. Furthermore, belief states are a sufficient statistic to satisfy the OP condition. See Section 4.4.1.2 for a discussion of generative models of this form.

The drawback of this approach is that in general it is very hard to predict future observations, at least in high dimensional settings like images. Fortunately, such prediction is not necessary for optimal behavior. Thus we now turn our attention to other training objectives.

4.4.2.3 Reward prediction

We can also train the latent encoder to predict the reward. Formally, we want to ensure we can satisfy the following condition, which we call **RP** for “reward prediction”:

$$\exists R \text{ s.t. } \mathbb{E}_{R^*}[r|\mathbf{h}, a] = \mathbb{E}_R[r|\phi(\mathbf{h}), a] \quad \forall \mathbf{h}, a \quad (4.30)$$

See Figure 4.6 for an illustration. In [Ni+24], they prove that a representation that satisfies ZP and RP is enough to satisfy value equivalence (sufficiency for Q^*).

4.4.2.4 Value prediction

Let $\mathbf{h}_t = (\mathbf{h}_{t-1}, \mathbf{a}_{t-1}, r_{t-1}, \mathbf{o}_t)$ be all the visible data (history) at time t , and let $\mathbf{z}_t = \phi(\mathbf{h}_t)$ be a latent representation (compressed encoding) of this history, where ϕ is called an encoder or a **state abstraction**

function. We will train the policy $\mathbf{a}_t = \pi(\mathbf{z}_t)$ in the usual way, so our focus will be on how to learn good latent representations.

An optimal representation $\mathbf{z}_t = \phi(\mathbf{h}_t)$ is a sufficient statistic for the optimal action-value function Q^* . Thus it satisfies the **value equivalence** principle [LWL06; Cas11; Gri+20b; GF21; GBS22; AP23; ARKP24], which says that two states s_1 and s_2 are value equivalent (given a policy) if $V^\pi(s_1) = V^\pi(s_2)$. In particular, if the representation is optimal, it will satisfy value equivalence wrt the optimal policy, i.e., if $\phi(\mathbf{h}_i) = \phi(\mathbf{h}_j)$ then $Q^*(\mathbf{h}_i, a) = Q^*(\mathbf{h}_j, a)$. We can train such a representation function by using its output $\mathbf{z} = \phi(\mathbf{h})$ as input to the Q function or to the policy. (We call such a loss **VP**, for value prediction.) This will cause the model to focus its representational power on the relevant parts of the observation history.

Note that there is a stronger property than value equivalence called **bisimulation** [GDG03]. This says that two states s_1 and s_2 are bisimilar if $P(s'|s_1, a) \approx P(s'|s_2, a)$ and $R(s_1, a) = R(s_2, a)$. From this, we can derive a continuous measure called the **bisimulation metric** [FPP04]. This has the advantage (compared to value equivalence) of being policy independent, but the disadvantage that it can be harder to compute [Cas20; Zha+21], although there has been recent progress on computationally efficient methods such as MICo [Cas+21] and KSMe [Cas+23].

4.4.2.5 Policy prediction

The value function and reward losses may be too sparse to learn efficiently. Although self-prediction loss can help somewhat, it does not use any extra information from the environment as feedback. Consequently it is natural to consider other kinds of prediction targets for learning the latent encoder (and dynamics). When using MCTS, it is possible compute what the policy should be for a given state, and this can be used as a prediction target for the reactive policy $\mathbf{a}_t = \pi(\mathbf{z}_t)$, which in turn can be used as a feedback signal for the latent state. This method is used by MuZero (Section 4.2.2.2) and EfficientZero (Section 4.2.2.2).

4.4.2.6 Self prediction (self distillation)

Unfortunately, in problems with sparse reward, predicting the value or policy may not provide enough of a feedback signal to learn quickly. Consequently it is common to augment the training with a **self-prediction** loss where we train ϕ to ensure the following condition hold:

$$\exists M \text{ s.t. } \mathbb{E}_{M^*}[\mathbf{z}'|\mathbf{h}, a] = \mathbb{E}_M[\mathbf{z}'|\phi(\mathbf{h}), a] \quad \forall \mathbf{h}, a \quad (4.31)$$

where the LHS is the predicted mean of the next latent state under the true model, and the RHS is the predicted mean under the learned dynamics model. We call this the **EZP**, which stands for expected \mathbf{z} prediction.⁹

4.4.2.7 Avoiding self-prediction collapse using frozen targets

A trivial way to minimize the self-prediction loss is to learn an embedding that maps everything to a constant vector, say $E(\mathbf{h}) = \mathbf{0}$, in which case \mathbf{z}_{t+1} will be trivial for the dynamics model M to predict. However this is not a useful representation. This problem is **representational collapse** [Jin+22]. Fortunately, we can provably prevent collapse (at least for linear encoders) by using a frozen target network [TCG21; Tan+23; Ni+24]. That is, we use the following auxiliary loss

$$\mathcal{L}_{\text{EZP}}(\phi, \theta; \mathbf{h}, a, \mathbf{h}') = \|M_\theta(E_\phi(\mathbf{h}, a)) - \text{sg}(E_{\bar{\phi}}(\mathbf{h}'))\|_2^2 \quad (4.32)$$

where

$$\bar{\phi}_t = \rho \phi_t + (1 - \rho) \bar{\phi}_{t-1} \quad (4.33)$$

is the exponential moving average (EMA) of the encoder weights ϕ . (If we use a frozen (old) copy of the weights instead, this is called a detached network.) See Figure 4.8(a) for an illustration. This approach means

⁹In [Ni+24], they also describe the ZP loss, which requires predicting the full distribution over \mathbf{z}' using a stochastic transition model. This is strictly more powerful, but somewhat more complicated, so we omit it for simplicity.

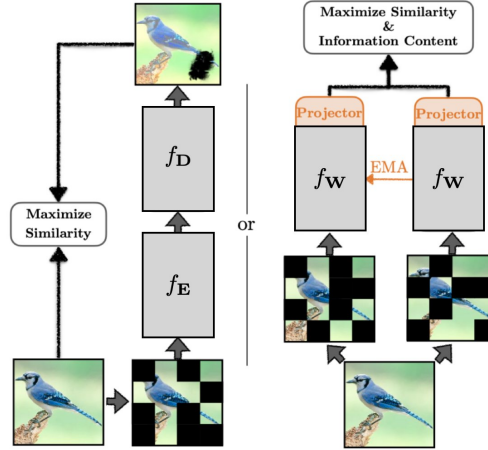


Figure 4.7: Self-supervised learning methods. Left: masked auto-encoder. Right: self-prediction. Figure from https://x.com/hugues_va/status/1991635221884891373?s=58m based on [VA+25].

the “goalposts” (the target representations) evolve slowly and consistently over time, guided by the progress of the encoder and predictor. This adds stability to the training process, ensuring the target representations don’t change erratically from one step to the next, which would make the predictor’s job impossible.

The above approach is used in many papers, such as **BYOL** [Gri+20a] (BYOL stands for “bootstrap your own latent”), **SimSiam** [CH20], **DinoV2** [Oqu+24], **JEPA** (Joint embedding Prediction Architecture) [LeC22], **I-JEPA** [Ass+23], **V-JEPA** [Bar+24; Ass+25], **Image World Models** [Gar+24], **Predictron** [Sil+17b], **Value Prediction Networks** [OSL17], **Self Predictive Representations** (SPR) [Sch+21c], **Efficient Zero** (Section 4.2.2.2), etc.

Note that an alternative to predicting the encoding of the next frame is to mask out the current frame in a random way, and then use this as the target. See Figure 4.7 for an illustration. For more details on such **self-supervised learning** or **SSL** methods, see e.g., [SZL24; Uel+25; VA+25].

Minimizing the self-prediction objective (with the stop-gradient term) has been proven to be theoretically sound for the case of linear encoders [Tan+23] with a fixed policy, which they call BYOL- π . In this case, the encoder converges learns the singular vectors of the transition matrix induced by the policy. In [Khe+25], they present BYOL-AC, which learns the singular vectors of the transition matrix for each (discrete) action, which makes the method independent of the policy. The analysis methods in these papers create an ODE where the encoder learns more slowly (using gradient steps) than the latent dynamics (which are solved to optimality at each inner step).

Proving convergence in the general nonlinear case is an open problem. It is possible that by adopting a game theoretic perspective, one can show that solving such an asymmetric Stackelberg game (see Section 5.2.13) can converge to an equilibrium, even if it does not correspond to the minimum of a single loss function. In particular, suppose the target encoder player (which acts as the slow-moving leader) makes its move by updating the EMA of the encoder network; then the prediction player (which acts as the fast-moving follower) makes it move by minimizing the EJP prediction error in Equation (4.32) using the leader’s encoder. Then it may be possible to modify the results from papers such as [Bai+21; GP23; Li+24b] to show that self-predictive learning for world models will converge in the general case.

4.4.2.8 Avoiding self-prediction collapse using information-theoretic regularization

An alternative way to avoid the latent collapse problem is to add regularization terms that try to maximize the information content in z_t and z_{t+1} (see Figure 4.8(b)), while also minimizing the prediction error. That

is, the objective becomes

$$J(\phi) = E_{\mathbf{o}_t, \mathbf{a}_t, \mathbf{o}_{t+1}, \epsilon_t} (\|\mathbf{z}_{t+1} - \hat{\mathbf{z}}_{t+1}\|_2^2 - \lambda I(\mathbf{z}_t) - \lambda I(\mathbf{z}_{t+1}))$$

where $\mathbf{z}_t = E(\mathbf{o}_t; \phi)$, $\mathbf{z}_{t+1} = E(\mathbf{o}_{t+1}; \phi)$, $\hat{\mathbf{z}}_{t+1} = \mathcal{M}(\mathbf{z}_t, \mathbf{a}_t, \epsilon_t; \theta)$

(4.34)

(Note that \mathbf{z}_{t+1} may correspond to the embedding of the next frame, \mathbf{o}_{t+1} , or to the embedding of a data augmented version of \mathbf{o}_t (e.g., a cropped or rotated version).

Various methods have been proposed to approximate the information content $I(\mathbf{z}_t)$, mostly based on some function of the outer product matrix $\sum_t \mathbf{z}_t \mathbf{z}_t^\top$, which captures second order moments.

The **Barlow Twins** [Zbo+21] method aims to make the cross-correlation matrix between the representations of two embeddings as close to the identity matrix as possible. This simultaneously encourages similar representations for similar inputs while decorrelating the features within the representation. The **VICReg** [BPL22a; BPL22b] method (Variance-Invariance-Covariance Regularization) uses three loss terms: one to maintain variance in the representations (avoid collapse), one to decorrelate the different variables in the latent vector (reduce redundancy), and one to make them invariant to data augmentations (by bringing different views closer together in embedding space). The VICReg approach can be thought of as constrastive across dimensions, whereas standard contrastive methods, such as SimCLR, are contrastive across samples (see [Gar+23]). Unfortunately, the methods in these papers do not provide a lower bound on $I(\mathbf{z}_t)$ and $I(\mathbf{z}_{t+1})$, which is needed to optimize Equation (4.34). (Since we are trying to minimize $-I$, we need an upper bound on $-I$, but these methods actually provide an upper bound on I .)

Note, however, that it is possible to use information theoretic regularizers when training generative models, that predict future observables. The idea is to create an **information bottleneck**, that remembers as little about the inputs (past) as possible, while still being able to predict the future. See e.g., [Lee+20; DR24].

4.4.2.9 Preventing self-prediction collapse using game-theoretic approaches

An arguably more theoretically sound approach to self-supervised learning, known as **CTRL** (closed-loop transcription) presented in [Dai+22], tries to solve a minmax game (similar to a GAN), rather than optimize a (regularized) objective. There are two players: the encoder E , and the generator or decoder G . The generator’s job is to minimize distortion. But rather than measuring this distance in pixel space, it is evaluated in embedding space:

$$D(\mathbf{Z}, \mathbf{Z}') = \|\mathbf{Z} - \mathbf{Z}'\|_F^2$$
(4.35)

where $\mathbf{Z} = E(\mathbf{X})$ is the embedding of a batch of inputs \mathbf{X} and

$$\mathbf{Z}' = E(G(\mathbf{Z})) = E(G(E(\mathbf{X})))$$
(4.36)

is the embedding of their reconstructions. A good encoder should also minimize the distortion.

However, to prevent the learning a trivial embedding, the encoder should also try to maximize its own **rate** (information content). In CTRL, they assume the embeddings are Gaussianly distributed, which gives the following tractable equation for the rate:

$$R(\mathbf{Z}) = \log \det(C(\mathbf{Z}))$$
(4.37)

where $C(\mathbf{Z}) = \mathbf{Z}^\top \mathbf{Z}$ is the covariance matrix of the data (assuming it is centered, for simplicity).¹⁰ Minimizing this single term achieves two goals at once: Decorrelation, since it forces the off-diagonal elements of the covariance matrix to be zero, meaning the learned features are independent and not redundant; and variance preservation, since it encourages the diagonal elements (the variance of each feature) to be non-zero and balanced, preventing the representation from collapsing into a lower-dimensional space. This is very similar in spirit to VICReg’s separate variance and covariance terms, but combined into one expression.

¹⁰In practice, computing the covariance matrix of the entire dataset, and the corresponding log determinant, is not tractable. We can handle a large number of samples by computing the rate over a minibatch of size n using $C(\mathbf{Z}_b) = \frac{1}{n} (\mathbf{Z}_b - \bar{\mathbf{Z}}_b)^\top (\mathbf{Z}_b - \bar{\mathbf{Z}}_b)$. We then use matrix determinant lemma to replace the objective with $R(\mathbf{Z}_b) = \log \det(\mathbf{Z}_b \mathbf{Z}_b^\top)$, which takes $O(n^3)$ time to compute rather than $O(d^3)$. In practice, we add a regularizer for numerical stability, and use $R(\mathbf{Z}_b) = \log \det(\mathbf{I}_n + \frac{d}{n\epsilon^2} \mathbf{Z}_b \mathbf{Z}_b^\top)$.

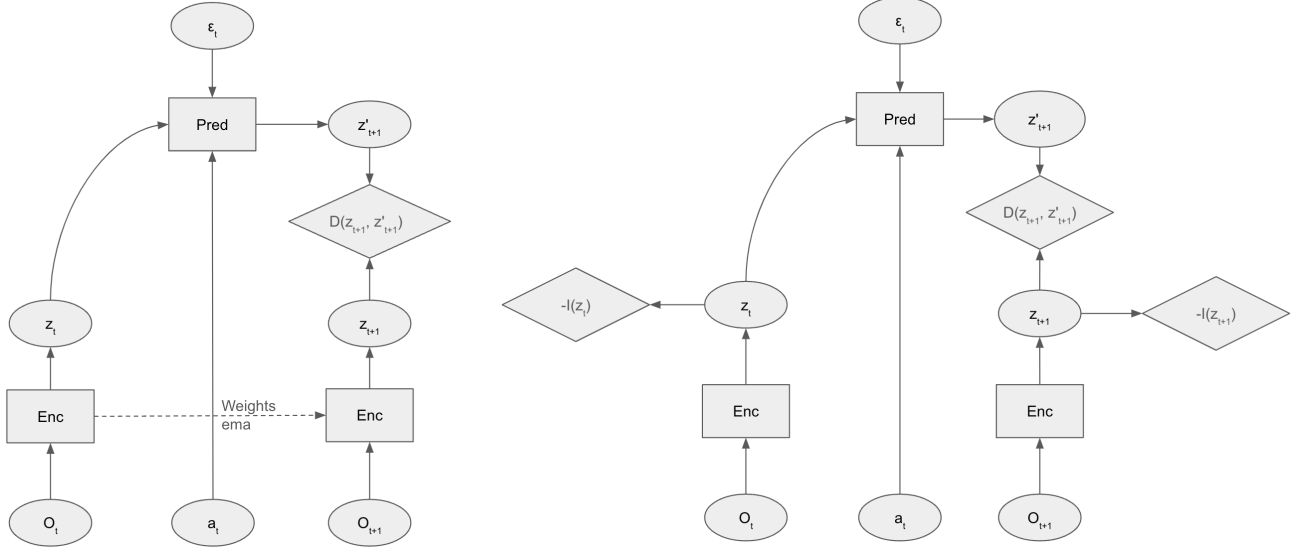


Figure 4.8: Illustration of the JEPA (Joint embedding Prediction Architecture) world model approach using two different approaches to avoid latent collapse: (a) self-distillation; (b) information theoretic regularizers. Diamonds represent fixed cost functions, squares represent learnable functions, circles are variables.

We can thus define the CTRL game as follows:

$$\min_G D(\mathbf{Z}, \mathbf{Z}'; E, G) \quad (4.38)$$

$$\max_E R(\mathbf{Z}; E) - \lambda D(\mathbf{Z}, \mathbf{Z}'; E, G) \quad (4.39)$$

Since the rate only depends on the encoder, there is an imbalance in the game. Hence, rather than just having the encoder maximizing the rate of the original data, we can train the encoder to maximize the **rate reduction**, and the decoder to minimize it, where the rate reduction is given by

$$\Delta R(\mathbf{Z}, \mathbf{Z}'; E, G) = R(\mathbf{Z}; E) - R(\mathbf{Z}'; E, G) \geq 0 \quad (4.40)$$

The new game becomes

$$\min_G \Delta R(\mathbf{Z}, \mathbf{Z}'; E, G) + \lambda D(\mathbf{Z}, \mathbf{Z}'; E, G) \quad (4.41)$$

$$\max_E \Delta R(\mathbf{Z}, \mathbf{Z}'; E, G) - \lambda D(\mathbf{Z}, \mathbf{Z}'; E, G) \quad (4.42)$$

To see why this works, note that the encoder E has to maximize $R(\mathbf{Z})$, so it must find features in the real data X that are diverse and fill the latent space; but it also must minimize $R(\mathbf{Z}')$, so it must be a good critic that can find the flaws and repetitive patterns in the Generator's fake \mathbf{X}' , encoding them into a collapsed, low-rate, redundant representation \mathbf{Z}' . Conversely the best way for the generator to minimize the rate reduction is to ensure that it makes \mathbf{X}' identical to \mathbf{X} , so then \mathbf{Z}' will be identical to \mathbf{Z} , and the rate reduction will be zero.

Although the above approach is elegant, it is designed for a set of iid static images. The best way to extend it to sequential data, as arises when learning a world model, is an open research problem.

4.4.2.10 Example: JEPA

In this section, we discuss the **JEPA** (Joint embedding Prediction Architecture) approach to world modeling, first proposed in [LeC22]. The basic idea is to jointly embed the current and following observations, to

compute $\mathbf{z}_t = E(\mathbf{o}_t)$ and $\mathbf{z}_{t+1} = E(\mathbf{o}_{t+1})$, and then to compare the actual latent embedding \mathbf{z}_{t+1} to its prediction $\mathbf{z}'_{t+1} = \mathcal{M}(\mathbf{z}_t, \mathbf{a}_t; \epsilon_t)$, where ϵ_t is a random noise source, and M is the deterministic world model. We then train the encoder to minimize the difference between \mathbf{z}_t and \mathbf{z}'_t .

To prevent the encoder collapsing to a trivial function, such as $E(\mathbf{o}) = \mathbf{0}$, two different classes of methods have been considered. The first is based on using a frozen EMA version of the encoder, as discussed in Section 4.4.2.7. JEPA also leverages the fact that the encoder is a low-dimensional embedding of the input, and the predictor is a shallow network, to create an information bottleneck. (For example, if we use a ViT encoder, an image of size $224 \times 224 \times 3 = 150,528$ gets compressed into a much smaller number of tokens, which often amounts to an embedding size of ~ 1024 .) The **I-JEPA** method of [Ass+23], designed for images, and the **V-JEPA** method of [Bar+24; Ass+25], designed for videos, also trains on masked versions of the inputs to create a harder learning problem, since the target shares less information with the input.

An alternative way to avoid the latent collapse problem is to add regularization terms that try to maximize the information content in \mathbf{z}_t and \mathbf{z}_{t+1} , while also minimizing the prediction error, as discussed in Section 4.4.2.8. See Figure 4.8 for an illustration to these two approaches. (They can also be combined. For example, [MT24] combines I-JEPA with VICReg.)

4.4.2.11 Example: DinoWM

In the case where the observations are high-dimensional, such as images, it is natural to use a pre-trained representation, $\mathbf{z}_t = \phi(\mathbf{o}_t)$, as input to the world model (or policy). The representation function ϕ can be pretrained on a large dataset using a non-reconstructive loss, such as the DINOv2 method [Oqu+24]. Although this can sometimes give gains (as in the **DinoWM** and **Dino-World** papers [Zho+24a; Bal+25b]), in other cases, better results are obtained by training the representation from scratch [Han+23; Sch+24]. The performance is highly dependent on the similarity or differences between the pretraining distribution and the agent’s distribution, the form of the representation function, and its training objective.

4.4.2.12 Example: TD-MPC

In this section, we describe **TD-MPC2** [HSW24], which is an extension of **TD-MPC** of [HSW22]. This learns the following functions:

- Encoder: $\mathbf{e}_t = E(\mathbf{o}_t)$
- Latent dynamics (for rollouts): $\mathbf{z}'_t = \mathcal{M}(\mathbf{z}_{t-1}, \mathbf{a}_t)$
- Latent update (after each observation): $\mathbf{z}_t = \mathcal{U}(\mathbf{z}_{t-1}, \mathbf{e}_t, \mathbf{a}_t) = \mathbf{e}_t$
- Reward: $\hat{r}_t = R(\mathbf{z}_t, \mathbf{a}_t)$
- Value: $\hat{q}_t = Q(\mathbf{z}_t, \mathbf{a}_t)$
- Policy prior: $\hat{\mathbf{a}}_t = \pi_{\text{prior}}(\mathbf{z}_{t-1})$

The model is trained using the following VP+ZP loss applied to trajectories sampled from the replay buffer:

$$\mathcal{L}(\theta) = \mathbb{E}_{(\mathbf{o}, \mathbf{a}, r, \mathbf{o}')_{0:H} \sim \mathcal{B}} \left[\sum_{t=0}^H \lambda^t (||\mathbf{z}'_t - \text{sg}(E(\mathbf{o}'_t))||_2^2 + \text{CE}(\hat{r}_t, r_t) + \text{CE}(\hat{q}_t, q_t)) \right] \quad (4.43)$$

We use cross-entropy loss on a discretized representation of the reward and Q value in a log-transformed space, in order to be robust to different value scales across time and problem settings (see Section 7.3.2). The target value for the Q function update is defined by

$$q_t = r_t + \gamma \overline{Q}(\mathbf{z}'_t, \pi_{\text{prior}}(\mathbf{z}'_t)) \quad (4.44)$$

where \overline{Q} is the EMA for the Q function.

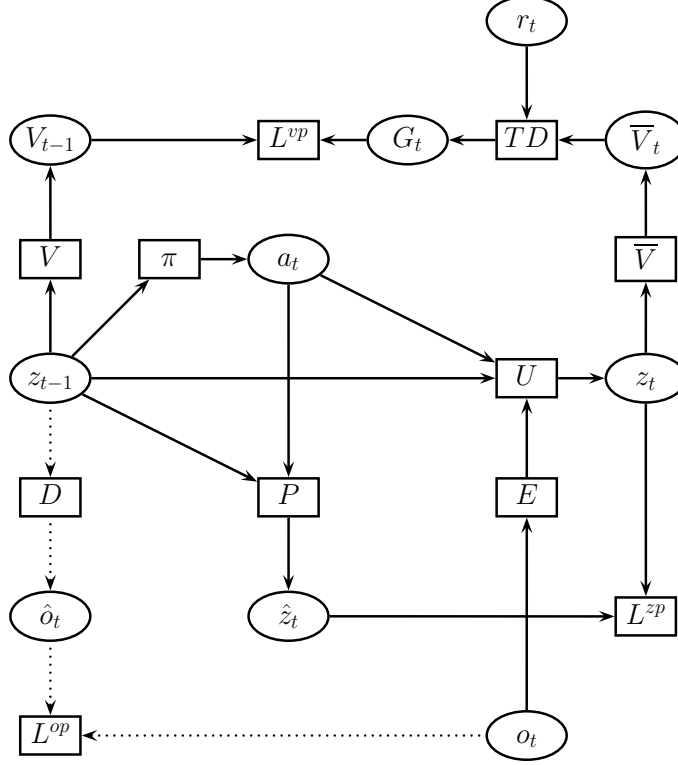


Figure 4.9: Illustration of (a simplified version of) the BYOL-Explore architecture, represented as a factor graph (so squares are functions, circles are variables). The dotted lines represent an optional observation prediction loss. The map from notation in this figure to the paper is as follows: $U \rightarrow h^c$ (closed-loop RNN update), $P \rightarrow h^o$ (open-loop RNN update), $D \rightarrow g$ (decoder), $E \rightarrow f$ (encoder). We have unrolled the forwards prediction for only 1 step. Also, we have omitted the reward prediction loss. The \bar{V} node is the EMA version of the value function. The TD node is the TD operator.

The policy is trained using the SAC objective (see Section 3.6.8) on imaginary rollouts in latent space using observations and actions from the replay buffer:

$$\mathcal{L}_\pi(\theta) = \mathbb{E}_{(\mathbf{o}, \mathbf{a})_{0:H} \sim \mathcal{B}} \left[\sum_{t=0}^H \lambda^t [\alpha Q(\mathbf{z}_t, \pi_{\text{prior}}(\mathbf{z}_t)) - \beta \mathbb{H}(\pi_{\text{prior}}(\cdot | \mathbf{z}_t))] \right], \mathbf{z}_{t+1} = \mathcal{M}(\mathbf{z}_t, \mathbf{a}_t), \mathbf{z}_0 = E(\mathbf{o}_0) \quad (4.45)$$

This policy is used as a proposal (prior), in conjunction with the MPPI trajectory planning method (Section 4.2.4.6) to select actions at run time.

In [Esc+25], they extended TD-MPC to work with sparse rewards by assuming access to demonstrations, partitioning the overall task into stages, and then converting this semi-sparse stage reward to a dense reward, which is needed to learn the world model. In [Han+25], they avoid the need for demonstration data, by using a latent prior on human motion, derived from mocap data, to learn a whole body world model and controller.

4.4.2.13 Example: BYOL

In [Gri+20a], they present **BYOL** (Build Your Own Latents), which uses the ZP and VP loss. See Figure 4.9 for the computation graph, which we see is slightly simpler than the Dreamer computation graph in Figure 4.5 due to the lack of stochastic latents.

In [Guo+22], they present **BYOL-Explore**, which extends BYOL by using the self-prediction error to define an intrinsic reward. This encourages the agent to explore states where the model is uncertain. See Section 7.4 for further discussion of this topic.

4.4.2.14 Example: Imagination-augmented agents

In [Web+17], they train a model to predict future states and rewards, and then use the hidden states of this model as additional context for a policy-based learning method. This can help overcome partial observability. They call their method **imagination-augmented agents**.

4.4.3 World models that are trained to help planning

One solution to the objective mismatch problem is to use **differentiable planning**, in which we combine model learning and policy learning together, and train them jointly end-to-end, rather than in an alternating fashion. In particular, we can solve try to optimize

$$\min_{\hat{M}, Q} \mathbb{E}_{(s, a, s') \sim \mathcal{D}} [(R(s, a) + \gamma V(s') - Q(s, a))^2]$$

where $s' = \hat{M}(s, a)$ is the learned dynamics model, subject to the constraint that the value function is derived from the model using

$$V(s) = \operatorname{argmax}_{a(0:K)} \mathbb{E}_{\hat{M}} \left[\sum_{k=0}^{K-1} \gamma^k R(s_k, a_k) + \gamma^K V(s_K) \mid S_0 = s \right].$$

This bilevel optimization problem was first proposed in the **Value Iteration Network** paper of [Tam+16], and extended in the **TreeQN** paper [Far+18]. In [ML25], they propose **D-TSN** (differentiable tree search network), which is similar to TreeQN, but constructs a best-first search tree, rather than a fixed depth tree, using a stochastic tree expansion method.

In [Nik+22; Ban+23] they propose to use implicit differentiation to avoid explicitly unrolling the inner optimization.

4.4.4 Dealing with model errors and uncertainty

The theory in Section 4.3.1 tells us that the model-as-leader approach, which trains a new policy in imagination at each inner iteration while gradually improving the model in the outer loop, will converge to the optimal policy, provided the model converges to the true model (or one that is value equivalent to it, see Section 4.4.2.4). This can be assured provided the model is sufficiently powerful, and the policy explores sufficiently widely to collect enough diverse but task-relevant data. Nevertheless, models will inevitably have errors, and it can be useful for the policy learning to be aware of this (see e.g., [Aru+18]). We discuss some approaches to this below.

4.4.4.1 Avoiding compounding errors in rollouts

In MBRL, we have to rollout imaginary trajectories to use for training the policy. It makes intuitive sense to start from a previously visited real-world state, since the model will likely be reliable there. We should start rollouts from different points along each real trajectory, to ensure good state coverage, rather than just expanding around the initial state [Raj+17]. However, if we roll out too far from a previously seen state, the trajectories are likely to become less realistic, due to **compounding errors** from the model [LPC22].

In [Jan+19a], they present the MBPO method, which uses short rollouts (inside Dyna) to prevent compounding error (an approach which is justified in [Jia+15]). [Fra+24b] is a recent extension of MBPO which dynamically decides how much to roll out, based on model uncertainty.

Another approach to mitigating compounding errors is to learn a trajectory-level dynamics model, instead of a single-step model, see e.g., [Zho+24b] which uses diffusion to train $p(s_{t+1:t+H} | s_t, a_{t:t+H-1})$, and uses this inside an MPC loop.

If the model is able to predict a reliable distribution over future states, then we can leverage this uncertainty estimate to compute an estimate of the expected reward. For example, PILCO [DR11; DFR15] uses Gaussian processes as the world model, and uses this to analytically derive the expected reward over

trajectories as a function of policy parameters, which are then optimized using a deterministic second-order gradient-based solver. In [Man+19], they combine the MPO algorithm (Section 3.6.5) for continuous control with **uncertainty sets** on the dynamics to learn a policy that optimizes for a worst case expected return objective.

4.4.4.2 Unified model and planning variational lower bound

In [Eys+22], they propose a method called **Mismatched No More** (MNM) to solve the objective mismatch problem. They define an optimality variable (see Section 3.6) based on the entire trajectory, $p(O = 1|\tau) = R(\tau) = \sum_{t=1}^{\infty} \gamma^t R(s_t, a_t)$. This gives rise to the following variational lower bound on the log probability of optimality:

$$\log p(O = 1) = \log \int_{\tau} P(O = 1, \tau) = \log \mathbb{E}_{P(\tau)} [P(O = 1|\tau)] \geq \mathbb{E}_{Q(\tau)} [\log R(\tau) + \log P(\tau) - \log Q(\tau)]$$

where $P(\tau)$ is the distribution over trajectories induced by policy applied to the true world model, $P(\tau) = \mu(s_0) \prod_{t=0}^{\infty} M(s_{t+1}|s_t, a_t) \pi(a_t|s_t)$, and $Q(\tau)$ is the distribution over trajectories using the estimated world model, $Q(\tau) = \mu(s_0) \prod_{t=0}^{\infty} \hat{M}(s_{t+1}|s_t, a_t) \pi(a_t|s_t)$. They then maximize this bound wrt π and \hat{M} .

In [Ghu+22] they extend MNM to work with images (and other high dimensional states) by learning a latent encoder $\hat{E}(z_t|o_t)$ as well as latent dynamics $\hat{M}(z_{t+1}|z_t, a_t)$, similar to other self-predictive methods (Section 4.4.2.6). They call their method **Aligned Latent Models**.

4.4.4.3 Dynamically switching between MFRL and MBRL

One problem with the above methods is that, if the model is of limited capacity, or if it learns to model “irrelevant” aspects of the environment, then any MBRL method may be dominated by a MFRL method that directly optimizes the true expected reward. A safer approach is to use a model-based policy only when the agent is confident it is better, but otherwise to fall back to a model-free policy. This is the strategy proposed in the **Unified RL** method of [Fre+24].

4.4.5 Exploration for learning world models

In Section 1.3.5, we discussed the exploration-exploitation tradeoff, which contrasts the need to (1) collect diverse experiences (by trying many new actions in many new states) so as to learn a better policy to help long-run performance with (2) the need to stay in familiar parts of the state space where the optimal policy has already been learned, so as to ensure short-term rewards. When using MBRL, the need for diverse data becomes even more important, to ensure we learn the correct underlying world model (which is then used to train the policy, or for online planning).

One popular approach to this is to use posterior sampling RL, which applies Thompson sampling to the MDP parameters (i.e., the world model), as explained in Section 7.2.2.2. This was applied to MBRL in [WCM24].

If we are in the reward-free setting (see Section 7.4), we can view the problem of learning a world model as similar to the scientist’s job of trying to create a **causal model** of the world, which can explain the effects of actions (interventions). This requires designing and carrying out experiments in order to collect informative trajectories for model fitting (see e.g., [Sek+20; Cao+25a; Cao+25c; Cao+25b]). Recently it has become popular to use LLMs to help with this problem; this can be thought of an **AI scientist** (see e.g., [Gan+25b]). It is also possible to combine LLMs as hypothesis generators with Bayesian inference and information theoretic reasoning principles for a more principled approach to the problem (see e.g., [Pir+24]).

4.5 Beyond one-step models: predictive representations

The “world models” we described in Section 4.4 are **one-step models** of the form $p(s'|s, a)$, or $p(z'|z, a)$ for $z = \phi(s)$, where ϕ is a state-abstraction function. However, such models are problematic when it comes

to predicting many kinds of future events, such as “will a car pull in front of me?” or “when will it start raining?”, since it is hard to predict exactly when these events will occur, and these events may correspond to many different “ground states”. In principle we can roll out many possible long term futures, and apply some abstraction function to the resulting generated trajectories to extract features of interest, and thus derive a predictive model of the form $p(t', \phi(s_{t+1:t'})|s_t, \pi)$, where t' is the random duration of the sampled trajectory, and ϕ maps from state trajectories to features. However, it would be more efficient if we could directly predict this distribution without having to know the value of t' , and without having to predict all the details of all the intermediate future states, many of which will be irrelevant after we pass them into the abstraction function ϕ . This motivates the study of multi-step world models, that predict multiple steps into the future, either at the state level, or at the feature level. These are called **predictive representations**, and are a compromise between standard model-based RL and model-free RL, as we will see. Our presentation on this topic is based on [Car+24]. (See also Section 7.5, where we discuss the related topic of temporal abstraction from a model-free perspective.)

4.5.1 General value functions

The value function is based on predicting the sum of expected discounted future rewards. But the reward is just one possible signal of interest we can extract from the environment. We can generalize this by considering a **cumulant** $C_t \in \mathbb{R}$, which is some scalar of interest derived from the state or observation (e.g., did a loud bang just occur? is there a tree visible in the image?). We then define the **general value function** or **GVF** as follows [Sut95; Sut+11; Com+18; Rin21; Xu+22; GGL22]:

$$V^{\pi, C, \gamma}(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t C(s_{t+1}) | s_0 = s, a_{0:\infty} \sim \pi \right] \quad (4.46)$$

If $C(s_{t+1}) = R_{t+1}$, this reduces to the value function.¹¹ If we define the cumulant to be the observation vector, then the GVF will learn to predict future observations at multiple time scales; this is called **nexting** [MS14; MWS14; Whi15]. Predicting the GVFs for multiple cumulants can be useful as an auxiliary task while solving the main task (e.g., as a form of auxiliary input to the policy, or just to “densify” the training signal), as shown in [Jad+17].

[Vee+19] present an approach (based on meta-gradients) to learn which cumulants are worth predicting. In the inner loop, the model f predicts the policy π_t and value function V_t , as usual, and also predicts the GVFs \mathbf{y}_t for the specified cumulants; the function f is called the answer network, and is denoted by $(\pi_t, V_t, \mathbf{y}_t) = f_{\theta}(\mathbf{o}_{t-i-1:t})$. In the outerloop, the model g learns to extract the cumulants and their discounts given future observations; this called the question network and is denoted by $(\mathbf{c}_t, \gamma_t) = g_{\eta}(\mathbf{o}_{t+1:t+j})$. The outer update to η is based on the gradient of the RL loss after performing K inner updates to θ using the RL loss and auxiliary loss.

4.5.2 Successor representations

In this section we consider a variant of GVF where the cumulant corresponds to a state occupancy vector $C_{\tilde{s}}(s_{t+1}) = \mathbb{I}(s_{t+1} = \tilde{s})$, which provides a dense feedback signal. Computing this for each possible state \tilde{s} gives us the **successor representation** or **SR** [Day93; Mac+23]:

$$M^{\pi}(s, \tilde{s}) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \mathbb{I}(s_{t+1} = \tilde{s}) | S_0 = s \right] \quad (4.47)$$

If we define the policy-dependent state-transition matrix by

$$T^{\pi}(s, s') = \sum_a \pi(a|s) T(s'|s, a) \quad (4.48)$$

¹¹This follows the convention of [SB18], where we write $(s_t, a_t, r_{t+1}, s_{t+1})$ — as opposed to (s_t, a_t, r_t, s_{t+1}) — to represent the transitions, since r_{t+1} and s_{t+1} are both generated by applying a_t in state s_t .

then the SR matrix can be rewritten as

$$\mathbf{M}^\pi = \sum_{t=0}^{\infty} \gamma^t [\mathbf{T}^\pi]^{t+1} = \mathbf{T}^\pi (\mathbf{I} - \gamma \mathbf{T}^\pi)^{-1} \quad (4.49)$$

Thus we see that the SR replaces information about individual transitions with their cumulants, just as the value function replaces individual rewards with the reward-to-go.

Like the value function, the SR obeys a Bellman equation

$$M^\pi(s, \tilde{s}) = \sum_a \pi(a|s) \sum_{s'} T(s'|s, a) (\mathbb{I}(s' = \tilde{s}) + \gamma M^\pi(s', \tilde{s})) \quad (4.50)$$

$$= \mathbb{E} [\mathbb{I}(s' = \tilde{s}) + \gamma M^\pi(s', \tilde{s})] \quad (4.51)$$

Hence we can learn an SR using a TD update of the form

$$M^\pi(s, \tilde{s}) \leftarrow M^\pi(s, \tilde{s}) + \underbrace{\eta (\mathbb{I}(s' = \tilde{s}) + \gamma M^\pi(s', \tilde{s}) - M^\pi(s, \tilde{s}))}_{\delta} \quad (4.52)$$

where s' is the next state sampled from $T(s'|s, a)$. Compare this to the value-function TD update in Equation (2.16):

$$V^\pi(s) \leftarrow V^\pi(s) + \underbrace{\eta (R(s') + \gamma V^\pi(s') - V^\pi(s))}_{\delta} \quad (4.53)$$

However, with an SR, we can easily compute the value function for any reward function (given a fixed policy) as follows:

$$V^{R, \pi} = \sum_{\tilde{s}} M^\pi(s, \tilde{s}) R(\tilde{s}) \quad (4.54)$$

See Figure 4.10 for an example.

We can also make a version of SR that depends on the action as well as the state to get

$$M^\pi(s, a, \tilde{s}) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \mathbb{I}(s_{t+1} = \tilde{s}) \mid s_0 = s, a_0 = a, a_{1:\infty} \sim \pi \right] \quad (4.55)$$

$$= \mathbb{E} [\mathbb{I}(s' = \tilde{s}) + \gamma M^\pi(s', a, \tilde{s}) \mid s_0 = s, a_0 = a, a_{1:\infty} \sim \pi] \quad (4.56)$$

This gives rise to a TD update of the form

$$M^\pi(s, a, \tilde{s}) \leftarrow M^\pi(s, a, \tilde{s}) + \underbrace{\eta (\mathbb{I}(s' = \tilde{s}) + \gamma M^\pi(s', a', \tilde{s}) - M^\pi(s, a, \tilde{s}))}_{\delta} \quad (4.57)$$

where s' is the next state sampled from $T(s'|s, a)$ and a' is the next action sampled from $\pi(s')$. Compare this to the (on-policy) SARSA update from Equation (2.28):

$$Q^\pi(s, a) \leftarrow Q^\pi(s, a) + \underbrace{\eta (R(s') + \gamma Q^\pi(s', a') - Q^\pi(s, a))}_{\delta} \quad (4.58)$$

However, from an SR, we can compute the state-action value function for any reward function:

$$Q^{R, \pi}(s, a) = \sum_{\tilde{s}} M^\pi(s, a, \tilde{s}) R(\tilde{s}) \quad (4.59)$$

This can be used to improve the policy as we discuss in Section 4.5.3.1.

We see that the SR representation has the computational advantages of model-free RL (no need to do explicit planning or rollouts in order to compute the optimal action), but also the flexibility of model-based RL (we can easily change the reward function without having to learn a new value function). This latter property makes SR particularly well suited to problems that use intrinsic reward (see Section 7.4), which often changes depending on the information state of the agent.

Unfortunately, the SR is limited in two key ways: (1) it assumes a finite, discrete state space; (2) it depends on a given policy. We discuss ways to overcome limitation 1 in Section 4.5.3 and Section 4.5.4, and limitation 2 in Section 4.5.5.

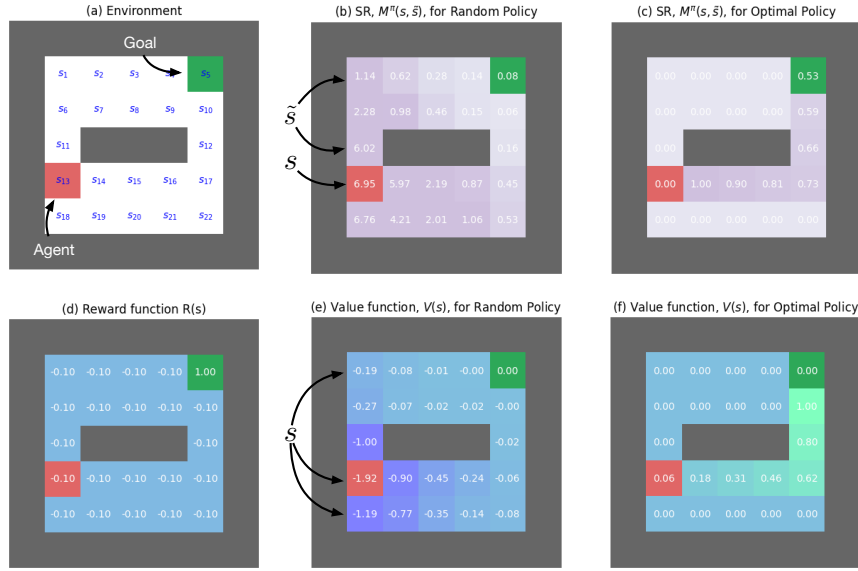


Figure 4.10: Illustration of successor representation for the 2d maze environment shown in (a) with reward shown in (d), which assigns all states a reward of -0.1 except for the goal state which has a reward of 1.0. In (b-c) we show the SRs for a random policy and the optimal policy. In (e-f) we show the corresponding value functions. In (b), we see that the SR under the random policy assigns high state occupancy values to states which are close (in Manhattan distance) to the current state s_{13} (e.g., $M^\pi(s_{13}, s_{14}) = 5.97$) and low values to states that are further away (e.g., $M^\pi(s_{13}, s_{12}) = 0.16$). In (c), we see that the SR under the optimal policy assigns high state occupancy values to states which are close to the optimal path to the goal (e.g., $M^\pi(s_{13}, s_{14}) = 1.0$) and which fade with distance from the current state along that path (e.g., $M^\pi(s_{13}, s_{12}) = 0.66$). From Figure 3 of [Car+24]. Used with kind permission of Wilka Carvalho. Generated by https://github.com/wcarvalho/jazneurorl/blob/main/successor_representation.ipynb.

4.5.3 Successor features

SRs require defining expectations or distributions over the entire future state vector, which can be problematic in high dimensional and continuous spaces. In [Bar+17] they introduced **successor features**, that generalize SRs by working with features $\phi(s)$ instead of primitive states. In particular, if we define the cumulant to be $C(s_{t+1}) = \phi(s_{t+1})$, we get the following definition of SF:

$$\psi^{\pi, \phi}(s) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_{t+1}) | s_0 = s, a_{0:\infty} \sim \pi \right] \quad (4.60)$$

We will henceforth drop the ϕ superscript from the notation, for brevity.

By analogy to Equation (4.49), we can write the successor features in matrix form as follows:

$$\Psi^{\pi} = \sum_{t=0}^{\infty} (\gamma \mathbf{T}^{\pi})^t \Phi = (\mathbf{I} - \gamma \mathbf{T}^{\pi})^{-1} \Phi \quad (4.61)$$

where Φ is the $S \times D$ matrix of features for each state.

SFs also obey a Bellman equation

$$\psi(s) = \mathbb{E} [\phi(s') + \gamma \psi(s')] \quad (4.62)$$

If we assume the reward function can be written as

$$R(s, \mathbf{w}) = \phi(s)^{\top} \mathbf{w} \quad (4.63)$$

then we can derive the value function for any reward as follows:

$$V^{\pi, \mathbf{w}}(s) = \mathbb{E} [R(s_1) + \gamma R(s_2) + \dots | s_0 = s] \quad (4.64)$$

$$= \mathbb{E} [\phi(s_1)^{\top} \mathbf{w} + \gamma \phi(s_2)^{\top} \mathbf{w} + \dots | s_0 = s] \quad (4.65)$$

$$= \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_{t+1}) | s_0 = s \right]^{\top} \mathbf{w} = \psi^{\pi}(s)^{\top} \mathbf{w} \quad (4.66)$$

Similarly we can define an action-conditioned version of SF as

$$\psi^{\pi, \phi}(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \phi(s_{t+1}) | s_0 = s, a_0 = a, a_{1:\infty} \sim \pi \right] \quad (4.67)$$

$$= \mathbb{E} [\phi(s') + \gamma \psi(s', a')] \quad (4.68)$$

We can learn this using a TD rule

$$\psi^{\pi}(s, a) \leftarrow \psi^{\pi}(s, a) + \underbrace{\eta (\phi(s') + \gamma \psi^{\pi}(s', a') - \psi^{\pi}(s, a))}_{\delta} \quad (4.69)$$

And we can use it to derive a state-action value function:

$$Q^{\pi, \mathbf{w}}(s) = \psi^{\pi}(s, a)^{\top} \mathbf{w} \quad (4.70)$$

This allows us to define multiple Q functions (and hence policies) just by changing the weight vector \mathbf{w} , as we discuss in Section 4.5.3.1.

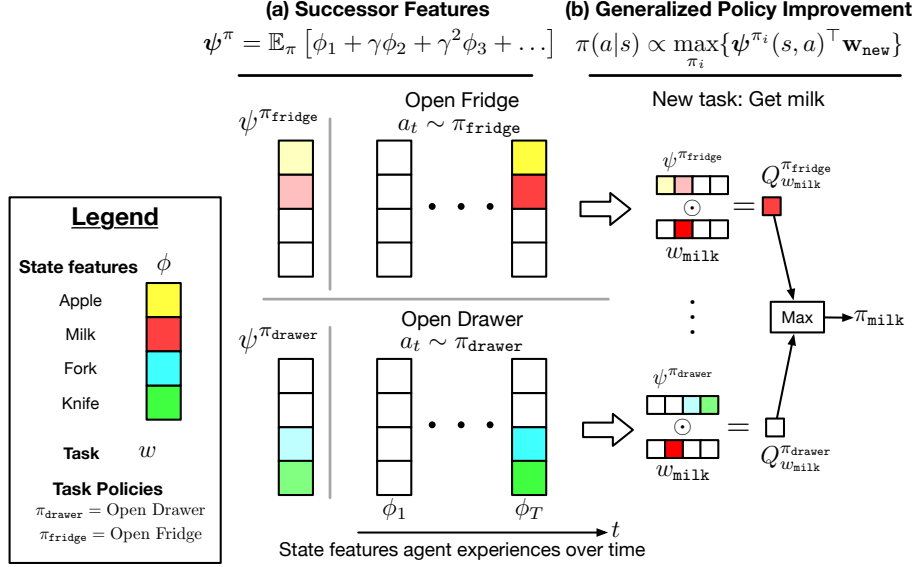


Figure 4.11: Illustration of successor features representation. (a) Here $\phi_t = \phi(s_t)$ is the vector of features for the state at time t , and ψ^π is the corresponding SF representation, which depends on the policy π . (b) Given a set of existing policies and their SFs, we can create a new one by specifying a desired weight vector \mathbf{w}_{new} and taking a weighted combination of the existing SFs. From Figure 5 of [Car+24]. Used with kind permission of Wilka Carvalho.

4.5.3.1 Generalized policy improvement

So far, we have discussed how to compute the value function for a new reward function but using the SFs from an existing known policy. In this section we discuss how to create a new policy that is better than an existing set of policies, by using **Generalized Policy Improvement** or **GPI** [Bar+17; Bar+20].

Suppose we have learned a set of N (potentially optimal) policies π_i and their corresponding SFs ψ^{π_i} for maximizing rewards defined by \mathbf{w}_i . When presented with a new task \mathbf{w}_{new} , we can compute a new policy using GPI as follows:

$$a^*(s; \mathbf{w}_{\text{new}}) = \underset{a}{\operatorname{argmax}} \max_i Q^{\pi_i}(s, a, \mathbf{w}_{\text{new}}) = \underset{a}{\operatorname{argmax}} \max_i \psi^{\pi_i}(s, a)^\top \mathbf{w}_{\text{new}} \quad (4.71)$$

If \mathbf{w}_{new} is in the span of the training tasks (i.e., there exist weights α_i such that $\mathbf{w}_{\text{new}} = \sum_i \alpha_i \mathbf{w}_i$), then the GPI theorem states that $\pi(a|s) = \mathbb{I}(a = a^*(s, \mathbf{w}_{\text{new}}))$ will perform at least as well as any of the existing policies, i.e., $Q^\pi(s, a) \geq \max_i Q^{\pi_i}(s, a)$ (c.f., policy improvement in Section 3.3). See Figure 4.11 for an illustration.

Note that GPI is a model-free approach to computing a new policy, based on an existing library of policies. In [Ale+23], they propose an extension that can also leverage a (possibly approximate) world model to learn better policies that can outperform the library of existing policies by performing more decision-time search.

4.5.3.2 Option keyboard

One limitation of GPI is that it requires that the reward function, and the resulting policy, be defined in terms of a fixed weight vector \mathbf{w}_{new} , where the preference over features is constant over time. However, for some tasks we might want to initially avoid a feature or state and then later move towards it. To solve this, [Bar+19; Bar+20] introduced the **option keyboard**, in which the weight vector for a task can be computed dynamically in a state-dependent way, using $\mathbf{w}_s = g(s, \mathbf{w}_{\text{new}})$. (Options are discussed in Section 7.5.1.) Actions can then be chosen as follows:

$$a^*(s; \mathbf{w}_{\text{new}}) = \underset{a}{\operatorname{argmax}} \max_i \psi^{\pi_i}(s, a)^\top \mathbf{w}_s \quad (4.72)$$

Thus \mathbf{w}_s induces a set of policies that are active for a period of time, similar to playing a chord on a piano.

4.5.3.3 Learning SFs

A key question when using SFs is how to learn the cumulants or state-features $\phi(s)$. Various approaches have been suggested, including leveraging meta-gradients [Vee+19]; image reconstruction [Mac+18b]; maximizing the mutual information between task encodings and the cumulants that an agent experiences when pursuing that task [Han+19]; and reward prediction methods [Chu+24]. The cumulants are encouraged to satisfy the linear reward constraint by minimizing

$$\mathcal{L}_r = \|r - \phi_\theta(s)^\top \mathbf{w}\|_2^2 \quad (4.73)$$

Once the cumulant function is known, we have to learn the corresponding SF. The standard approach learns a different SF for every policy, which is limiting. In [Bor+19] they introduced **Universal Successor Feature Approximators** which takes as input a policy encoding \mathbf{z}_w , representing a policy π_w (typically we set $\mathbf{z}_w = \mathbf{w}$). We then define

$$\psi^{\pi_w}(s, a) = \psi_\theta(s, a, \mathbf{z}_w) \quad (4.74)$$

The GPI update then becomes

$$a^*(s; \mathbf{w}_{\text{new}}) = \underset{a}{\operatorname{argmax}} \max_{\mathbf{z}_w} \psi_\theta(s, a, \mathbf{z}_w)^\top \mathbf{w}_{\text{new}} \quad (4.75)$$

so we replace the discrete over a finite number of policies, \max_i , with a continuous optimization problem $\max_{\mathbf{z}_w}$, to be solved per state.

If we want to learn the policies and SFs at the same time, we can optimize the following losses in parallel:

$$\mathcal{L}_Q = \|\psi_\theta(s, a, \mathbf{z}_w)^\top \mathbf{w} - \mathbf{y}_Q\|, \quad \mathbf{y}_Q = R(s'; \mathbf{w}) + \gamma \psi_\theta(s', a^*, \mathbf{z}_w)^\top \mathbf{w} \quad (4.76)$$

$$\mathcal{L}_\psi = \|\psi_\theta(s, a, \mathbf{z}_w) - \mathbf{y}_\psi\|, \quad \mathbf{y}_\psi = \phi(s') + \gamma \psi_\theta(s', a^*, \mathbf{z}_w) \quad (4.77)$$

where $a^* = \underset{a'}{\operatorname{argmax}} \psi_\theta(s', a', \mathbf{z}_w)^\top \mathbf{w}$. The first equation is standard Q learning loss, and the second is the TD update rule in Equation (4.69) for the SF. In [Car+23b], they present the **Successor Features Keyboard**, that can learn the policy, the SFs and the task encoding \mathbf{z}_w , all simultaneously. They also suggest replacing the squared error regression loss in Equation (4.76) with a cross-entropy loss, where each dimension of the SF is now a discrete probability distribution over M possible values of the corresponding feature. (c.f. Section 7.3.2).

4.5.3.4 Choosing the tasks

A key advantage of SFs is that they provide a way to compute a value function and policy for any given reward, as specified by a task-specific weight vector \mathbf{w} . But how do we choose these tasks? In [Han+19] they sample \mathbf{w} from a distribution at the start of each task, to encourage the agent to learn to explore different parts of the state space (as specified by the feature function ϕ). In [LA21] they extend this by adding an intrinsic reward that favors exploring parts of the state space that are surprising (i.e., which induce high entropy), c.f., Section 7.4.

4.5.4 Successor measures

In this section, we discuss the **successor model** (also called a γ -model, or **geometric horizon models**), which is a probabilistic extension of SR, [JML20]: Rather than just working with expectations, we can simulate future state trajectories by sampling. This allows us to generalize SR to work with continuous states and actions, as we will see.

The basic idea is to define the cumulant as the k -step conditional distribution $C(s_{k+1}) = P(s_{k+1} = \tilde{s} | s_0 = s, \pi)$, which is the probability of being in state \tilde{s} after following π for k steps starting from state s . (Compare this to the SR cumulant, which is $C(s_{k+1}) = \mathbb{I}(s_{k+1} = \tilde{s})$.) The SM is then defined as

$$\boldsymbol{\mu}^\pi(\tilde{s}|s) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_{t+1} = \tilde{s} | s_0 = s) \quad (4.78)$$

where the $1 - \gamma$ term ensures that $\boldsymbol{\mu}^\pi$ integrates to 1. (Recall that $\sum_{t=0}^{\infty} \gamma^t = \frac{1}{1-\gamma}$ for $\gamma < 1$.) In the tabular setting, the SM is just the normalized SR, since

$$\boldsymbol{\mu}^\pi(\tilde{s}|s) = (1 - \gamma) M^\pi(s, \tilde{s}) \quad (4.79)$$

$$= (1 - \gamma) \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t \mathbb{I}(s_{t+1} = \tilde{s}) | s_0 = s, a_{0:\infty} \sim \pi \right] \quad (4.80)$$

$$= (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_{t+1} = \tilde{s} | s_0 = s, \pi) \quad (4.81)$$

Thus $\boldsymbol{\mu}^\pi(\tilde{s}|s)$ tells us the probability that \tilde{s} can be reached from s within a horizon determined by γ when following π , even though we don't know exactly when we will reach \tilde{s} .

SMs obey a Bellman-like recursion

$$\boldsymbol{\mu}^\pi(\tilde{s}|s) = \mathbb{E}[(1 - \gamma)T(\tilde{s}|s, a) + \gamma \boldsymbol{\mu}^\pi(\tilde{s}|s')] \quad (4.82)$$

We can use this to perform policy evaluation by computing

$$V^\pi(s) = \frac{1}{1 - \gamma} \mathbb{E}_{\boldsymbol{\mu}^\pi(\tilde{s}|s)} [R(\tilde{s})] \quad (4.83)$$

We can also define an action-conditioned SM

$$\boldsymbol{\mu}^\pi(\tilde{s}|s, a) = (1 - \gamma) \sum_{t=0}^{\infty} \gamma^t P(s_{t+1} = \tilde{s} | s_0 = s, a_0 = a) \quad (4.84)$$

$$= (1 - \gamma)T(\tilde{s}|s, a) + \gamma \mathbb{E}[\boldsymbol{\mu}^\pi(\tilde{s}|s', a', \pi)] \quad (4.85)$$

Hence we can learn an SM using a TD update of the form

$$\boldsymbol{\mu}^\pi(\tilde{s}|s, a) \leftarrow \boldsymbol{\mu}^\pi(\tilde{s}|s, a) + \underbrace{\eta((1 - \gamma)T(s'|s, a) + \gamma \boldsymbol{\mu}^\pi(\tilde{s}|s', a') - \boldsymbol{\mu}^\pi(\tilde{s}|s, a))}_{\delta} \quad (4.86)$$

where s' is the next state sampled from $T(s'|s, a)$ and a' is the next action sampled from $\pi(s')$. With an SM, we can compute the state-action value for any reward:

$$Q^{R, \pi}(s, a) = \frac{1}{1 - \gamma} \mathbb{E}_{\boldsymbol{\mu}^\pi(\tilde{s}|s, a)} [R(\tilde{s})] \quad (4.87)$$

This can be used to improve the policy as we discuss in Section 4.5.3.1.

4.5.4.1 Learning SMs

Although we can learn SMs using the TD update in Equation (4.86), this requires evaluating $T(s'|s, a)$ to compute the target update δ , and this one-step transition model is typically unknown. Instead, since $\boldsymbol{\mu}^\pi$ is a conditional density model, we will optimize the cross-entropy TD loss [JML20], defined as follows

$$\mathcal{L}_\mu = \mathbb{E}_{(s, a) \sim p(s, a), \tilde{s} \sim (T^\pi \boldsymbol{\mu}^\pi)(\cdot | s, a)} [\log \boldsymbol{\mu}_\theta(\tilde{s} | s, a)] \quad (4.88)$$

where $(T^\pi \mu^\pi)(\cdot|s, a)$ is the Bellman operator applied to μ^π and then evaluated at (s, a) , i.e.,

$$(T^\pi \mu^\pi)(\tilde{s}|s, a) = (1 - \gamma)T(s'|s, a) + \gamma \sum_{s'} T(\tilde{s}|s, a) \sum_{a'} \pi(a'|s') \mu^\pi(\tilde{s}|s', a') \quad (4.89)$$

We can sample from this as follows: first sample $s' \sim T(s'|s, a)$ from the environment (or an offline replay buffer), and then with probability $1 - \gamma$ set $\tilde{s} = s'$ and terminate. Otherwise sample $a' \sim \pi(a'|s')$ and then create a bootstrap sample from the SM using $\tilde{s} \sim \mu^\pi(\tilde{s}|s', a')$.

There are many possible density models we can use for μ^π . In [Tha+22], they use a VAE. In [Tom+24], they use an autoregressive transformer applied to a set of discrete latent tokens, which are learned using VQ-VAE or a non-reconstructive self-supervised loss. They call their method **Video Occupancy Models**. Recently, [Far+25] proposed to use diffusion (flow matching) to learn SMs.

An alternative approach to learning SMs, that avoids fitting a normalized density model over states, is to use contrastive learning to estimate how likely \tilde{s} is to occur after some number of steps, given (s, a) , compared to some randomly sampled negative state [ESL21; ZSE24]. Although we can't sample from the resulting learned model (we can only use it for evaluation), we can use it to improve a policy that achieves a target state (an approach known as goal-conditioned policy learning, discussed in Section 1.2.3).

4.5.4.2 Jumpy models using geometric policy composition

In [Tha+22], they propose **geometric policy composition** or GPC as a way to learn a new policy by sequencing together a set of N policies, as opposed to taking N primitive actions in a row. This can be thought of as a **jumpy model**, since it predicts multiple steps into the future, instead of one step at a time (c.f., [Zha+23a]).

In more detail, in GPC, the agent picks a sequence of n policies π_i for $i = 1 : n$, and then samples states according to their corresponding SMs: starting with (s_0, a_0) , we sample $s_1 \sim \mu_{\gamma'}^{\pi_1}(\cdot|s_0, a_0)$, then $a_1 \sim \pi_1(\cdot|s_1)$, then $s_2 \sim \mu_{\gamma'}^{\pi_2}(\cdot|s_1, a_1)$, etc. This continues for $n - 1$ steps. Finally we sample $s_n \sim \mu_{\gamma'}^{\pi_n}(\cdot|s_{n-1}, a_{n-1})$, where $\gamma' > \gamma$ represents a longer horizon SM. The reward estimates computed along this sampled path can then be combined to compute the value of each candidate policy sequence.

4.5.4.3 Other related work

In [Far+23a], they introduce **proto-value networks**, which is a way to define auxiliary tasks based on successor measures.

In [TO21; TRO23], they propose the **forwards-backwards** representations, which provides a general framework for learning SMs.

4.5.5 Connection between options and successor representations

In Section 7.5.1, we discuss **options**, which are temporally extended actions. These are closely related to the successor representation [Mac+23]. Indeed, the eigenvectors of the SR matrix \mathbf{M}^π correspond to the eigenvectors of the **graph Laplacian**, defined as

$$\mathbf{L} = \mathbf{D}^{-\frac{1}{2}}(\mathbf{D} - \mathbf{A})\mathbf{D}^{-\frac{1}{2}} \quad (4.90)$$

where \mathbf{A} is the adjacency matrix corresponding to T^π , and \mathbf{D} is a diagonal matrix whose entries are the row-sums of \mathbf{A} . These eigenvectors are known as **proto value functions** [MM07], or **eigen-options** [Mac+18b].

We can use this connection to define a “universal” form of successor representation that is independent of a specific policy. Rather than constructing an option for every eigenvector of the graph Laplacian, a single option based on the second eigenvector is sufficient. This is called a **covering option** [Jin+19; Jin+20; KM23], since it minimizes the cover time of the underlying MDP, which loosely refers to how long it takes for a random high-level policy to visit all states.

Chapter 5

Multi-agent RL

In this section, we give a brief introduction to **multi-agent RL** or **MARL**. Our presentation is based on [ACS24]. MARL is closely related to game theory (see e.g. [LBS08]) and multi-agent systems design (see e.g. [SLB08]), as we will see. For other surveys on MARL, see e.g. [HLKT19; YW20; Won+22; GD22].

5.1 Games

Multi-agent environments are often called **games**, even if they represent “real-world” problems such as multi-robot coordination (e.g., a fleet of autonomous vehicles) or agent-based trading. In this section, we discuss different kinds of games that have been proposed, summarized in Figure 5.1.

In the **game theory** community, the rules of the game (i.e., the environment dynamics, and the reward function, aka **payoff function**) are usually assumed known, and the focus is on computing **strategies** (i.e., policies) for each **player** (i.e., agent), whereas in MARL, we usually assume the environment is unknown and the agents have to learn just by interacting with it. (This is analogous to the distinction between DP methods, that assume a known MDP, and RL methods, that just assume sample-based access to the MDP.)

5.1.1 Normal-form games

A **normal-form game** defines a single interaction between $n \geq 2$ agents. In particular, we have a finite set of agents $\mathcal{I} = \{1, \dots, n\}$ (we assume that n is fixed). For each agent $i \in \mathcal{I}$ we have a finite set of actions \mathcal{A}_i and a reward function $\mathcal{R}_i : \mathcal{A}_{1:n} \rightarrow \mathbb{R}$, where $\mathcal{A}_{1:n} = \mathcal{A}_1 \times \dots \times \mathcal{A}_n$. A single round of the game proceeds

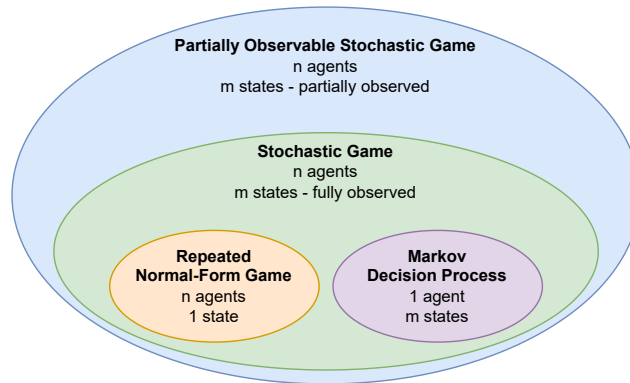


Figure 5.1: Hierarchy of games. From Fig 3.1 of [ACS24]. Used with kind permission of Stefano Albrecht.

as follows. Each agent samples an action $a_i \in \mathcal{A}_i$ with probability $\pi_i(a_i)$, then the resulting **joint action** $\mathbf{a} = (a_1, \dots, a_n)$ is taken and the reward $\mathbf{r} = (r_1, \dots, r_m)$ is given to each player, where $r_i = \mathcal{R}_i(\mathbf{a})$.

Games can be classified based on the type of rewards they contain. In **zero-sum games**, we have $\sum_i \mathcal{R}_i(\mathbf{a}) = 0$ for all \mathbf{a} . (For a two-player zero-sum game, **2p0s**, we must have $R_1(\mathbf{a}) = -R_2(\mathbf{a})$.) In **common-payoff games** (aka common-reward games), we have $\mathcal{R}_i(\mathbf{a}) = \mathcal{R}_j(\mathbf{a})$ for all \mathbf{a} . And in **general-sum games**, there are no restrictions on the rewards.

In zero-sum games, the agents must compete against each other, whereas in common-reward games, the agents generally must cooperate (although they may compete with each other over a shared resource). In general-sum games, there can be a mix of cooperation and competition. Although common-reward games can be easier to solve than general-sum games, it can be challenging to disentangle the contribution of each agent to the shared reward (this is a multi-agent version of the **credit assignment** problem), and coordinating actions across agents can also be difficult.

Normal-form games with 2 agents are called **matrix games** because they can be defined by a 2d reward matrix. We give some well-known examples in Table 5.1.

- In **rock-paper-scissors**, Rock can blunt scissors, Paper can cover rock, and Scissors can cut paper; from these constraints, we can determine which player wins or loses. This is a zero-sum game.
- In the **battle of the sexes** games, a male-female couple want to choose a shared activity. They both have different individual preferences (eg row player prefers Opera, column player prefers Football), both they would both rather spend time together than alone. This is an example of a **coordination game**.
- In the **Prisoner's dilemma**, which is a general-sum game, the players (who are prisoners being interrogated independently in different cells) can either cooperate with each other (by both “staying mum”, i.e., denying they committed the crime), or one can defect on the other (by claiming the other person committed the crime). If they both cooperate, they only have to serve 1 year in jail each, based on weak evidence. If they both defect, they each serve 3 years. But if the row player cooperates (stays silent) and the column player defects (implicates his partner), the row player gets 5 years and the column player gets out of jail free. This leads to an incentive for both players to defect, even though they would be better off if they both cooperated. We discuss this example in more detail in Section 5.2.4.

	R	P	S
R	0,0	-1,1	1,-1
P	1,-1	0,0	-1,1
S	-1,1	1,-1	0,0

(a) Rock-Paper-Scissors

	O	F
O	2,)	0,0
F	0,0	1,2

(b) Battle of the sexes

	C	D
C	-1,-1	-5,0
D	0,5	-3,-3

(c) Prisoner's Dilemma

Table 5.1: Three different matrix games. Here the notation (x, y) in cell (i, j) refers to the row player receiving x and the column player receiving y in response to the joint action (i, j) .

Suppose we consider matrix games with just 2 actions each. In this case, we can represent the game as follows:

$$\begin{pmatrix} a_{11}, b_{11} & a_{12}, b_{12} \\ a_{21}, b_{21} & a_{22}, b_{22} \end{pmatrix} \quad (5.1)$$

where a_{ij} is the reward to player 1 (row player) and b_{ij} is the reward to player 2 (column player) if player 1 picks action i and player 2 picks action j . Suppose we further restrict attention to strictly ordinal games, meaning that each agent ranks the 4 possible outcomes from 1 (least preferred) to 4 (most preferred)). In this case, there are 78 structurally distinct games [RG66]. These can be grouped into two main kinds. In

no-conflict games, both players have the same set of most preferred outcomes, whereas in **conflict games**, the players disagree about what is best. If we consider general ordinal 2×2 games (where one or both players may have equal preference for two or more outcomes), we find that there are 726 of them [KF88].

A **repeated matrix game** is the multi-agent analog of a multi-armed bandit problem, discussed in Section 1.2.5. In this case, the policy has the form $\pi_i(a_t^i | \mathbf{h}_t)$, where $\mathbf{h}_t = (\mathbf{a}_0, \dots, \mathbf{a}_{t-1})$ is the history of joint-actions. In some cases, the agent may choose to ignore the history, or only look at the last n joint actions. For example, in the **tit-for-tat** strategy in the prisoner’s dilemma, the policy for agent i at step t is to do the same action that agent $-i$ did at step $t-1$ (where $-i$ means the agent other than i), so the policy is conditioned on \mathbf{a}_{t-1} . (Note that this strategy will punish players who defect, and can lead to the evolution of cooperative behavior, even in selfish agents [AH81; Axe84].)

5.1.2 Stochastic games

A **stochastic game** is a multi-agent version of an MDP, and was first proposed in [Sha53b]. It is defined by a finite set of agents $\mathcal{I} = \{1, \dots, n\}$; a finite set of states \mathcal{S} , of which a subset $\bar{\mathcal{S}} \subset \mathcal{S}$ are terminal; a finite action set \mathcal{A}_i for each agent $i \in \mathcal{I}$; a reward function $\mathcal{R}_i(s, a, s')$ for each agent $i \in \mathcal{I}$ ¹; a state transition distribution $\mathcal{T}(s_{t+1} | s_{1:t}, \mathbf{a}_t) \in [0, 1]$; and an initial state distribution $\mu(s_0) \in [0, 1]$. Typically the transition distribution is Markovian (i.e., $\mathcal{T}(s_{t+1} | s_{1:t}, \mathbf{a}_t) = \mathcal{T}(s_{t+1} | s_t, \mathbf{a}_t)$, in which case this is called a **Markov game** [Lit94].) See Figure 5.2 for an example.

The policy for each agent in such a game has the form $\pi_i(a_t^i | \mathbf{h}_t)$ where $\mathbf{h}_t = (s_0, \mathbf{a}_1, \dots, s_t)$ is the state-action history. (We omit rewards from the definition of history for notational simplicity.) The overall **joint policy** is denoted by $\boldsymbol{\pi} = (\pi_1, \dots, \pi_n)$; if the agents make their decisions independently (which we assume), then this has the form

$$\boldsymbol{\pi}(\mathbf{a}_t | \mathbf{h}_t) = \prod_i \pi_i(a_t^i | \mathbf{h}_t) \quad (5.2)$$

Often we assume the policies are Markovian, in which case they can be written as $\pi_i(a_t^i | s_t)$.

Note that, from the perspective of each agent i , the environment transition function has the form

$$\mathcal{T}_i(s_{t+1} | s_t, a_t^i) = \sum_{\mathbf{a}_t^{-i}} \mathcal{T}(s_{t+1} | s_t, (a_t^i, \mathbf{a}_t^{-i})) \prod_{j \neq i} \pi_j(a_t^j | s_t) \quad (5.3)$$

Thus \mathcal{T}_i depends on the policies of the other players, which are often changing, which makes these local/agent-centric transition matrices non-stationary, even if the underlying environment is stationary. Typically agent i does not know the policies of the other agents j , so it has to learn them, or it can just treat the other agents as part of the environment (i.e., as another source of unmodeled “noise”) and then use single agent RL methods (see Section 5.3.2).

5.1.3 Partially observed stochastic games (POSG)

A **Partially Observed Stochastic Game** or **POSG** is a multi-agent version of a POMDP (See e.g., [HBZ04].) We augment the stochastic game with the observation distributions $\mathcal{O}_i(o_{t+1}^i | s_{t+1}, \mathbf{a}_t) \in [0, 1]$ for each agent i . (Alternatively, the i ’th observation distribution may just depend on i ’s actions.) Let $\mathbf{o}_t = (o_t^1, \dots, o_t^n)$ be the **joint observation** generated by the product distribution $\mathcal{O}_{1:n}(\mathbf{o}_t | s_t, \mathbf{a}_{t-1})$. The policy for each agent in such a game has the form $\pi_i(a_t^i | \mathbf{h}_t^i)$ where $\mathbf{h}_t^i = (o_0^i, a_0^i, o_1^i, a_1^i, \dots, o_t^i)$ is the **action observation history** for agent i , and $\mathbf{h}_t = \mathbf{h}_t^{1:n}$ is the joint observation history. (Note that the environment decides what is included in each observation; for example, it may or may not contain information about the other agent’s actions.) Note that a **Decentralized POMDP** or **Dec-POMDP** is a special case of a POMDP where the reward function is the same for all agents (thus it can only capture cooperative behavior). See [OA16] for more details.

¹Here $R(s, a, s')$ is the reward we receive if we take action a in state s and end up in s' . As explained in [ACS24, Sec 2.8], we can convert from $R(s, a, s')$ to the more common $R(s, a)$ notation, representing the expected reward, by noting that $R(s, a) = \sum_{s'} \mathcal{T}(s' | s, a) R(s, a, s')$.

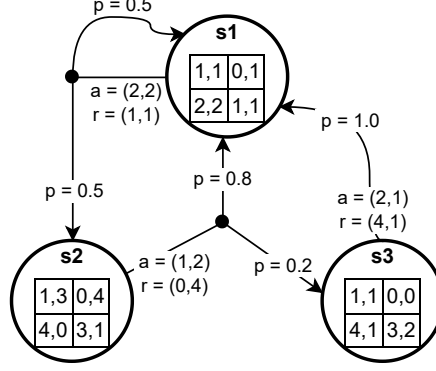


Figure 5.2: Example of a two-player general-sum stochastic game. Circles represent the states, inside of which we show the reward function in matrix form. Only one of the 4 possible transitions out of each state are shown. The little black dots are called the **after states**, and correspond to an intermediate point where a joint action has been decided by the players, but nature hasn't yet sampled the next transition, which occurs with the specified probabilities. From Fig 3.3(b) of [ACS24]. Used with kind permission of Stefano Albrecht.

5.1.3.1 Data generating process

The data generating process for a POSG proceeds as follows. First the environment samples an initial state from $\mu(s_0)$ and generates an initial observation from $\mathcal{O}_{1:n}^0(o_0|s_0)$. Then for $t = 0, 1, \dots$, we repeat the following

1. Each agent generates action from $\pi_i(a_t^i|h_t^i)$.
2. Environment generates next state from $\mathcal{T}(s_{t+1}|s_t, \mathbf{a}_t)$.
3. Environment generates observation from $\mathcal{O}_i(o_{t+1}^i|s_{t+1}, \mathbf{a}_t)$ for each i .
4. Environment generates reward from $\mathcal{R}_i(s_t, \mathbf{a}_t, s_{t+1})$ for each i . (For simplicity, we often just use $\mathcal{R}_i(s_t, \mathbf{a}_t)$, which we define as $\mathbb{E}_{s_{t+1}}[\mathcal{R}_i(s_t, \mathbf{a}_t, s_{t+1})]$.)
5. Each agent updates its history using $h_{t+1}^i = f_i(h_t^i, a_t^i, o_{t+1}^i)$, where f_i is agent i 's update function (e.g., list concatenation, or an RNN).

5.1.3.2 Objective

We define the sum of rewards as $G = \sum_t \mathcal{R}_i(s_t, \mathbf{A}_t)$, where we use capital letters for random variables, and bold face for everything that is joint across all agents. The objective of player i is to maximize $\mathcal{J}_i(\pi_i) = \mathbb{E}_{\pi_i}[G]$. We can compute this using Bellmans equations, as follows. First define the expected state value for agent i under joint policy π given joint history \mathbf{h}_t as

$$v_i^\pi(\mathbf{h}_t) = \mathbb{E}_\pi[G_{\geq t}|\mathbf{h}_t] = \mathbb{E}_\pi\left[\sum_{t' \geq t} \mathcal{R}_i(s_{t'}, \mathbf{A}_{t'})|\mathbf{h}_t\right] \quad (5.4)$$

The expected state value for agent i under the joint policy π and its local history h_t^i is

$$v_i^\pi(h_t^i) = \mathbb{E}_\pi[v_i^\pi(\mathbf{H}_t)|h_t^i] \quad (5.5)$$

Similarly define the expected state-action value given the joint history as

$$q_i^\pi(\mathbf{h}_t, a_t^i) = \mathbb{E}_\pi[G_{\geq t}|\mathbf{h}_t, a_t^i] = \mathbb{E}_\pi[\mathcal{R}_i(s_t, \mathbf{A}_t) + v_i^\pi(\mathbf{H}_{t+1})|\mathbf{h}_t, a_t^i] \quad (5.6)$$

and the expected state-action value given the local history as

$$q_i^\pi(h_t^i, a_t^i) = \mathbb{E}_\pi[q_i^\pi(\mathbf{H}_t, a_t^i)|h_t^i] \quad (5.7)$$

5.1.3.3 Single agent perspective

From the perspective of agent i , it just observes a sequence of observations generated by the following “sensor stream distribution”, (which is non-Markovian [LMLFP11]):

$$p_i(o_{t+1}^i | \mathbf{h}_t^i, \mathbf{a}_t^i) = \sum_{s_{t+1}} \sum_{\mathbf{a}_t^{-i}} \hat{\mathcal{O}}_i(o_{t+1}^i | s_{t+1}, \mathbf{a}_t) p_i(\mathbf{a}_t^{-i} | \mathbf{h}_t^i) p_i(s_{t+1} | \mathbf{h}_t^i, \mathbf{a}_t) \quad (5.8)$$

$$p_i(\mathbf{a}_t^{-i} | \mathbf{h}_t^i) = \prod_{j \neq i} \hat{\pi}_i^j(\mathbf{a}_t^j | \mathbf{h}_t^i) \quad (5.9)$$

$$p_i(s_{t+1} | \mathbf{h}_t^i, \mathbf{a}_t) = \sum_{s_t} \hat{\mathcal{T}}_i(s_{t+1} | s_t, \mathbf{a}_t) b_i(s_t | \mathbf{h}_t^i) \quad (5.10)$$

where $\hat{\pi}_i^j$ in Equation (5.9) is i ’s estimate of j ’s policy; $\hat{\mathcal{T}}_i$ is i ’s estimate of \mathcal{T} based on \mathbf{h}_t^i ; $\hat{\mathcal{O}}_i$ is i ’s estimate of \mathcal{O}_i based on \mathbf{h}_t^i ; and $b_i(s_t | \mathbf{h}_t^i)$ is i ’s **belief state** (i.e., its posterior distribution over the underlying latent state given its local observation history). The agent can either learn a policy given this “collapsed” representation, treating the other agents as part of the environment, or it can explicitly try to learn the true joint world model \mathcal{T} , local observation model \mathcal{O}_i and other agent policies π_i^j , so it can reason about the other agents. In this section, we follow the latter approach.

5.1.3.4 Factored Observation Stochastic Games (FOSG)

[Kov+22] propose a formalism called **Factored Observation Stochastic Games** or **FOSG** that extends POSGs by partitioning the observation for each player into public and private.² (We say that information is public if it is visible to all players, and all players know this; thus it is a form of **common knowledge**.) Explicitly distinguishing these two kinds of information is important in order to tractably solve certain kinds of games, like Poker or Hanabi (see e.g., [Sok+21]).

5.1.4 Extensive form games (EFG)

In the game theory literature, it is common to use the **extensive form game** representation. Rather than representing a sequence of world states that evolve over time, we represent a tree of possible choices or actions take by each player (and optionally a **chance player**, if the game is stochastic, e.g., backgammon). Each node represents a unique sequence (history) of actions leading up to that point.

In the context of EFGs, some additional terminology is commonly used. If all the nodes are observed (including chance nodes), we say the game has perfect and complete information. If the moves of some players are not visible and/or the state of the game is not fully known (e.g., poker), the game has **imperfect information**. In this case, we define an **information set** as the set of nodes that an agent cannot distinguish between. This is analogous to having a distribution over the hidden states in a POSG.

If an agent does not know the other player’s type or payoff function (e.g. in an auction, or playing against players with unknown skill level), then the game has **incomplete information**. In this case, the agent should maintain a Bayesian belief state about the unknown factors. This is analogous to having a distribution over the parameters of the POSG itself, similar to a multi-agent version of a Bayes Adaptive POMDP [RCdP07].

Note that in theoretical work, a useful result is it is possible to convert any EFG into an equivalent (stateless) NFG, where the actions of the NFG correspond to the deterministic policies of the EFG, and the payoffs for a joint action are the expected returns of the corresponding joint policy in the EFG.

5.1.4.1 Example: Kuhn Poker as EFG

In this section, we give an example of an EFG formulation of the game of **Kuhn Poker**, introduced in [Kuh51]. We first define the rules of the game, following [Kov+22]:

²Note that this kind of factorization is different from factoring the state vector or reward function; the latter is called a factored POSG.

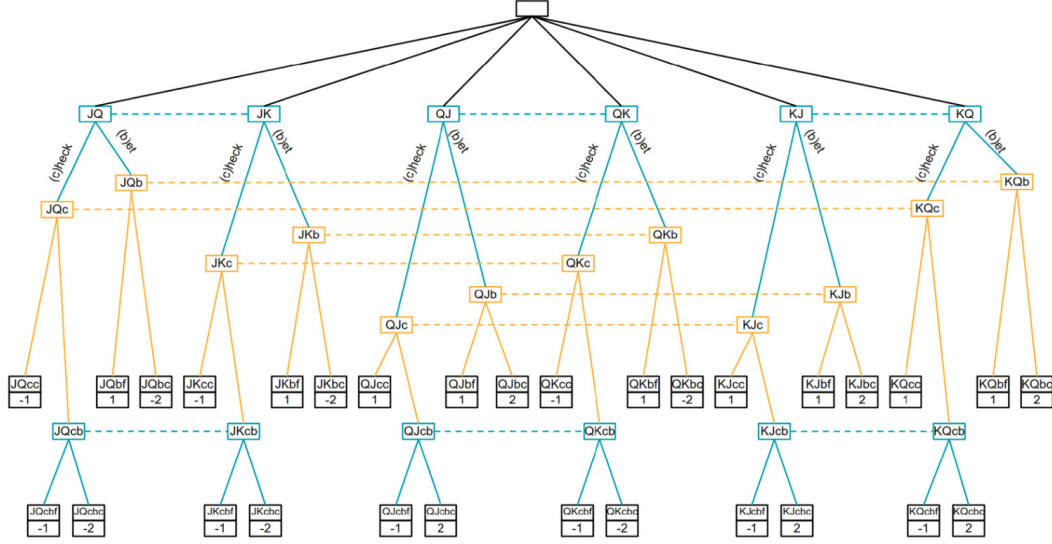


Figure 5.3: EFG for Kuhn Poker. The dashed lines connect histories in the same information set. From Fig 2 of [Kov+22]. Used with kind permission of Vilim Lisy.

Kuhn poker is a form of (two player) poker where the deck includes only three cards: Jack, Queen, and King. First, each player places one chip into the pot as the initial forced bet (ante). Each player is then privately dealt one card (the last card isn't revealed). This is followed by a betting phase (explained below). The game ends either when one player folds (forfeiting all bets made so far to their opponent) or there is a showdown, where the private cards are revealed and the higher card's owner receives the bets. At the start of the betting, player one can either check/pass or raise/bet (one chip). If they check, player two can also check/pass — leading to a showdown — or bet. If one of the players bets, their opponent must either call (betting one chip to match the opponent's bet), followed by a showdown, or fold.

The EFG for this is shown in Figure 5.3. To interpret this figure, consider the left part of the tree, where the state of nature is JQ (as determined by the chance player's first two dealing actions). Suppose the first player checks, leading to state JQc. The second player can either check, leading to JQcc, resulting in a showdown with a reward of -1 to player 1 (since $J < Q$); or bet, leading to JQcb. In the latter case, player 1 must then either fold, leading to JQcbf with a reward to player 1 of -1 (since player 1 only put in one chip); or player 1 must call, leading to JQcbcb with a reward to player 1 of -2 (since player 1 put in two chips).

5.1.4.2 Converting FOSG to EFG

We can convert an FOSG into an EFG by “unrolling it”. First we define the information set for a given information state as the set of consistent world state trajectories. By applying the policy of each agent to the world model, we can derive a tree of possible world states (trajectories) and corresponding information sets for each agent, and thus can construct a corresponding (augmented) EFG. See [Kov+22] for details.

5.2 Solution concepts

In the multi-agent setting the definition of “optimality” is much more complex than in the single agent setting, as we will see. That is, there are multiple **solution concepts**.

5.2.1 Notation and definitions

First we define some notation. Let $\hat{\mathbf{h}}_t = \{(s_k, \mathbf{o}_k, \mathbf{a}_k)_{k=1}^{t-1}, s_t, \mathbf{o}_t\}$ be the **full history**, containing all the past states, joint observations, and joint actions. Let $\sigma(\hat{\mathbf{h}}_t) = \mathbf{h}_t = (\mathbf{o}_1, \dots, \mathbf{o}_t)$ be the history of joint observations, and $\sigma_i(\hat{\mathbf{h}}_t) = \mathbf{h}_t^i = (o_1^i, \dots, o_t^i)$ be the history of observations for agent i . (This typically also includes the actions chosen by agent i .)

We define the expected return for agent i under joint policy $\boldsymbol{\pi}$ by

$$U_i(\boldsymbol{\pi}) = \sum_{\hat{\mathbf{h}}_t} p(\hat{\mathbf{h}}_t | \boldsymbol{\pi}) u_i(\hat{\mathbf{h}}_t) \quad (5.11)$$

where the distribution over full histories is given by

$$p(\hat{\mathbf{h}}_t | \boldsymbol{\pi}) = \mu(s_0) \mathcal{O}_{1:n}^0(\mathbf{o}_0 | s_0) \prod_{k=1}^{t-1} \pi(\mathbf{a}_k | \hat{\mathbf{h}}_k) \mathcal{T}(s_{k+1} | s_k, \mathbf{a}_k) \mathcal{O}_{1:n}(\mathbf{o}_{k+1} | s_{k+1}, \mathbf{a}_k) \quad (5.12)$$

and $u_i(\hat{\mathbf{h}}_t)$ is the discounted actual return for agent i in a given full history

$$u_i(\hat{\mathbf{h}}_t) = \sum_{k=0}^{t-1} \gamma^k \mathcal{R}_i(s_k, \mathbf{a}_k, s_{k+1}) \quad (5.13)$$

We can also derive the following Bellman-like equations:

$$V_i^\pi(\hat{\mathbf{h}}) = \sum_{\mathbf{a}} \pi(\mathbf{a} | \sigma(\hat{\mathbf{h}})) Q_i^\pi(\hat{\mathbf{h}}, \mathbf{a}) \quad (5.14)$$

$$Q_i^\pi(\hat{\mathbf{h}}, \mathbf{a}) = \sum_{s'} \mathcal{T}(s' | s(\hat{\mathbf{h}}), \mathbf{a}) \left[\mathcal{R}_i(s(\hat{\mathbf{h}}), \mathbf{a}, s') + \gamma \sum_{\mathbf{o}'} \mathcal{O}_{1:n}(\mathbf{o}' | \mathbf{a}, s') V_i^\pi((\hat{\mathbf{h}}, \mathbf{a}, s', \mathbf{o}')) \right] \quad (5.15)$$

where $s(\hat{\mathbf{h}})$ extracts the last state from $\hat{\mathbf{h}}$. With this, we can define the expected return using

$$U_i(\boldsymbol{\pi}) = \mathbb{E}_{\mu(s_0) \mathcal{O}_{1:n}^0(\mathbf{o}_0 | s_0)} [V_i^\pi((s_0, \mathbf{o}_0))] \quad (5.16)$$

Finally, we define the **best response policy** for agent i as the one that maximizes the expected return for agent i against a given set of policies for all the other agents, $\boldsymbol{\pi}_{-i} = (\pi_1, \dots, \pi_{i-1}, \pi_{i+1}, \dots, \pi_n)$. That is,

$$\text{BR}_i(\boldsymbol{\pi}_{-i}) = \underset{\pi_i}{\operatorname{argmax}} U_i((\pi_i, \boldsymbol{\pi}_{-i})) \quad (5.17)$$

5.2.2 Minimax

The **minimax** solution is defined for two-agent zero-sum games. Its existence for normal-form games was first proven by John von Neumann in 1928. We say that joint policy $\boldsymbol{\pi} = (\pi_i, \pi_j)$ is a minimax solution if

$$U_i(\boldsymbol{\pi}) = \max_{\pi'_i} \min_{\pi'_j} U_i(\pi'_i, \pi'_j) \quad (5.18)$$

$$= \min_{\pi'_j} \max_{\pi'_i} U_i(\pi'_i, \pi'_j) \quad (5.19)$$

$$= -U_j(\boldsymbol{\pi}) \quad (5.20)$$

In other words, $\boldsymbol{\pi}$ is a minimax solution iff $\pi_i \in \text{BR}_i(\pi_j)$ and $\pi_j \in \text{BR}_j(\pi_i)$. We can solve for the minimax solution using linear programming.

Minimax solutions also exist for two-player zero-sum stochastic games with finite episode lengths, such as chess and Go. In the case of perfect information games, the problems are Markovian, and so dynamic programming can be used to solve them. In generally this may be slow, but minimax search (a depth-limited version of DP that requires a heuristic function) can be used.

5.2.3 Exploitability

In the case of a 2 player zero-sum game, we can measure how close we are to a minimax solution by computing the **exploitability** score, defined as

$$\text{exploitability}(\boldsymbol{\pi}) = \frac{1}{2} [\max_{\pi'_1} J(\pi'_1, \pi_2) - \min_{\pi'_2} J(\pi_1, \pi'_2)] \quad (5.21)$$

where $J(\boldsymbol{\pi})$ is the expected reward for player 1 (which is the loss for player 2). Exploitability is the expected return of π_i playing against a best response to π_i , averaged over both players $i \in 1, 2$. Joint policies with exploitability zero are Nash equilibria (see Section 5.2.4).

Note that computing the exploitability score requires computing a best response to any given policy, which can be hard in general. However, one can use standard deep RL methods for this (see e.g., [Tim+20]).

5.2.4 Nash equilibrium

The **Nash equilibrium** generalized the idea of mutual best response to general-sum games with two or more agents. That is, we say that $\boldsymbol{\pi}$ is a Nash equilibrium (NE) if no agent i can improve its expected returns by changing its policy π_i , assuming the other agents policies remain fixed:

$$\forall, \pi'_i. U_i(\pi'_i, \boldsymbol{\pi}_{-i}) \leq U_i(\boldsymbol{\pi}) \quad (5.22)$$

John Nash proved the existence of such a solution for general-sum non-repeated normal form games in 1950.

Below we discuss the kinds of equilibria that exist for the games shown in Table 5.1.

- For the rock-paper-scissors game, the only NE is the **mixed strategy** (i.e., stochastic policy) where each agent chooses actions uniformly at random, so $\pi_i = (1/3, 1/3, 1/3)$. This yields an expected return of 0.
- For the battle-of-the sexes game, there are two pure strategy Nash equilibria: (Opera, Opera) and (Football, Football). There's also a mixed strategy equilibrium but it involves randomness and gives lower expected payoffs.
- For the Prisoner's Dilemma game, the only NE is the pure strategy of (D,D), which yields an expected return of (-3,-3). Note that this is worse than the maximum possible expected return, which is (-1,-1) given by the strategy of (C,C). However, such a strategy is not an NE, since each player could improve its return if it unilaterally deviates from it (i.e., defects on its partner).

Interestingly, it can be shown that two agents that use rational (Bayesian) learning rules to update their beliefs about the opponent's strategy (based on observed outcomes of earlier games), and then compute a best response to this belief, will eventually converge to Nash equilibrium [KL93].

Computing a Nash equilibrium involves identifying a set of strategies where no player can improve their outcome by unilaterally changing their own strategy. This requires solving a system of simultaneous equations or inequalities, which can be computationally intensive, especially as the number of players and available strategies increases. In fact, for general games, finding a Nash equilibrium is a famously hard problem, falling into a complexity class called **PPAD-complete**, which suggests that there is no efficient, universal algorithm for finding one.

5.2.5 Approximate Nash equilibrium

It is possible to relax the definition of exact inequality by defining an ϵ -Nash equilibrium as a joint policy that satisfies

$$\forall, \pi'_i. U_i(\pi'_i, \boldsymbol{\pi}_{-i}) - \epsilon \leq U_i(\boldsymbol{\pi}) \quad (5.23)$$

Unfortunately, the expected return from a ϵ -Nash equilibrium can be very different from the expected return from a true NE. For example, consider this matrix game:

$$\begin{array}{c|cc}
 & C & D \\
 \hline
 A & 100, 100 & 0, 0 \\
 B & 1, 2 & 1, 1
 \end{array} \tag{5.24}$$

The unique NE is (A,C), but the ϵ -NE with $\epsilon = 1$ is either (A,C) or (B,D), which clearly have very different rewards.

Despite the above drawback, much computational work focuses on approximate Nash equilibria. Indeed we can measure the rate of convergence to such a state by defining

$$\text{NashConv}(\pi) = \sum_i \delta_i(\pi) \tag{5.25}$$

where $\delta_i(\pi)$ is the amount of incentive that i has to deviate to one of its best responses away from the joint policy:

$$\delta_i(\pi) = u_i(\pi_i^b, \pi_{-i}) - u_i(\pi), \pi_i^b \in \text{BR}(\pi_{-i}) \tag{5.26}$$

5.2.6 Entropy regularized Nash equilibria (aka Quantal Response Equilibria)

In this section, we discuss **quantal response equilibria** or **QRE** [MP95; MP98]. These are like Nash equilibria except the best response policy is a “soft” entropy-regularized policy (see Section 3.6.8). This kind of equilibrium reflects the fact that players may not always choose the best response with certainty, but instead they make choices based on a probability distribution over actions, based on the relative expected utility of each action. Thus it is a Bayesian equilibrium. This can be useful for modeling human behavior that deviates from the predictions of Nash (studied in the field of **behavioral game theory**). In addition, it is useful for developing algorithms that converge to a unique equilibrium, as we discuss in Section 5.3.10.1.

For single agent problems with a single state (i.e., bandit problems), we say that a policy π is α -soft optimal in the normal sense if it satisfies

$$\pi = \underset{\pi' \in \Delta(A)}{\text{argmax}} \mathbb{E}_{A \sim \pi'} q(A) + \alpha \mathbb{H}(\pi') \tag{5.27}$$

where $\Delta(A)$ is the action simplex, and q is the action-value function. If this holds for all states (decision points) s , we say that π is α -soft optimal in the behavioral sense.

For two-player zero-sum NFGs (which have a single decision point or state), we say that a policy is a QRE if each player’s policy is soft optimal in the normal sense conditioned on the other player not changing its policy [MP95]. For two-player zero-sum games with multiple states (i.e., EFGs), we say that a policy is a **agent QRE** if each player’s policy is soft optimal in the behavioral sense conditioned on the other player not changing its policy [MP98].

5.2.7 Correlated equilibrium

The concept of a Nash equilibrium assume the policies are independent, which can limit the expected returns. A **correlated equilibrium** (CE) allows for correlated policies. Specifically, we assume there is a central policy π_c that defines a distribution over joint actions. Agents can follow this recommended policy, or can choose to deviate from it by using an action modified $\xi_i : \mathcal{A}_i \rightarrow \mathcal{A}_i$. We then say that π_c is a CE if for all i and ξ_i we have

$$\sum_a \pi_c(a) \mathcal{R}_i((\xi_i(a^i), a^{-i})) \leq \sum_a \pi_c(a) \mathcal{R}_i(a) \tag{5.28}$$

That is, player i has no incentive to deviate from the recommendation, after receiving it. It can be shown that the set of correlated equilibria contains the set of Nash equilibria. In particular, since Nash equilibrium

is a special case of correlated equilibrium in which the joint policy π_c is factored into independent agent policies with $\pi_c(\mathbf{a}) = \prod_i \pi_i(a_i)$.

To see how a correlated equilibrium can give higher returns than a Nash equilibrium, consider the **Chicken game**. This models two agents that are driving towards each other. Each agent can either stay on course (S) or leave (L) and avoid a crash. The payoff matrix is as follows:

	S	L
S	0,0	7,2
L	2,7	6,6

(5.29)

This reflects the fact that if they both stay on course, then they both die and get reward 0; if they both leave, they both survive and get reward 6; but if player i chooses to stay and the other one leaves, then i gets a reward of 7 for being brave, and $-i$ only gets a reward of 2 for chickening out.

We can represent π_i by the scalar $\pi_i(S)$, since $\pi_i(L) = 1 - \pi_i(S)$. Hence π can be defined by the tuple (π_1, π_2) . There are 3 uncorrelated NEs: $\pi = (1, 0)$ with return $(7, 2)$; $\pi = (0, 1)$ with return $(2, 7)$; and $\pi = (\frac{1}{3}, \frac{1}{3})$ with return $(4.66, 4.66)$. There is 1 CE, namely $\pi_c(L, L) = \pi_c(S, L) = \pi_c(L, S) = \frac{1}{3}$ and $\pi_c(S, S) = 0$. The central policy has an expected return of

$$7 \cdot \frac{1}{3} + 2 \cdot \frac{1}{3} + 6 \cdot \frac{1}{3} = 5 \quad (5.30)$$

which we see is higher than the NE of 4.66. This is because it avoids the deadly joint (S,S) action. To show that this is a CE, consider the case where i (e.g., row player) receives recommendation L ; they know that j (column player) will choose either S or L with probability 0.5 (because the central policy is uniform). If i sticks with the recommendation, its expected return is $0.5 \cdot 2 + 0.5 \cdot 6 = 4$; this is greater than deviating from the recommendation and picking S , which has expected return of $0.5 \cdot 0 + 0.5 \cdot 7 = 3.5$. Thus π_c is a CE.

In [Aum87] they show that the CE solution corresponds to the behavior of a rational Bayesian agent. The correlated equilibrium solution can be computed via linear programming.

5.2.8 Limitations of equilibrium solutions

Equilibrium solutions have several limitations. First, they do not always maximize expected returns. For example, in Prisoner's Dilemma, (D,D) is Nash but (C,C) yields higher returns. Second, there can be multiple (even infinitely many) equilibria, each with different expected returns, as we have seen. Third, equilibria for sequential games don't specify what to do if the history deviates from the equilibrium path, i.e., they do not define the policy for full histories where $p(\hat{\mathbf{h}}|\pi) = 0$; this can be problematic when the agents are learning, or the environment is changing in some other way. Consequently it is common to define additional solution requirements, as we discuss below.

5.2.9 Pareto optimality

We say a joint policy π is **Pareto optimal** if it is not **Pareto dominated** by any other joint policy π' . We say that π is Pareto dominated by π' if π' improves the expected return for at least one agent:

$$\forall i. U_i(\pi') \geq U_i(\pi) \text{ and } \exists i. U_i(\pi') > U_i(\pi) \quad (5.31)$$

and if it does not decrease the payoff for any agents.

Figure 5.4 illustrates the expected joint rewards and the Pareto frontier for all feasible policies (up to quantization error) applied to the Chicken game in Equation (5.29). We see that the two pure NEs are on the Pareto frontier (as are many other policies that are not Nash), but the mixed NE is not Pareto optimal.

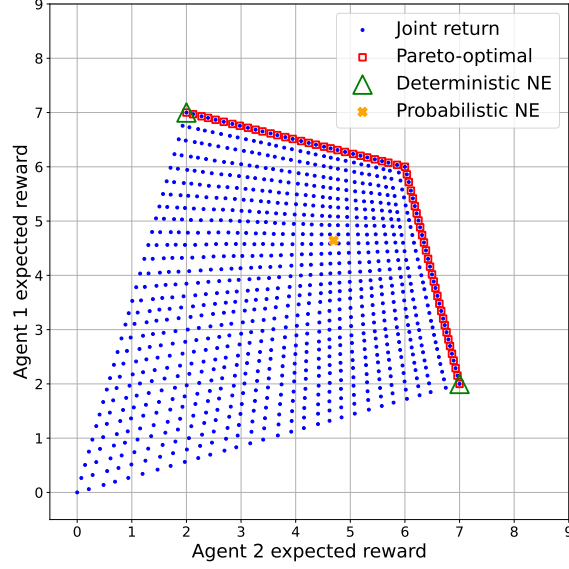


Figure 5.4: Space of (discretized) joint policies for Chicken game. From Fig 4.4 of [ACS24]. Used with kind permission of Stefano Albrecht.

5.2.10 Social welfare and fairness

Pareto optimality ensures there is no other solution in which at least one agent is better off, without making other agents worse off. However, it does not make any guarantees about the total rewards, or their distribution amongst agents. For example, along the Pareto frontier in Figure 5.4, the joint returns vary from (7,2) to (6,6) to (2,7).

To further constrain the space of desirable solutions, we can consider additional concepts. For example, we define **welfare optimality** as

$$W(\pi) = \sum_i U_i(\pi) \quad (5.32)$$

A joint policy is welfare-optimal if $\pi \in \operatorname{argmax}_{\pi'} W(\pi')$. One can show that welfare optimality implies Pareto optimality, but not (in general) vice versa.

Similarly, we define **fairness optimality** as

$$F(\pi) = \prod_i U_i(\pi) \quad (5.33)$$

A joint policy is fairness-optimal if $\pi \in \operatorname{argmax}_{\pi'} F(\pi')$.

In the battle-of-the-sexes game in Table 5.1, the only fair outcome is the joint distribution over $Pr(F, F) = Pr(O, O) = 0.5$, which means the couple spend half their time watching football and half going to the opera.

In the Chicken game in Figure 5.4, there is only one solution that is both welfare-optimal and fairness-optimal, namely the joint policy with expected return of (6,6). Note, however, that this is not a Nash policy.

5.2.11 No regret

The quantity known as **regret** measures the difference between the rewards an agent received versus the maximum rewards it could have received if it had chosen a different action. For a non-repeated normal-form

general-sum game, played over E episodes, this is defined as

$$\text{Regret}_i^E = \max_{a^i} \sum_{e=1}^E [\mathcal{R}_i((a^i, \mathbf{a}_e^{-i})) - \mathcal{R}_i(\mathbf{a}_e)] \quad (5.34)$$

We can generalize the definition of regret to stochastic games and POSGs by defining the regret over policies instead of actions. That is,

$$\text{Regret}_i^E = \max_{\pi^i} \sum_{e=1}^E [U_i((\pi^i, \pi_e^{-i})) - U_i(\pi_e)] \quad (5.35)$$

In all these cases, an agent is said to have no-regret if

$$\forall i. \lim_{E \rightarrow \infty} \frac{1}{E} \text{Regret}_i^E \leq 0 \quad (5.36)$$

5.2.12 Shapley values

The **Shapley value** [Sha53a] allows one to estimate the marginal contribution of a single agent to a common reward (cooperative) game, which can help ameliorate the **credit assignment problem**. Specifically, suppose there are N players, and let S be a subset of $\leq N$ players that form a team. Let $v(S)$ be the expected value obtained by that team, and $v(S \setminus \{i\})$ be the value of the team when i is absent. Then we define i 's Shapley value as

$$\phi(i) = \sum_{S \subseteq N \setminus \{i\}} w(S) [v(S \cup \{i\}) - v(S)] \quad (5.37)$$

where the weighting term is given by

$$w(S) = \frac{|S|!(|N| - |S| - 1)!}{|N|!} = \frac{1}{|N| \cdot \binom{|N|-1}{|S|}} \quad (5.38)$$

which represents the probability that, for a random ordering of all players, the players in coalition S come before player i , and the remaining players come after player i .

An interesting application of Shapley values arises in the “explainable AI” literature, where one of the goals is to estimate the importance of individual predictors to an overall prediction. This can be done by using the **SHAP** (SHapley Additive exPlanations) framework of [LL17]. Similarly, people use Shapley values to assess the significance of individual tokens in a prompt (see e.g., [Eno+24]).

5.2.13 Stackelberg equilibrium

So far we have mostly focused on games where the players make their moves simultaneously. We can also consider **sequential games**, where the players take turns. For such games, the concept of **Stackelberg equilibrium**, developed by economist Heinrich Freiherr von Stackelberg, becomes relevant. Here, one player, called the “leader”, makes their decision first. The other player, the “follower”, observes the leader’s choice and then makes their own decision. This sequential structure gives the leader a significant advantage. The leader, knowing how the follower will react to any given move, can choose a strategy that maximizes their own payoff, anticipating the follower’s subsequent best response. The follower, in turn, optimizes their outcome based on the leader’s committed action. The resulting set of strategies, where the leader has chosen their optimal move and the follower has responded optimally, constitutes the Stackelberg equilibrium.

The concept of Stackelberg equilibrium is useful for analysing **setter-solver** problems, which arise in unsupervised environment design/ curriculum learning (Section 7.4.2.2). Here one player (the setter) designs a task or environment that is challenging for the other player (the solver); as the solver gets better, the setter gradually makes the problem more challenging. See e.g., [Den+20; Rac+20; PH+23].

Note that Stackelberg equilibria can be easier to compute than Nash equilibria. To find the Stackelberg equilibrium, you can use a method called **backward induction**. You start by determining the follower’s

best response for every possible action the leader could take. Then, the leader, knowing how the follower will react, simply chooses the action that will lead to the best outcome for themselves. This process is often straightforward and can be solved with a single optimization problem, including using gradient-based methods [CMS07; ZS22; GP23].

5.3 Algorithms

In this section, we discuss various MARL algorithms.

5.3.1 Centralized learning

The simplest way to solve a MARL problem is to reduce it to a single agent RL (SARL) problem. In **central learning**, we learn a single joint policy over the joint action space. This requires that we can transform the joint reward $\mathbf{r}_t = (r_t^1, \dots, r_t^n)$ into a scalar r_t . This is easy to do in common reward games, where the agents must cooperate. However, for general sum games, it may be impossible to define a single scalar reward across all agents. And even if we can define such a shared reward, the resulting method may not scale well with the number of agents, and learns a policy that requires global access to all of the observations for each agent.

5.3.2 Independent learning

In **independent learning**, each agent treats all other agents as part of the environment, and then uses any standard single-agent RL algorithm for training. This is done in parallel across all agents.

5.3.2.1 Independent Q learning

For example, if we use Q learning for each agent, the method is known as independent Q-learning or **IQL**; see Algorithm 18 for the pseudocode.

Algorithm 18: Independent Q learning (DQN for multiple independent agents)

```

1 Initialize  $n$  value networks with random parameters  $\theta_1, \dots, \theta_n$ ;
2 Initialize  $n$  target networks with parameters  $\bar{\theta}_1 = \theta_1, \dots, \bar{\theta}_n = \theta_n$ ;
3 Initialize a replay buffer for each agent  $D_1, D_2, \dots, D_n$ ;
4 for time step  $t = 0, 1, 2, \dots$  do
5   Collect current observations  $o_t^1, \dots, o_t^n$ ;
6   for agent  $i = 1, \dots, n$  do
7     With probability  $\epsilon$ : choose random action  $a_t^i$ ;
8     Otherwise: choose  $a_t^i \in \operatorname{argmax}_{a_i} Q(h_t^i, a_i; \theta_i)$ ;
9   Apply actions  $(a_t^1, \dots, a_t^n)$ ; collect rewards  $r_t^1, \dots, r_t^n$  and next observations  $o_{t+1}^1, \dots, o_{t+1}^n$ ;
10  for agent  $i = 1, \dots, n$  do
11    Store transition  $(h_t^i, a_t^i, r_t^i, h_{t+1}^i)$  in replay buffer  $D_i$ ;
12    Sample random mini-batch of  $B$  transitions  $(h_k^i, a_k^i, r_k^i, h_{k+1}^i)$  from  $D_i$ ;
13    if  $s_{k+1}^i$  is terminal then
14      Targets  $y_k^i \leftarrow r_k^i$ ;
15    else
16      Targets  $y_k^i \leftarrow r_k^i + \gamma \max_{a'_i \in \mathcal{A}_i} Q(h_{k+1}^i, a'_i; \bar{\theta}_i)$ ;
17    Loss  $\mathcal{L}(\theta_i) \leftarrow \frac{1}{B} \sum_{k=1}^B \left( y_k^i - Q(h_k^i, a_k^i; \theta_i) \right)^2$ ;
18    Update parameters  $\theta_i$  by minimizing the loss  $\mathcal{L}(\theta_i)$ ;
19    In a set interval, update target network parameters  $\bar{\theta}_i$ ;

```

5.3.2.2 Independent Actor Critic

Instead of using value-based methods, we can also use policy learning methods. The multi-agent version of the policy gradient theorem in Equation (3.6) is the following (see e.g., [ACS24] for derivation):

$$\nabla_{\theta_i} J(\theta_{1:n}) \propto \mathbb{E}_{\hat{\mathbf{h}} \sim p(\hat{\mathbf{h}}|\pi), a^i \sim \pi_i, \mathbf{a}^{-i} \sim \pi^{-i}} \left[Q_i^\pi(\hat{\mathbf{h}}, (a^i, \mathbf{a}^{-i})) \nabla_{\theta_i} \log \pi(a_i | h_i = \sigma_i(\hat{\mathbf{h}}); \theta_i) \right] \quad (5.39)$$

where $\hat{\mathbf{h}}_t = \{(s_k, \mathbf{o}_k, \mathbf{a}_k)_{k=1}^{t-1}, s_t, \mathbf{o}_t\}$ is the **full history** (containing all the past states, joint observations, and joint actions), and $\sigma_i(\hat{\mathbf{h}}_t) = \mathbf{h}_t^i = (o_1^i, \dots, o_t^i)$ is the history of observations for agent i .

In practice, we usually subtract a baseline term from Q , to reduce the variance. If we use the value function for the baseline, then the first term inside the expectation becomes

$$Q_i^\pi(\hat{\mathbf{h}}, (a^i, \mathbf{a}^{-i})) - V_i^\pi(\hat{\mathbf{h}}) = \text{Adv}_i^\pi(\hat{\mathbf{h}}, \mathbf{a}) \quad (5.40)$$

where Adv is the advantage, as we discussed in Section 3.2.1. This can be used inside a multi-agent version of the advantage actor critic or A2C method (known as **MAA2C**) shown in Algorithm 19. (To combat the fact that we cannot use replay buffers with an on-policy method, we assume instead that we can parallelize over multiple (synchronous) environments, to ensure we have a sufficiently large minibatch to estimate the loss function at each step.)

Algorithm 19: Multi-agent Advantage Actor-Critic (MAA2C)

```

1 Initialize  $n$  actor networks with random parameters  $\theta_1, \dots, \theta_n$ 
2 Initialize  $n$  critic networks with random parameters  $w_1, \dots, w_n$ 
3 Initialize  $K$  parallel environments
4 Initialize histories  $h_0^{i,k}$  for each agent  $i$  and environment  $k$ 
5 for time step  $t = 0 \dots$  do
6   for environment  $k = 1, \dots, K$  do
7     Sample actions:  $\{a_t^{i,k} \sim \pi(\cdot | h_t^{i,k}, \theta^i)\}_{i=1}^n$ 
8     Sample next state:  $s_{t+1}^k \sim T(\cdot | s_t^k, \mathbf{a}_t^{1:n,k})$ 
9     Sample observations:  $\{o_{t+1}^{i,k} \sim O_i(\cdot | s_t^k, a_t^{i,k})\}_{i=1}^n$ 
10    Sample rewards:  $\{r_t^{i,k} \sim R_i(\cdot | s_t^k, \mathbf{a}_t^{1:n,k}, s_{t+1}^k)\}_{i=1}^n$ 
11    Update histories:  $\{h_{t+1}^{i,k} = (h_t^{i,k}, o_{t+1}^{i,k})\}_{i=1}^n$ 
12  for agent  $i = 1, \dots, n$  do
13    if  $s_{t+1}^k$  is terminal then
14       $\text{Adv}(h_t^{i,k}, a_t^{i,k}) \leftarrow r_t^{i,k} - V(h_t^{i,k}; w_i);$ 
15      Critic target  $y_t^{i,k} \leftarrow r_t^{i,k};$ 
16    else
17       $\text{Adv}(h_t^{i,k}, a_t^{i,k}) \leftarrow r_t^{i,k} + \gamma V(h_{t+1}^{i,k}; w_i) - V(h_t^{i,k}; w_i);$ 
18      Critic target  $y_t^{i,k} \leftarrow r_t^{i,k} + \gamma V(h_{t+1}^{i,k}; w_i);$ 
19    Actor loss:
20      Critic loss:
21      Update parameters  $\theta_i$  by minimizing the actor loss  $\mathcal{L}(\theta_i);$ 
22      Update parameters  $w_i$  by minimizing the critic loss  $\mathcal{L}(w_i);$ 

```

$$\mathcal{L}(\theta_i) \leftarrow \frac{1}{K} \sum_{k=1}^K \text{Adv}(h_t^{i,k}, a_t^{i,k}) \log \pi(a_t^{i,k} | h_t^{i,k}; \theta_i)$$

$$\mathcal{L}(w_i) \leftarrow \frac{1}{K} \sum_{k=1}^K \left(y_t^{i,k} - V(h_t^{i,k}; w_i) \right)^2$$

5.3.2.3 Independent PPO

We can implement an independent version of PPO (known as **IPPO**) in a similar way, by updating all the policies in parallel [Wit+20].

5.3.2.4 Learning dynamics of multi-agent policy gradient methods

In general, applying policy gradient methods to multiple agents in parallel may not result in convergence [CP19; Blo+15]. To illustrate this, consider a non-repeated normal-form general-sum game with two players and two actions. (This is an imperfect information game since each player does not know the other's actions when they make their decision.) Denote the reward matrices by

$$\mathcal{R}_i = \begin{pmatrix} r_{11} & r_{12} \\ r_{21} & r_{22} \end{pmatrix}, \mathcal{R}_j = \begin{pmatrix} c_{11} & c_{12} \\ c_{21} & c_{22} \end{pmatrix} \quad (5.41)$$

Denote the policies by

$$\pi^i = (\alpha, 1 - \alpha), \pi^j = (\beta, 1 - \beta) \quad (5.42)$$

The expected reward for agent i , given the joint policy $\pi = (\alpha, \beta)$, is given by

$$U_i(\alpha, \beta) = \alpha\beta r_{11} + \alpha(1 - \beta)r_{12} + (1 - \alpha)\beta r_{21} + (1 - \alpha)(1 - \beta)r_{22} \quad (5.43)$$

The expression $U_j(\alpha, \beta)$ is analogous, with r_{ij} replaced with c_{ij} . (Note that computing U_i requires knowledge of \mathcal{R}_i but also of π_j (and vice versa for computing U_j); we will relax this assumption below.) Finally, we can learn the policies using gradient ascent:

$$\alpha_{k+1} = \alpha_k + \kappa \frac{\partial U_i(\alpha_k, \beta_k)}{\partial \alpha_k}, \beta_{k+1} = \beta_k + \kappa \frac{\partial U_j(\alpha_k, \beta_k)}{\partial \beta_k} \quad (5.44)$$

where κ is the learning rate.

We can analyse the dynamics of the above procedure as $\kappa \rightarrow 0$; this is known as **infinitesimal gradient ascent** or **IGA**. One can show that (α, β) does not always converge (depending on the values in \mathcal{R}_i and \mathcal{R}_j), but if it does, the resulting converged joint policy is a NE [SKM00]. However, there is a method called **Win or Learn Fast** or **WoLF** from [BV02] which can ensure that IGA policies always converge to a NE for two-agent two-action normal-form games. The trick is to learn slow (by using smaller κ) when winning (i.e., if $U_i(\alpha_k, \beta_k) > U_i(\alpha_e, \beta_k)$ where α_e is a policy from some NE), and to learn fast (by using larger κ) when losing (i.e., when not winning). This approach can be extended to stochastic games, without requiring knowledge of reward functions or policies. The resulting method is called **WoLF-PHC**, which stands WoLF with Policy Hill Climbing [BV02]. See Figure 5.5 for an example.

5.3.3 Centralized training of decentralized policies (CTDE)

We can improve performance beyond independent learning by using a paradigm known as Centralized Training and Decentralized Execution (**CTDE**), in which the learning algorithm has access to all the information (from all agents) at training time, but at test time, agents only observe their own local observations. The **central information** can contain the joint action taken by all agents, and/or the joint observation vector, even such joint information is not available at execution time.

We can modify the multi-agent A2C algorithm in Algorithm 19 to exploit this assumption by writing the i 'th value/ advantage function as $V(h_t^i, c_t; \mathbf{w}_i) / \text{Adv}(h_t^i, c_t, a_t^i)$, where c_t is the shared central information. It is perfectly valid for the *critics* to have this kind of central information, as long as the *policies* do not rely on this, since the critics are not used during execution. This is known as **centralized critics** with **decentralized actors**. We can create a CTDE version of PPO in a similar way; this is known as **MAPPO** [Yu+22].

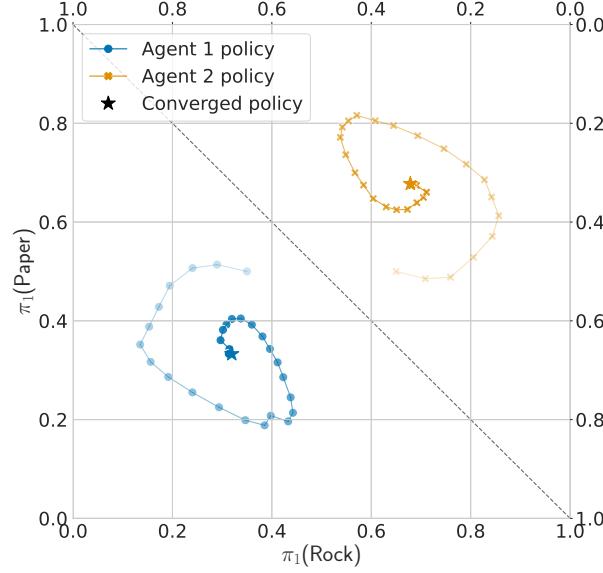


Figure 5.5: Learning dynamics of two policies for the rock-paper-scissors game. The upper and lower triangles illustrate the policies for each agent over time as a point in the 2d simplex (noting that $\pi_e^i(S) = 1 - (\pi_e^i(P) + \pi_e^i(R))$, where S is the scissors action, P is paper action, and R is rock action). The update rule is $\pi_{e+1} = LR(\mathcal{D}_e, \pi_e)$, where LR is the learning rule known as WoLF-PHC (see Section 5.3.2.4). From Fig 5.5 of [ACS24]. Used with kind permission of Stefano Albrecht.

5.3.3.1 Application to Diplomacy (Cicero)

In this section, we describe the **Cicero** system from [Met+22], which achieved human-level performance in the complex natural language 7-player strategy game called **Diplomacy**, which requires both cooperative and competitive behavior.³ Cicero used CTDE, combining an LLM for generating and interpreting dialog with a mix of self-play RL, imitation learning, opponent modeling, and policy generation using regret minimization. The system uses imitation learning on human games to warm-start the initial policy and language model, and then is refined using RL with self-play. The system uses explicit belief state modeling over the opponents' intents and plans; this is learning via supervised learning over past dialogues and game outcomes, and refined during self-play. (For a more recent method, which does not need human data, and which is based on LLMs augmented with memory and self-reflection, see [Gua+24].)

5.3.4 Value decomposition methods for common-reward games

In this section, we discuss methods for deriving a policy from a centralized state-action value function $Q(\mathbf{h}, c, \mathbf{a})$, where c is the central information (see Section 5.3.3), \mathbf{h} is the shared state (history), and \mathbf{a} is the joint action. To ensure that the per-agent policy can be implemented using only locally available information, we need to use **value decomposition** methods, which assume that the global value function can be decomposed into separate value functions, one per agent. (This is only possible if we are solving a common-reward or cooperative game.) This decomposition is valid provided the separate value functions satisfy a property known as the **individual global max** or **IGM** property, which says that

$$\forall \mathbf{a}. \mathbf{a} \in A^*(\mathbf{h}, c; \boldsymbol{\theta}) \Leftrightarrow \forall i. a_i \in A_i^*(h^i; \boldsymbol{\theta}_i) \quad (5.45)$$

where $A^*(\mathbf{h}, c; \boldsymbol{\theta}) = \operatorname{argmax}_{\mathbf{a}} Q(\mathbf{h}, c, \mathbf{a}; \boldsymbol{\theta})$ and $A_i^*(h^i; \boldsymbol{\theta}_i) = \operatorname{argmax}_{a_i} Q(h^i, a_i; \boldsymbol{\theta}_i)$. This ensures that picking actions locally for each agent will also be optimal globally.

³A simpler version of the game, called **No-press diplomacy**, without the natural language negotiation phase, is discussed in [Bak+21].

5.3.4.1 Value decomposition network (VDN)

For example, consider the **value decomposition network** or **VDN** method of [Sun+17]. This assumes a linear decomposition

$$Q(\mathbf{h}_t, c_t, \mathbf{a}_t; \boldsymbol{\theta}) = \sum_i Q(h_t^i, a_t^i; \boldsymbol{\theta}_i) \quad (5.46)$$

This clearly satisfies IGM.

5.3.4.2 QMIX

A more general method, known as **QMIX**, is presented in [Ras+18]. This assumes

$$Q(\mathbf{h}_t, c_t, \mathbf{a}_t; \boldsymbol{\theta}) = f_{\text{mix}}(Q(h_t^1, a_t^1; \boldsymbol{\theta}_1), \dots, Q(h_t^n, a_t^n; \boldsymbol{\theta}_n)) \quad (5.47)$$

where f_{mix} is a neural network that is constructed so that it is monotonically increasing in each of its arguments. (This is ensured by requiring all the weights of the mixing network to be non-negative; the weights themselves are predicted by another “hyper-network”, conditioned on the state h_t^i .) This satisfies IGM since

$$\max_{\mathbf{a}} Q(\mathbf{h}_{t+1}, c_{t+1}, \mathbf{a}; \bar{\boldsymbol{\theta}}) = f_{\text{mix}} \left(\max_{a^1} Q(h_{t+1}^1, a^1; \bar{\boldsymbol{\theta}}_1), \dots, \max_{a^n} Q(h_{t+1}^n, a^n; \bar{\boldsymbol{\theta}}_n) \right) \quad (5.48)$$

Hence we can fit this Q function by minimizing the TD loss (with target network $\bar{\boldsymbol{\theta}}$):

$$\mathcal{L}(\boldsymbol{\theta}) = \frac{1}{B} \sum_{(\mathbf{h}_t, c_t, \mathbf{a}_t, r_t, \mathbf{h}_{t+1}, c_{t+1}) \in \mathcal{B}} \left(r_t + \gamma \max_{\mathbf{a}} Q(\mathbf{h}_{t+1}, c_{t+1}, \mathbf{a}; \bar{\boldsymbol{\theta}}) - Q(\mathbf{h}_t, c_t, \mathbf{a}; \boldsymbol{\theta}) \right)^2 \quad (5.49)$$

5.3.5 Policy learning with self-play

For symmetric zero-sum games, where $r_i(s) = -r_j(s)$, we can assume that each player uses the same policy, modulo rearrangement of the input state. That is, $\pi_j(\cdot|s) = \pi_i(\cdot|\psi(s))$, where π_i is the main policy for player i , and $\psi(s)$ deterministically modifies the state to reflect the symmetry.

For example, consider the game of chess, where the state is represented by $s = (x, y)$, where x is a vector containing the location of player 1’s pieces (or -1 if they are removed), and y is a vector containing the opponent’s pieces. Thus the policy for player 1, π_1 , just needs to access the x part of the state. For player 2, we can transform the state vector into $s' = \psi(s) = (y, x)$ and then apply π_1 to choose actions, so $\pi_2(s) = \pi_1(\psi(s))$.

The result is a single agent problem, in which, from i ’s perspective, the state transition function has the following form, in which we treat player j as part of the environment:

$$p^\pi(s'|s, a^i) = \sum_{a^j} \pi(a^j|\psi(s)) p(s'|s, a^i, a^j) \quad (5.50)$$

Thus let’s learn π using standard single-agent policy learning methods, such as PPO or policy improvement based on search (decision-time planning). This is known as **self-play**.

Self-play is used by AlphaZero (see Section 4.2.2.1) to learn to play perfect information games like Chess and Go at super-human level. In [Lin+23], they used multi-agent self-play to train soccer agents. And in [Zhe+24b; Liu+25a; Bo+25] they used self-play to train LLM policies (see Section 6.2.5.4).

For perfect information, zero-sum games, self-play can be proved to converge to a Nash equilibrium. Unfortunately, in general games (e.g., for imperfect information games, such as Poker or Hanabi, or for general sum games), self-play can lead to oscillating strategies or cyclical behavior, rather than converging to a Nash equilibrium, as we discuss in Section 5.3.2.4. Thus self-play can result in policies that are easily exploited. We discuss more stable learning methods below.

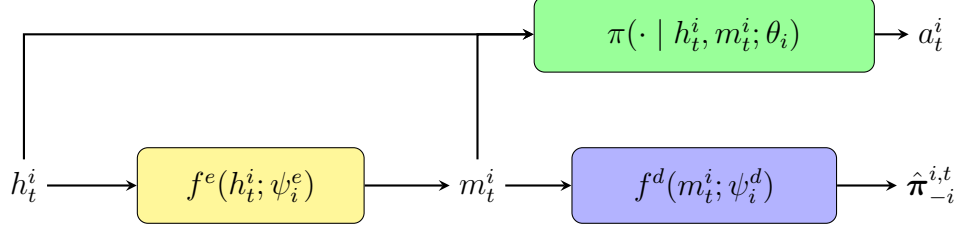


Figure 5.6: Encoder-decoder architecture for agent modeling. From Fig 9.21 of [ACS24]. Used with kind permission of Stefano Albrecht.

5.3.6 Policy learning with learned opponent models

Instead of using self-play, we can learn an **opponent model**. In the CTDE paradigm, where each agent sees the other agents actions, agent i can use supervised learning to predict the actions of agent j given i 's observations. There are many possible opponent models we can use (see e.g. [AS18] for a review). For example, we can train an encoder-decoder network to predict the actions of other agents via a bottleneck, and then pass this bottleneck embedding to the policy as side information, as proposed in [PCA21]. In more detail, let $m_t^i = f^e(h_t^i; \psi_i^e)$ be the encoding of i 's history, which is then passed to the decoder f^d to predict the other agents actions using $\hat{\pi}_{-i}^{i,t} = f^d(m_t^i; \psi_i^d)$. In addition, we pass m_t^i to i 's policy to compute the action using $a_t^i \sim \pi(\cdot | h_t^i, m_t^i; \theta_i)$. This is illustrated in Figure 5.6. (A similar method could be used to predict other properties of agent j , as long as their as observable by agent i .)

5.3.7 Best response

In this section, we discuss MARL algorithms that can provably converge to a Nash equilibrium, even for zero-sum, imperfect-information games, unlike basic policy learning methods based on self-play (or opponent modeling). The approach we use is built on the concept of a best response. This is the action (for a given state) that gives the highest expected reward for agent i , given that the policies for all other agents are fixed.

Specifically, let h^i be the information state for agent i (i.e., action-observation history). We compute its expected state-action value, given the joint policy, as follows:

$$AV_i^\pi(h^i, a^i) = \sum_{\mathbf{a}^{-i}} Q_i(h^i, (a^i, \mathbf{a}^{-i})) \prod_{j \neq i} \pi_j(a^j | h^j; \theta_j^i) \quad (5.51)$$

The best response is then given by

$$BR_i(h^i) = \operatorname{argmax}_{a^i} AV_i^\pi(h^i, a^i) \quad (5.52)$$

If there are a large number of actions, we can approximate the sum over \mathbf{a}^{-i} using Monte Carlo sampling. The only thing left is to specify how to learn the opponent policies. We discuss this below.

5.3.7.1 Fictitious play

In **fictitious play**, each agent i estimates the policies of the other players, based on their past actions. It then computes a best response. For example, imagine you're playing rock-paper-scissors repeatedly. If you notice your opponent plays "rock" 60% of the time so far, your best response is to play "paper" more often. You adjust your strategy based on the empirical frequency of their past moves. The method is called "fictitious" because each player is acting as if the opponents are playing a fixed strategy, even though they're actually adapting.

The method was originally developed for non-repeated normal-form games (which are stateless, so $h^i = []$ and $Q_i(\mathbf{a}) = \mathcal{R}_i(\mathbf{a})$). In this case, we can estimate the policies by counting and averaging. That is,

$$\hat{\pi}_j^t(a^j) = \frac{C_j^t(a^j)}{\sum_{a'} C_j^t(a')} \quad (5.53)$$

where $C_j^t(a^j)$ is the number of times agent j chose action a^j in episodes up to step t . We then compute the best response, given by

$$\text{BR}_i^t = \underset{\mathbf{a}^i}{\text{argmax}} \sum_{\mathbf{a}^{-i}} \mathcal{R}_i((a^i, \mathbf{a}^{-i}) \prod_{j \neq i} \hat{\pi}_j^t(a^j) \quad (5.54)$$

Equivalently we can say that π_i^t is the best response to $\bar{\pi}_{t-1} = \text{avg}(\pi_1, \dots, \pi_{t-1})$. For two-player zero-sum finite games, this procedure will converge to a NE. That is, the exploitability of the average $\bar{\pi}_t$ generated by FP converges to zero as t grows large.

5.3.7.2 Neural fictitious self play (NFSP)

We can extend FP to the partially observed, non-tabular setting as follows. If we assume each agent sees the other agents actions, then it is easy to learn an opponent model representing their average strategy. In particular, agent i will learn a model of j 's policy, given i 's state (history), which we denote by $\bar{\pi}_{j|i}^t(a^j|h^i)$. We fit this by minimizing the cross entropy loss

$$\mathcal{L}(\bar{\pi}_{j|i}^t) = \mathbb{E}_{k \sim U(1,t), (h_k^i, a_k^j) \in \mathcal{D}_t} \left[-\log \bar{\pi}_{j|i}^t(a_k^j|h_k^i) \right] \quad (5.55)$$

where \mathcal{D}_t is the replay buffer containing previous states and actions of all the players. In addition, we use DQN to learn $Q_i(h^i, \mathbf{a})$ for each agent. We then use this learned average policy, plus the Q functions, to compute AV_i , and hence the best response.

In the zero-sum two-player case, we can use self-play, so we just assume $\bar{\pi}_{-i}$ is the opposite of $\bar{\pi}_i$, which can be learned by supervised learning applied just to its own states and actions. This is called **fictitious self play** [HLS15]. It was extended to the neural net case in [HS16], who call it **neural fictitious self play**. If the Q function converges to the optimal function, then this process converges to a NE, like standard FP.

5.3.8 Population-based training

In Section 5.3.5, we discussed the concept of **self-play**, which is a way to train an agent to play a two-player game by modeling the opponent as using the same policy as the agent itself. To avoid overfitting, we typically train against multiple versions of the agent's own policy. This concept can be generalized to work with general-sum games with two or more players, by training against a population of different policies. This is called **population based training** [Jad+19].

5.3.8.1 PSRO (policy space response oracle)

In this section, we describe the **policy space response oracle** or **PSRO** method of [Lan+17], which is a game-theoretic instance of population based training, which can compute policies that satisfy various solution concepts for any kind of stochastic game, including partially observed, general sum games.

The idea behind PSRO is as follows. At generation k , each agent i has a finite set of policies it can use, denoted Π_i^k . We can define a normal-form **meta-game** M^k from this by letting each agent choose one of its policies, where the reward for the joint action $\mathbf{a} = \boldsymbol{\pi} = (\pi_1, \dots, \pi_n)$ is given by $\mathcal{R}_i(\boldsymbol{\pi}) = U_i(\boldsymbol{\pi})$ for $i = 1 : n$. These returns can be estimated empirically by simulating n agents interacting with each other according to these policies in the underlying game G . Once we have determined the reward matrix, we can solve for some kind of equilibrium solution (e.g., Nash equilibrium), using a **meta-strategy solver**. We can then extract the probability distributions over policies (aka strategy) σ_i^k for each agent. To ensure this distribution is diverse, we can enforce a lower bound that $\sigma_i^k(\pi_i) > \epsilon$, so each policy is used some fraction of the time.

k	Π_1^k	Π_2^k	σ_1^k	σ_2^k	π'_1	π'_2
1	<u>R</u>	<u>P</u>	1	1	S	P
2	R, <u>S</u>	P	(0, 1)	1	S	R
3	R,S	<u>R</u> ,P	$(\frac{2}{3}, \frac{1}{3})$	$(\frac{2}{3}, \frac{1}{3})$	P	R/P
4	R, <u>P</u> ,S	R,P	$(0, \frac{2}{3}, \frac{1}{3})$	$(\frac{1}{3}, \frac{2}{3})$	R	S
5	R,P,S	R,P, <u>S</u>	$(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$	$(\frac{1}{3}, \frac{1}{3}, \frac{1}{3})$	R/P/S	R/P/S

Figure 5.7: PSRO for rock-paper-scissors. We show the populations Π_i^k , distributions σ_i^k and best responses π'_i for both agents over generations $k = 1 : 5$. (We use the shorthand of R to denote the pure policy that always plays R , and similarly for S and P .) New entries added to the population are shown with an underline. At generation $k = 3$, there are 2 best responses, R and P , so the oracle can select either of them. Similarly, at generation $k = 5$, there are 3 best responses.

We can now expand the set of policies for each agent by using an oracle to compute a new policy π'_i and adding it to Π_i^k to create the set Π_i^{k+1} . For example, the oracle can compute a best response

$$\pi'_i \in \operatorname{argmax}_{\pi_i} \mathbb{E}_{\pi_{-i} \sim \sigma_{-i}^k} [U_i((\pi_i, \pi_{-i}))] \quad (5.56)$$

where $\sigma_{-i}^k(\pi_{-i}) = \prod_{j \neq i} \sigma_j^k(\pi_j)$. We can compute π'_i by using a single agent RL algorithm in the underlying game G . Since the policies of the other agents are uncertain, as well as the environment, we can use Bayesian RL methods, such as [OA14]. In this approach, in each episode the policies of the other agents $j \neq i$ are sampled from $\pi_j \sim \sigma_j^k$, and then standard RL is applied.

It can be shown that, if PSRO uses a meta-solver that computes exact Nash equilibria for the meta-game, and if the oracle computes the exact best-response policies in the underlying game G , then the distributions $\{\sigma_i^k\}_{i \in \mathcal{I}}$ converge to a Nash equilibrium of G . See Figure 5.7 for an example.

Note that PSRO can also be applied to general-sum, imperfect information games. For example, [Li+23b] uses (information set) MCTS, together with a learned world model, to compute the best response policy π'_i at each step of PSRO.

5.3.8.2 Application to StarCraft (AlphaStar)

The **AlphaStar** system of [Vin+19] used a PSRO-like method, combined with the (single agent) A2C RL algorithm, to achieve grandmaster status in the challenging real-time strategy game known as StarCraft II.⁴ In particular, it used the following steps: Build a pool of agents that represent different playstyles and skill levels (known as a league); Compute best responses to existing strategies; Update a meta-strategy to mix agents in a way that approximates a Nash equilibrium; select opponents from the Nash mixture to ensure robustness; and train a new agent against the weighted mixture of past opponents. See the paper for more details.

5.3.9 Counterfactual Regret Minimization (CFR)

In this section we describe **Counterfactual Regret Minimization (CFR)**, which is an algorithm for imperfect information, two-player, zero-sum game. In [Zin+07] they show that when using this procedure,

⁴In StarCraft II, the AI agent controls an entire army, and must defeat a human opponent, making this a zero-sum, two-player game, which can be solved using deep RL with self-play. This is different from the StarCraft Multi-Agent Challenge (**SMAC**) [Sam+19], which is a cooperative, partially observed multi-agent game, where the agents (corresponding to individual units) must work together as a team to defeat a fixed AI opponent in certain predefined battles.

the average policies converge to an epsilon-Nash.

5.3.9.1 Tabular case

Let $\tau = (s_0, \dots, s_t)$ be a trajectory of world states. Let $\eta^\pi(\tau)$ be the probability of this trajectory under the joint policy. (Note that stochastic dynamics of the world are modeled by the policy of the chance player.) We can decompose this as

$$\eta^\pi(\tau) = \eta_i^\pi(\tau) \eta_{-i}^\pi(\tau) \quad (5.57)$$

Similarly define $\eta^\pi(\tau, \mathbf{z})$ as the probability of the trajectory $\tau = (s_0, \dots, s_t)$ followed by $\mathbf{z} = (s_{t+1}, \dots, s_T)$, which is some continuation that ends in a terminal state. Let $Z(\mathbf{h}^i) = \{(\tau, \mathbf{z})\}$ be the set of trajectories and their terminal extensions which are compatible with \mathbf{h}^i , in the sense that $\mathcal{O}_i(\tau) = \mathbf{h}^i$, and which end in a terminal state. Also, let τ_{a^i} be the trajectory followed by action a^i . We then define the **counterfactual state-action value** for an information state as

$$q_{\pi,i}^c(\mathbf{h}^i, a^i) = \sum_{(\tau, \mathbf{z}) \in Z(\mathbf{h}^i)} \eta_{-i}^\pi(\tau) \eta^\pi(\tau_{a^i}, \mathbf{z}) u_i(\mathbf{z}) \quad (5.58)$$

The counterfactual state-value is

$$v_{\pi,i}^c(\mathbf{h}^i) = \sum_{a^i} \pi_i(a^i | \mathbf{h}^i) q_{\pi,i}^c(\mathbf{h}^i, a^i) \quad (5.59)$$

Finally, define the instantaneous **counterfactual regret** for player i at iteration k to be

$$r_i^k(\mathbf{h}^i, a^i) = q_{\pi^k,i}^c(\mathbf{h}^i, a^i) - v_{\pi^k,i}^c(\mathbf{h}^i) \quad (5.60)$$

Note that this is the counterfactual version of an advantage function, as explained in [Sri+18]. Similarly we define the cumulative counterfactual regret to be

$$R_i^k(\mathbf{h}^i, a^i) = \sum_{j=0}^k r_i^j(\mathbf{h}^i, a^i) \quad (5.61)$$

CFR starts with a uniform random joint policy π^0 and then updates it at each iteration by performing **regret matching** [HMC00; Far+23b]. That is, it updates the policy as follows

$$\pi_i^{k+1}(\mathbf{h}^i, a^i) = \begin{cases} \frac{R_i^{k,+}(\mathbf{h}^i, a^i)}{\sum_{a \in \mathcal{A}_i(\mathbf{h}^i)} R_i^{k,+}(\mathbf{h}^i, a)} & \text{if denominator is positive} \\ \frac{1}{|\mathcal{A}_i(\mathbf{h}^i)|} & \text{otherwise} \end{cases} \quad (5.62)$$

where $x^+ = \max(x, 0)$.

In [Zin+07] they show that the above procedure results in an ϵ -Nash equilibrium, where $\epsilon = O(\max_i |\mathcal{H}_i| \sqrt{|\mathcal{A}_i|} / \sqrt{t})$, for any two-player, zero-sum game (with perfect recall of past observations).

5.3.9.2 Deep CFR

In practice, the expectations over trajectories in Equation (5.58) can be approximated using Monte Carlo sampling [Lan+09]. In addition, we can approximate the tabular q , v and r terms with neural networks; this is called **Deep CFR** [Bro+19], which builds on the earlier Regression CFR method of [Wau+15].

5.3.9.3 Applications to Poker and other games

The first known combination of CFR with neural networks was DeepStack [Mor+17]. This was also was one of the first systems to beat professional players at a two-player poker variant called **heads-up no-limit Texas hold'em**. Another system that came out at the same time, and also beat humans at this game, was the

(neural-free) **Libratus** method of [BS17], based on regret matching. Libratus was later extended to make the **Pluribus** method of [BS19], which was able to beat human players at the six-player version of Texas hold'em.

In [Sch+21a], they proposed a method called **Student of Games**, that is a version of AlphaZero where CFR is the policy improvement operator. This was applied to various games, such as Chess, Go, Poker, and Scotland Yard.

5.3.10 Regularized policy gradient methods

In this section, we discuss policy gradient methods that incorporate a regularization term to ensure convergence, even in adversarial settings, such as 2p0s games.

5.3.10.1 Magnetic Mirror Descent (MMD)

In [Sok+22], they present the **Magnetic Mirror Descent** or **MMD** algorithm, which is designed for two-player zero-sum games (but which can also work well for single player games). MMD is a modification of policy gradient that adds additional regularizers to ensure it converges (unlike traditional PG methods, which can oscillate). In the tabular case, we use an update of the following form, applied at each decision point (state) s and for each agent i separately:

$$\pi_{k+1} = \operatorname{argmax}_{\pi} \langle \pi, q_k \rangle - \alpha D_{\text{KL}}(\pi, \rho) - \frac{1}{\eta} D_{\text{KL}}(\pi, \pi_k) \quad (5.63)$$

where $q_k(a) = q_k(s, A)$ is the value of action a in state s , $\pi(A) = \pi(A|s)$ is the agent policy, $\langle \pi, q_k \rangle = \mathbb{E}_{a \sim \pi} q_k(a)$ is an expectation, ρ is a magnet policy (designed to prevent oscillation), α is a regularization term (corresponding to entropy penalty if ρ is uniform), and η is a stepsize. For discrete actions, the optimal solution to the above is given by the following (computed elementwise)

$$\pi_{k+1} \propto [\pi_k \rho^{\alpha \eta} e^{\eta q_k}]^{\frac{1}{1+\alpha \eta}} \quad (5.64)$$

If we drop the magnet term, by setting $\alpha = 0$, the method is equivalent to the **mirror descent policy optimization** or **MDPO** algorithm of [Tom+20]. In this case, the optimal solution is given by

$$\pi_{k+1} \propto [\pi_k e^{\eta q_k}] \quad (5.65)$$

as in the exponentiated gradient algorithm.

In [Sok+22], they prove that this procedure (when used with a uniform magnet policy and applied to NFGs) will converge to a QRE (Section 5.2.6) exponentially fast. If the entropy term is annealed to 0, they can match the results of CFR (Section 5.3.9) in the case of tabular games. Their theory does not yet apply to the parametric case, but experimentally they still find fast convergence to the AQRE.

5.3.10.2 PPO

The MMD method of Section 5.3.10.1 is very similar to the PPO algorithm of Section 3.3.3. In particular, the KL penalized version of PPO uses the following loss

$$\mathbb{E}_{s_t, a_t} \left[\frac{\pi(a_t|s_t)}{\pi_{\text{old}}(a_t|s_t)} A_{\text{old}}(s_t, a_t) + \alpha \mathbb{H}(\pi(\cdot|s_t)) - \beta D_{\text{KL}}(\pi_{\text{old}}(\cdot|s_t), \pi(\cdot|s_t)) \right] \quad (5.66)$$

where $A_{\text{old}}(s, a) = q_{\text{old}}(s, a) - v_{\text{old}}(s)$ is the advantage function. By comparison, if we use a uniform magnet for ρ , the MMD loss in Equation (5.63) becomes

$$\mathbb{E}_{s_t, a_t} [\pi(a_t|s_t) q_{\text{old}}(s_t, a_t) + \alpha \mathbb{H}(\pi(\cdot|s_t)) - \beta D_{\text{KL}}(\pi(\cdot|s_t), \pi_{\text{old}}(\cdot|s_t))] \quad (5.67)$$

where $\beta = 1/\eta$ is the inverse stepsize. The main difference between these equations is just the use of a reverse KL instead of forwards KL. (The two expressions also differ by the scaling factor $1/\pi_{\text{old}}(a_t|s_t)$ and the offset term $v_{\text{old}}(s_t)$.)

Despite the similarities, in [Sok+22], PPO has been shown to perform worse than MMD on various 2p0s games. One possible reason for PPO’s poor performance is due to the use forwards vs reverse KL penalty. However, [HMDH20] compared the use of reverse KL regularization instead of forward KL in PPO for Mujoco, and found that the two yielded similar performance. The explanation suggested in [Rud+25] is simply that the hyper-parameters in PPO (in particular, the entropy penalty α) was not tuned properly for the 2p0s setting (the latter tending to require much larger values, such as 0.05-2.0, whereas single agent PPO implementations usually use 0-0.01).

They experimentally tested this hypothesis by comparing PPO with various other algorithms (including MMD, CFR (Section 5.3.9), PSRO (Section 5.3.8.1) and NFSP (Section 5.3.7.2)) on a set of imperfect information games (partially observed or “phantom”/“dark” versions of Tic-Tac-Toe and 3x3 Hex, where the agents actions are invisible to the non-acting player). They find that properly tuned policy gradient methods (including both PPO and MMD) performed the best in terms of having the lowest exploitability scores. (The exploitability score is defined in Equation (5.21), and was computed by exactly solving for the optimal opponent policy given a candidate learned policy.)

The above experimental result led the authors of [Rud+25] to propose the following “Policy Gradient Hypothesis”:

Appropriately tuned policy gradient methods that share an ethos with magnetic mirror descent are competitive with or superior to model-free deep reinforcement learning approaches based on fictitious play, double oracle [population-based training], or counterfactual regret minimization in two-player zero-sum imperfect-information games.

If true, this hypothesis would be very useful, since it means we can use standard single agent policy gradient methods, such as (suitably tuned) PPO, for multiplayer games, both cooperative (see [Yu+22]) and adversarial (see [Rud+25]).

5.3.11 Decision-time planning methods

In this section we focus on decision-time planning (DTP) methods, that improve upon a base policy (known as a blueprint policy) by doing some kind of forward search (from the current state) in a world model, as discussed in Section 4.2. We focus the **update-equivalent DTP** method of [Sok+23], which makes a connection between DTP and other policy update algorithms.

Recall from Equation (2.12) that the policy iteration algorithm can be viewed as performing an update to the policy at each step based on acting greedily wrt $Q(s, a)$:

$$\pi_{\text{new}}(s) = \operatorname{argmax}_a R(s, a) + \gamma \mathbb{E}[V_\pi(s')] = \operatorname{argmax}_a Q(s, a) \quad (5.68)$$

If we consider a single state, we can write this update as

$$\pi_{\text{new}} = U(\pi, q) = \operatorname{argmax}_{\pi' \in \Delta(A)} \langle \pi', q \rangle \quad (5.69)$$

One way to estimate the action values q for the current state is to perform **Monte Carlo search** or **MCS** [TG96], which unrolls possible futures using the current policy, as in DTP. Thus with enough samples, DTP (with the correct world model, and a suitable exploratory policy) combined with this update will give the same results as (asynchronous) policy iteration.

5.3.11.1 Magnetic Mirror Descent Search (MMDS)

In [Sok+23] they propose to generalize this idea to the multi-agent setting by using the MMD algorithm from Section 5.3.10.1 as the update operator. They call this **magnetic mirror descent search** or **MMDS**. The local policy update (for player i) has the form

$$\pi_{\text{new}} = U(\pi, q) = \operatorname{argmax}_{\pi' \in \Delta(A)} \langle \pi', q \rangle - \alpha D_{\text{KL}}(\pi', \rho) - \frac{1}{\eta} D_{\text{KL}}(\pi', \pi) \quad (5.70)$$

Algorithm 20: Magnetic Mirror Descent Search (MMDS)

```
1 Input: current state  $h_t^i$ , joint policy  $\pi$ , agent id  $i$ 
2  $q[a] = 0$ ,  $N[a] = 0$  for each action  $a \in \mathcal{A}_i$ 
3 repeat
4   Sample current world state using agent's local belief state:  $s_t \sim P_\pi(\cdot|h_t^i)$ 
5   for  $a \in \mathcal{A}_i$  do
6     Sample return  $G_{\geq t} \sim P_\pi(G_{\geq t}|s_t, a)$  by rolling out  $\pi$  in world model starting at  $s_t$ 
7      $q[a] = q[a] + G_{\geq t}$ 
8      $N[a] = N[a] + 1$ 
9    $q[a] = q[a]/N[a]$  for  $a \in \mathcal{A}_i$ 
10  Return  $U(\pi^i(h_t^i), q)$  by performing SGD on Equation (5.71).
11 until until search budget exhausted;
```

where π is the previous local (**blueprint**) policy and ρ is the local magnet policy (which can be taken as uniform). If h_t^i is the current state (root of search tree for player i), and the actions are discrete, we can equivalently perform an SGD step on the following parametric policy loss:

$$\mathcal{L}(\theta) = \sum_a \left[\pi_\theta(a|h_t^i) q(h_t^i, a) - \alpha \pi_\theta(a|h_t^i) \log \frac{\pi_\theta(a|h_t^i)}{\rho(a)} - \frac{1}{\eta} \pi_\theta(a|h_t^i) \log \frac{\pi_\theta(a|h_t^i)}{\pi_{\text{old}}(a|h_t^i)} \right] \quad (5.71)$$

See Algorithm 20 for the pseudocode.

Note that, if we use a uniform magnet, this is equivalent to adding an entropy regularizer. Also, for common-payoff games, we can drop the magnet term, which gives rise to the simpler **mirror descent search** method.

5.3.11.2 Belief state approximations

To implement this algorithm, we need to sample from $P_\pi(s_t|h_t^i)$, which is the distribution over world states given agent i 's local history. One approach to this is to use particle filtering, cf. [Lim+23]).

Another approach is to train a **belief model** to predict the other player's private information, and the underlying environment state, given the current player's history, i.e., we learn to predict $P(s_t, \{h_t^j\}|h_t^i)$. In the **learned belief search (LBS)** method of [Hu+21] (which was designed for Hanabi, which is a Dec-POMDP), rather than predicting the entire action-observation history for each agent, they just predict the private information (card hand) for each agent (represented as a sequence of tokens). This can be used (together with the shared public information) to reconstruct the environment state. They train this model (represented as a seq2seq LSTM) using supervised learning, where agent i learns to predict its own private information given its public history. At test time, agent i uses j 's public history as input to its model to sample j 's private information. (This assumes that j is using the same blueprint policy to choose actions that i used during training.) Given the imputed private information, it then reconstructs the environment state and performs rollouts, using the joint blueprint policy, in order to locally improve its own policy.

5.3.11.3 Experiments

In [Sok+23], they implemented the above method and applied it to several imperfect information games (using the true known world model) For the common-reward game of Hanbai (5 card and 7 card variants), they used PPO to pretrain the blueprint policy, and they pretrained a seq2seq belief model. At run time, they use 10k samples for each step of MDS to locally improve the policy (which takes about 2 seconds). They observed modest gains over rival methods. For the 2p0s games, they used the partially observed (dark/phantom) versions of 3x3 Hex and Tic-Tac-Toe. For belief state estimation, they use a particle filter with just 10 particles, for speed. As a blueprint policy they consider uniform and MMD (for 1M steps). They find that MMDS can improve the blueprint, and this combination beats baselines such as PPO and NFSP.

They also compare to MMD as a baseline. For the MMD-1M baseline, the blueprint matches the baseline (by construction), but the MMDS version beats it. However, the MMD-10M baseline beats MMDS, showing that enough offline computation can beat less online computation.

5.3.11.4 Open questions

It is an interesting open question how well this MMDS method will work when the world model needs to be learned, since this results in rollout errors, as discussed in Section 4.3.1. Similarly errors in the belief state approximation may adversely affect the estimate of q for the root node.

In addition, it is an open question to prove convergence properties of the generalized version of MMDS, that uses more than just action value feedback. For example, MCTS updates the local policy at internal nodes, not just the root node. In some cases, MCTS can work better than simple MCS, although this is not always the case (see e.g., [Ham+21]).

5.3.12 MARL for LLM agents

A recently growing trend is to use LLMs as agents (see Chapter 6), which can be made to interact with each other, via protocols such as **A2A** (agent-to-agent).⁵ It is possible to apply MARL techniques to optimize such systems, see e.g., [Lia+25; Zho+25c; Ren+25]. See also Section 6.2.5.4 for a discussion of how to use self-play to train LLMs, and [Sun+25] for a broader review of game theory and LLMs.

⁵See <https://developers.googleblog.com/en/a2a-a-new-era-of-agent-interoperability/>.

Chapter 6

LLMs and RL

6.1 Introduction

In this section, we discuss connections between RL and **foundation models**, also called **large language models** or **LLMs**. LLMs are generative models (usually based on auto-regressive transformers) which are trained on large amounts of web data.¹ More details on the connections between RL and LLMs can be found in e.g., [Pte+24; Hao+25; Zha+25a; Kha+25b], and more details on LLMs in general can be found in e.g., [XZ25; Bur25; Bro24].

6.2 RL for LLMs

In this section, we discuss how to use RL to improve the performance of LLMs. This is a fast growing field, so we only briefly mention a few highlights. For more details, see e.g. [Lam25; Bro24] and other references online.

6.2.1 RL fine tuning (RLFT)

LLMs are usually trained with behavior cloning, i.e., MLE on a fixed dataset, such as a large text corpus scraped from the web. This is called **pre-training**. We can then improve their performance using various **post-training** methods, which are designed to improve their capabilities and **alignment** with human preferences (see e.g., [Zen+25]), as opposed to just being generative models of the data seen on the web. A simple way to perform post-training is to use **instruction fine tuning**, also called **supervised fine-tuning** (or **SFT**), in which we collect human demonstrations of (prompt, response) pairs, and fine-tune the model on them. However, it is very difficult to collect sufficient quantities of such data. An alternative to demonstrating good behaviors is to use RL to train the model using a suitable reward function. (We discuss where these reward functions come from in Section 6.2.2.) This is called **reinforcement learning fine-tuning** or **RLFT**.

RLFT can be preferable to SFT for several reasons. First, it is often the case that verification is easier than generation (e.g., it is easier to ask people which answer they prefer rather than to ask them to generate good answers, an insight we exploit in Section 6.2.2.3). Second, RL can be used to learn a set of “thinking actions”, which are created in response to the question before generating the answer (see Section 6.2.3). For complex problems (e.g., in math), this tends to work much better than trying to directly learn an input-output mapping [PLG23]. (It is possible to use SFT on explicitly provided thinking traces, but it has been found that RL can generalize more reliably [Chu+25].) Finally, RL opens the path to super-human performance [SS25], going beyond whatever supervised examples humans can create.

¹When also trained on visual data, LLMs are sometimes called **vision language models** or **VLMs**, and when trained on action data, they are called **vision language action** models or **VLAs** [Kaw+25].

6.2.2 Reward models

In this section, we discuss different kinds of reward functions that are used for RLFT.

6.2.2.1 RL with verifiable rewards (RLVR)

For problems such as math and coding, it can be easy to determine if an answer is correct, by checking equality between the generated answer and the true answer (for math), or checking if a set of unit tests pass (for code). This allows us to define a binary reward signal. Using RL with such a reward is called “**RL with verifiable rewards**” or **RLVR** (see e.g., [ZWG22; Lam+24]). We will use this approach to train “thinking models” in Section 6.2.3.

6.2.2.2 Process vs outcome reward models

If the reward function $R(s_t)$ is defined on partial trajectories, it is called a **process reward model** or **PRM**. This provides a form of dense feedback. If the reward is just defined on the final sequence $R(s_T) = R(s_0, \mathbf{a}_{1:T})$, it is called an **outcome reward model** or **ORM**, and corresponds to a sparse reward. For example, suppose we are solving a math problem using a thinking model (see Section 6.2.3): if we just check the final answer, we have an ORM, but if we also check correctness of the intermediate proof steps, we have a PRM. Note that a PRM is related to a value function (that models expected future reward), and is typically harder to learn than an ORM.

6.2.2.3 Learning the reward model from human feedback (RLHF)

To train LLMs to do well in general tasks, such as text summarization or poetry writing, it is common to use **reinforcement learning from human feedback** or **RLHF**, which refers to learning a reward model from human data, and then using RL to train the LLM to maximize this.

The basic idea is as follows. We first generate a large number of (context, answer1, answer2) tuples either by a human or an LLM. We then ask human raters if they prefer answer 1 or answer 2. Let x be the prompt (context), and y_w be the winning (preferred) output, and y_l be the losing output. Let $r_\theta(x, y)$ be the reward assigned to output y . (This model is typically a shallow MLP on top of the last layer of a pretrained LLM.) We train the reward model by maximizing the likelihood of the observed preference data. The likelihood function is given by the **Bradley Terry choice model**:

$$p_\theta(y_w > y_l) = \frac{\exp(r_\theta(x, y_w))}{\exp(r_\theta(x, y_w)) + \exp(r_\theta(x, y_l))} \quad (6.1)$$

We thus need to maximize

$$J(\theta) = \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\frac{\exp(r_\theta(x, y_w))}{\exp(r_\theta(x, y_w)) + \exp(r_\theta(x, y_l))} \right] = \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\frac{1}{1 + \frac{\exp(r_\theta(x, y_l))}{\exp(r_\theta(x, y_w))}} \right] \quad (6.2)$$

$$= \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\sigma(r_\theta(x, y_w) - r_\theta(x, y_l))] \quad (6.3)$$

Equivalently we can minimize

$$\mathcal{L}(\theta) = \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \left(1 + e^{r_\theta(x, y_l) - r_\theta(x, y_w)} \right) \right] \quad (6.4)$$

$$= -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log (\sigma(r_\theta(x, y_w) - r_\theta(x, y_l)))] \quad (6.5)$$

In some cases, we ask human raters if they prefer answer 1 or answer 2, or if there is a tie, denoted $y \in \{1, 2, \emptyset\}$. In this case, we can optimize

$$\mathcal{L}(\theta) = \mathbb{E}_{(x, y_1, y_2, y) \sim \mathcal{D}} [\mathbb{I}(y = 1) \log p_\theta(y_1 > y_2 | x) + \mathbb{I}(y = 2) \log p_\theta(y_1 < y_2 | x)] \quad (6.6)$$

$$+ \mathbb{I}(y = \emptyset) \log p_\theta(y_1 > y_2 | x) p_\theta(y_1 < y_2 | x)] \quad (6.7)$$

For a discussion of some of the implementation details of RLHF, see [Lam25]. For a discussion of other ways of modeling preferences and rewards, based on **social choice theory**, see [Con+24; Ge+24].

6.2.2.4 Learning the reward model from AI feedback (RLAIF)

Instead of asking humans their preferences for each possible input example, we can ask an LLM to predict the preference. This is called **LLM as judge**. We can then fit the reward model to this synthetically labeled data, just as in RLHF. Alternatively, we can just ask the LLM to predict the reward directly. This is called **RLAIF**, which stands for RL from AI feedback. It is also common to use VLMs for RLAIF (see e.g., [Du+23a; Bau+23; Hua+23; Son+23; Roc+24]).

In order to specify how to judge things, the LLM needs to be prompted. Anthropic (which is the company that makes the Claude LLM) created a technique called **constitutional AI** [Ant22], where the prompt is viewed as a “constitution”, which specifies what kinds of responses are desirable or undesirable. With this method, the system can critique its own outputs, and thus self-improve.

6.2.2.5 Generative reward models (GRM)

A **generative reward model** or **GRM** predicts the reward for a given response, but also returns its chain of thought, thus providing richer textual feedback. For example, it might generate “This response is helpful but makes a factual error about X, so overall reward is 0.3”. In addition to passing the scalar reward to an RL algorithm, the textual feedback can be parsed by the LLM itself to decide how to improve the policy. For example, the **GEPA** algorithm of [Agr+25] uses an evolutionary algorithm to optimize prompts for a frozen LLM, by mutating them given textual feedback from a GRM.

6.2.3 Agents which “think”

In this section, we discuss how to leverage the power of LLMs to create agents that “think” before they act.

6.2.3.1 Chain of thought prompting

The quality of the output from an LLM can be improved by prompting it to “show its work” before presenting the final answer. These intermediate tokens are called a “**Chain of Thought**” [Wei+22]. Models that act in this way are often said to be doing “**reasoning**” or “**thinking**” (see e.g., [Zha+25b; Xu+25]), although in less anthropomorphic terms, we can think of them as just policies with dynamically unrolled computational graphs [Kam+25]. This is motivated by various theoretical results that show that such CoT can significantly improve the expressive power of transformers [MS24; Li+24c].

6.2.3.2 Training a thinking model using RL

Rather than just relying on prompting, we can explicitly train a model to think by letting it generate a variable number of tokens “in its head” before generating the final answer. Only the final outcome is evaluated, using a known reward function (as in the case of math and coding problems).

This approach was recently demonstrated by the **DeepSeek-R1-Zero** system [Dee25] (released by a Chinese company in January 2025). They started with a strong LLM base model, known as **DeepSeek-V3-Base** [Dee24], which was pre-trained on a large variety of data (including Chains of Thought). They then used a variant of PPO, known as GRPO (see Section 6.2.4.3) to do RLFT, using a set math and coding benchmarks where the ground truth answer is known. The resulting system got excellent performance on math and coding benchmarks.² The closed-source models **ChatGPT-o1** and **ChatGPT-o3** from OpenAI³ and the **Gemini 2.0 Flash Thinking** model from Google Deepmind⁴ are believed to follow similar principles

²Although DeepSeek-R1-Zero exhibited excellent performance on math and coding benchmarks, it did not work as well on more general reasoning benchmarks. So their final system, called DeepSeek-R1, combined RL training with more traditional SFT (on synthetically generated CoTs).

³See <https://openai.com/index/learning-to-reason-with-llms/>.

⁴See <https://deepmind.google/technologies/gemini/flash-thinking>.

to DeepSeek-R1, although the details are not public.⁵ For a recent review of these reasoning methods, see e.g., [Xu+25; Zha+25c].

6.2.3.3 Thinking as marginal likelihood maximization

Since we usually only care about maximizing the probability that the final answer is correct, and not about the values or “correctness” of the intermediate thoughts (since it can be hard to judge heuristic arguments), we can view training a thinking model as equivalent to maximizing the marginal likelihood $p(y|x) = \sum_z p(y, z|x)$, where z are the latent thoughts (see e.g., [Hof+23; Zel+24; TWM25]).

6.2.3.4 Can we bootstrap a model to think from scratch?

One reason DeepSeek-R1 got so much attention in the press is that during the training process, it seemed to “spontaneously” exhibit some “emergent abilities”, such as generating increasingly long sequence of thoughts, and using self-reflection to refine its thinking, before generating the final answer.

Note that the claim that RL “caused” these emergent abilities has been disputed by many authors (see e.g., [Liu+25b; Yue+25]). Instead, the general consensus is that the base model itself was already trained on datasets that contained some COT-style reasoning patterns. This is consistent with the findings in [Gan+25a], which showed that applying RL to a base model that had not been pre-trained on reasoning patterns (such as self-reflection) did not result in a final model that could exhibit such behaviors. However, RL can “expose” or “amplify” such abilities in a base model if they are already present to a certain extent. (See also [FMR25] for a detailed theoretical study of this issue.)

Recently **Absolute Zero Reasoner** from [Zha+25f] showed it is possible to automatically generate a curriculum of (programs, inputs, outputs), which is used to improve the math and coding abilities of the LLM using RL. In particular, the LLM is trained to perform induction (predict p given (i, o) pairs), deduction (predict o given p and i), and abduction (predict i given p and o). (See also related work in [Kub+25].)

6.2.3.5 Agentic AI

There is currently a lot of hype around **Agentic AI** systems, that consist of a set of interacting LLMs, often called “agents”, which are essentially different prompts, reflecting different roles or personas, which can be given to the shared LLM to make it act in different ways. Typically these prompts, and the way the different agents interact, are hand-designed — this is called a **workflow** or **scaffolding**. Such agents may process the input, some may access or process memory (see e.g., [Zha+24c; Sum+24]), and some may call tools such as web search (see e.g., [Luo+25; Gul25]). Note, however, that unlike the true multi-agent setup of Chapter 5, these “agents” do not maximize their own reward functions, and are really just a set of modules inside a single larger agent [Han+24; Yan25]. These workflows are usually hand-engineered. However, they can be improved using RL (see e.g., [Ngu+25; Hui+25; Mo+25]).

6.2.4 Algorithms for single-turn RL

In this section, we discuss RL methods for training LLMs to solve math and reasoning problems. In this setup, there is just a single state, namely the input prompt s ; the action is a sequence of tokens generated by the policy in response, and then the game ends. This is equivalent to a contextual bandit problem, with sequence-valued input (context) and output (action). (We consider the full multi-turn case in Section 6.2.5.)

6.2.4.1 Problem setup

Formally, the goal is to maximize

$$J(\theta) = \mathbb{E}_{s \sim \mathcal{D}, \mathbf{a} \sim \pi_\theta(\mathbf{a}|s)} [R(s, \mathbf{a})] \quad (6.8)$$

⁵However, shortly after the release of R1, the CRO of Open AI (Mark Chen) confirmed that o1 uses some of the same core ideas as R1: https://x.com/markchen90/status/1884303237186216272?s=46&xt=Vx_0-TgDXth-Mt_kw6ggqw.

where s is the context/prompt (sampled from the dataset), and \mathbf{a} is the generated sequence of actions (tokens) sampled from the policy:

$$\pi_{\theta}(\mathbf{a}|s) = \pi_{\theta}(a_1|s)\pi_{\theta}(a_2|s, a_1) \cdots \pi_{\theta}(a_T|s_0, a_1, \dots, a_{T-1}) = \prod_{t=1}^T \pi_{\theta}(a_t|\mathbf{a}_{1:t-1}, s) \quad (6.9)$$

Here $T = |\mathbf{a}|$ is the length of the generated output (which is terminated by generating an `<eos>` token).

We can convert this into an MDP by defining the following deterministic state transition

$$p(s_t|s_{t-1}, a_t) = \delta(s_t = \text{concat}(s_{t-1}, a_t)) = \delta(s_t|s_{t-1}, a_t) \quad (6.10)$$

with initial distribution $\delta(s_0 = s)$. Thus the state s_t is just the set of tokens from the initial prompt s plus the generated tokens up until time t . This definition of state restores the Markov property, and allows us to write the policy in the usual way as $\pi_{\theta}(\mathbf{a}|s) = \prod_{t=1}^T \pi_{\theta}(a_t|s_t)$.

We can also rewrite the objective in standard MDP form as follows:

$$J(\theta) = \mathbb{E}_{s_0 \sim \mathcal{D}} \left[\sum_{t=1}^T \sum_{s_t} \sum_{a_t} \pi_{\theta}(a_t|s_t) \delta(s_t|s_{t-1}, a_t) R(s_t, a_t) \right] \quad (6.11)$$

where $R(s_t, a_t)$ is a sparse reward function defined as follows:

$$R(s_t, a_t) = \begin{cases} 0 & \text{if } t < T \\ R(s_T, a_T = \text{eos}) = R(s, a_1, \dots, a_T) & \text{if } t = T \end{cases} \quad (6.12)$$

where $R(s, a_1, \dots, a_T) = R(\tau)$ is the final trajectory level reward.

In practice, the above approach can overfit to the reward function, so we usually regularize the problem to ensure the policy π_{θ} remains close to the base pre-trained LLM π_{ref} . We can do this by adding a penalty of the form $-\beta D_{\text{KL}}(\pi_{\theta}(a_t|s_t) \parallel \pi_{\text{ref}}(a_t|s_t))$ to the per-token reward $R(s_t, a_t)$.

Below we discuss various algorithms to solve this optimization problem.

6.2.4.2 PPO

A natural approach to training the LLM policy is to use PPO (Section 3.3.3). In the bandit case, we can write the objective as follows:

$$J_{\text{ppo}}(\theta) = \mathbb{E}_{s_n \sim \mathcal{D}} \mathbb{E}_{\mathbf{a}_n \sim \pi_{\text{old}}(\cdot|s_n)} \min(\rho_n(\theta) \cdot A_n \text{ clip}(\rho_n(\theta) \cdot A_n)) \quad (6.13)$$

where we define

$$\rho_n(\theta) = \frac{\pi_{\theta}(\mathbf{a}_n|s_n)}{\pi_{\text{old}}(\mathbf{a}_n|s_n)} \quad // \text{ likelihood ratio} \quad (6.14)$$

$$s_n + \mathbf{a}_n = \text{concat}(s_n, a_{n1}, \dots, a_{nT}) \quad // \text{ full sequence} \quad (6.15)$$

$$A_n = A(s_n + \mathbf{a}_n) = R_n - b_n \quad // \text{ advantage} \quad (6.16)$$

$$R_n = R(s_n + \mathbf{a}_n) \quad // \text{ trajectory level reward} \quad (6.17)$$

$$b_n = b(s_n + \mathbf{a}_n) = V(s_n + \mathbf{a}_n) \quad // \text{ baseline} \quad (6.18)$$

For more details on PPO for LLMs, see [Hua+24a].

6.2.4.3 GRPO

Learning an actor (policy) and a value function (critic) takes twice as much time and memory as just learning a policy. This is problematic for large LLMs. Therefore there have been a bunch of recent methods that replace the parametric value function with MC estimators.⁶

⁶The disadvantage is that estimating the value with MC rollouts requires that we can reset the environment, so that we can generate multiple responses (action trajectories) given the same initial state. This is fine for question answering, but much harder for multi-turn RL, which we discuss below. Also, there is no credit assignment to intermediate states, since we are only estimating the value of the initial state. Thus the method is statistically quite inefficient.

The Group Relative PPO or **GRPO** algorithm of [Sha+24], which was used to train DeepSeek-R1-Zero (discussed in Section 6.2.3.2), is a variant of PPO which replaces the critic network with a Monte Carlo estimate of the value function. In more detail, for each prompt s_n , we generate J answers $\mathbf{a}_n^j \sim \pi_{\text{old}}(\cdot|s_n)$ (called a **group**, often of size $J \sim 8$) which give final rewards R_n^j . We then compute the advantage by subtracting the group average and dividing by the group standard deviation:

$$\hat{A}_n^j = \frac{R_n^j - \mu_n}{\sigma_n} \quad (6.19)$$

where $\mu_n = \text{mean}(R_n^j : j = 1 : J)$ and $\sigma_n = \text{std}(R_n^j : j = 1 : J)$. The use of the normalization term ensures the rewards are **calibrated**, so that a hard problem with low average reward may still result in an update if the deviation from this low mean is large [Bal+25a; Mro25]. (This is related to the use of **reward centering** in continual RL [Nai+24].)

Since the policy generates a sequence, we can expand out the loss for each sequence into a sum of per-token losses. In GRPO, we set the step-level advantage to be equal to the normalized trajectory-level advantage, $\hat{A}_{nt}^j = \hat{A}_n^j$, and we define the likelihood ratio as

$$\rho_{nt}^j(\theta) = \frac{\pi_{\theta}(\mathbf{a}_{nt}^j | a_{n,<t}, s_n)}{\pi_{\text{old}}(\mathbf{a}_{nt}^j | a_{n,<t}, s_n)} \quad (6.20)$$

We then normalize the clipped advantage by the length of each sequence, to ensure that long sequences don't dominate the loss:⁷

$$J_{GRPO}(\theta) = \frac{1}{N} \sum_{n=1}^N \frac{1}{J} \sum_{j=1}^J \frac{1}{|\mathbf{a}_n^j|} \sum_{t=1}^{|\mathbf{a}_n^j|} \min \left(\rho_{nt}^j(\theta) \hat{A}_{nt}^j, \text{clip}(\rho_{nt}(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_{nt}^j \right) \quad (6.21)$$

6.2.4.4 DAPO

In the **DAPO** paper of [Yu+25], they suggest an asymmetric clipping of the likelihood ratio term:

$$\text{clip}(\rho_{nt}(\theta), 1 - \epsilon_{\text{low}}, 1 + \epsilon_{\text{high}}) \hat{A}_{nt}^j \quad (6.22)$$

In particular, they suggest using $\epsilon_{\text{high}} = 0.28 > \epsilon_{\text{low}} = 0.2$ in the clipping term, so that actions which are low-probability under the previous model, and which therefore get large likelihood ratio ρ_{it} , are not clipped as much, which would suppress exploration and result in “entropy collapse” (see also [Liu+25e, Sec 4.2]). In [Xi+25], they propose a way to automatically adapt the clipping ratios.

6.2.4.5 GSPO

In the **GSPO** (Group Sequence Policy Optimization) paper [Zhe+25b], they point out a flaw with GRPO, due to the fact that the importance sampling correction ρ_{nt}^j is applied to each token, even though the reward is evaluated at the trajectory level. This can result in unstable training. They therefore propose to use the following sequence-level objective, following the contextual bandit formulation:

$$J_{GSPO}(\theta) = \frac{1}{N} \sum_{n=1}^N \frac{1}{J} \sum_{j=1}^J \min \left(\rho_n^j(\theta) \hat{A}_n^j, \text{clip}(\rho_n^j(\theta)) \hat{A}_n^j \right) \quad (6.23)$$

where the importance ratio is given by the ratio of sequence level likelihoods [Zhe+23], where they also normalize by length to ensure the magnitude (and hence clipping value) is comparable across sequences:

$$\rho_n^j(\theta) = \left(\frac{\pi_{\theta}(\mathbf{a}_n^j | s_n)}{\pi_{\text{old}}(\mathbf{a}_n^j | s_n)} \right)^{1/|\mathbf{a}_n^j|} = \exp \left(\frac{1}{|\mathbf{a}_n^j|} \sum_{t=1}^{|\mathbf{a}_n^j|} \log \frac{\pi_{\theta}(\mathbf{a}_{nt}^j | \mathbf{a}_{n,<t}^j, s_n)}{\pi_{\text{old}}(\mathbf{a}_{nt}^j | \mathbf{a}_{n,<t}^j, s_n)} \right) \quad (6.24)$$

⁷Some arguments against this kind of normalization can be found in [Liu+25c; Liu+25b].

6.2.4.6 RLOO

It turns out that the division by the standard deviation used by GRPO when normalizing the advantage terms induces a difficulty bias [Liu+25b], in which a very easy or hard prompt s_n may have a low group-level standard deviation σ_n of the corresponding rewards, and dividing by a small σ_n can result in unstable gradients. The problem can be solved using the **Dr GRPO** (GRPO Done Right) method of [Liu+25c], where they just drop the denominator, giving

$$\hat{A}_{nj}^{\text{DrGRPO}} = R_{nj} - \mu_n, \mu_n = \frac{1}{J} \sum_{j=1}^J R_{nj} \quad (6.25)$$

Another approach is to use the Reinforce Leave-One-Out method (**RLOO**) of [Ahm+24], they propose the following baseline:

$$\hat{A}_{nj}^{\text{RLOO}} = R_{nj} - \mu_n, \mu_n = \frac{1}{J-1} \sum_{j=1, j \neq i}^J R_{nj} \quad (6.26)$$

which is the average reward for all the samples in the batch, excluding the current sample.

The RLOO expression is identical (up to a scaling factor of $J/(J-1)$) to the DrGPO expression. To see this, note that

$$\frac{J}{J-1} \hat{A}_{nj}^{\text{DrGRPO}} = \frac{J}{J-1} (R_{nj} - \frac{1}{J} \sum_{j=1}^J R_{nj}) = \frac{J}{J-1} R_{nj} - \frac{1}{J-1} (\sum_{j=1, j \neq n}^J R_{nj} + R_{nj}) \quad (6.27)$$

$$= \frac{J}{J-1} R_{nj} - \mu_n - \frac{1}{J-1} R_{nj} = R_{nj} - \mu_j = \hat{A}_{nj}^{\text{RLOO}} \quad (6.28)$$

6.2.4.7 REINFORCE++

Unfortunately, RLOO estimates a different baseline μ_n for each prompt s_n , which can lead to overfitting on simpler prompts and vulnerability to reward hacking. The **REINFORCE++** algorithm of [Hu+25] computes the mean and standard deviation globally across the batch; this is unbiased and improves the training stability. That is, it uses the following estimate for the advantage:

$$\hat{A}_n^j = \frac{R_n^j - \mu}{\sigma}, \mu = \text{mean}(R_n^j : n = 1 : N, j = 1 : J), \sigma = \text{std}(R_n^j : n = 1 : N, j = 1 : J) \quad (6.29)$$

Another variant is to first subtract off a local group mean and then do global normalization to get

$$\tilde{A}_n^j = R_n^j - \mu_n, \mu_n = \text{mean}(R_n^j : j = 1 : J) \quad (6.30)$$

$$\hat{A}_n^j = \frac{\tilde{A}_n^j - \mu}{\sigma}, \mu = \text{mean}(\tilde{A}_n^j : n = 1 : N, j = 1 : J), \sigma = \text{std}(\tilde{A}_n^j : n = 1 : N, j = 1 : J) \quad (6.31)$$

6.2.4.8 VinePPO

GRPO and related methods computes an unbiased estimate of the value of each token based on rolling out multiple trajectories from the same starting state (prompt). All intermediate states (action tokens) are treated equally, as illustrated in Figure 6.1. For long chains of thought, this can be a poor estimate of the true RTG (value) of an intermediate token. In [Kaz+25], they propose **VinePPO**, which exploits the fact that in language-based environments, it is possible to reset directly to any intermediate state simply by refeeding the partial context; this enables multiple MC rollouts from any state in the trajectory, thus providing a more accurate value estimate without fitting a value function network (which can be biased). Unfortunately this technique is slow, and not general purpose.

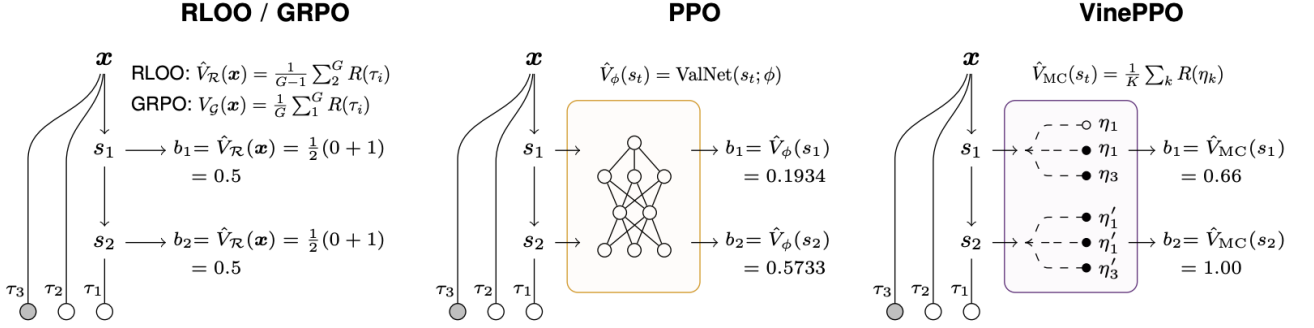


Figure 6.1: Three different ways to estimate the baseline function b_t (needed to compute the advantage function A_t) for intermediate tokens s_t based on $G=3$ trajectory rollouts. (a) MC estimate $b_t = V_{MC}(s_0)$ derived from the initial state. (b) Parametric value function, $b_t = V_{\phi}(s_t)$. (c) MC estimate $b_t = V_{MC}(s_t)$ derived from K rollouts $\eta_k \sim \pi(\cdot|s_t)$ from state s_t . (This exploits the fact that language environments are deterministic and resettable.) Figure from [Kaz+25].

6.2.4.9 Adding a KL regularizer

It is common to add a KL penalty to the per-step reward, to prevent the policy from deviating too far from the base (SFT) LLM:

$$\hat{R}_{n,t}^j = R_n^j - \beta D_{\text{KL}} \left(\pi_{\text{old}}(a_{nt}^j | s_n, \mathbf{a}_{n,<t}^j) \parallel \pi_{\text{ref}}(a_{nt}^j | s_n, \mathbf{a}_{n,<t}^j) \right) \quad (6.32)$$

In GRPO (and many other papers), they use a low-variance MC estimator of KL divergence proposed in <http://joschu.net/blog/kl-approx.html>. The naive estimator of $KL(q, p) = E_{q(a)}[\log(q(a)/p(a))]$ has the form

$$k_1 = \log(r) = \log q(a) - \log p(a) \quad (6.33)$$

where $a \sim q$ and $r(a) = q(a)/p(a)$. Another estimator, which is biased but has lower variance, is the following:

$$k_2 = \frac{1}{2}(\log(r))^2 \quad (6.34)$$

Finally, an unbiased estimator with low variance is

$$k_3 = (r - 1) - \log(r) \quad (6.35)$$

Although k_3 is unbiased, [TM25] and [Hu+25, App B.1] show that its gradient is biased, so they recommend k_2 .

6.2.4.10 DPO

Rather than first fitting a reward model from preference data using the Bradley-Terry model in Section 6.2.2.3, and then optimizing the policy to maximize this, it is possible to optimize the preferences directly, using the **DPO** (Direct Preference Optimization) method of [Raf+23]. (This is sometimes called **direct alignment**.)

We now derive the DPO method. To simplify notation, we use x for the input prompt (initial state s) and y for the output (answer sequence \mathbf{a}). The objective for KL-regularized policy learning is to maximize the following:⁸

$$J(\pi) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(y|x)} \left[R(x, y) - \beta \log \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)} \right] \quad (6.36)$$

⁸In [Hua+25b], they recently showed that it is better to replace the KL divergence with a χ^2 divergence, which quantifies uncertainty more effectively than KL-regularization. The resulting algorithm is provably robust to overoptimization, unlike standard DPO.

Equivalently we can minimize the loss

$$\mathcal{L}(\pi) = \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(y|x)} \left[\log \frac{\pi(y|x)}{\pi_{\text{ref}}(y|x)} - \frac{1}{\beta} R(x, y) \right] \quad (6.37)$$

$$= \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi(y|x)} \left[\log \frac{\pi(y|x)}{\frac{1}{Z(x)} \pi_{\text{ref}}(y|x) \exp(\frac{1}{\beta} R(x, y))} - \log Z(x) \right] \quad (6.38)$$

where

$$Z(x) = \sum_y \pi_{\text{ref}}(y|x) \exp(\frac{1}{\beta} R(x, y)) \quad (6.39)$$

is the partition function. We can now rewrite the loss as

$$\mathcal{L}(\pi) = \mathbb{E}_{x \sim \mathcal{D}} \left[D_{\text{KL}} \left(\pi(y|x) \parallel \frac{1}{Z(x)} \pi_{\text{ref}}(y|x) \exp(\frac{1}{\beta} R(x, y)) \right) - \log Z(x) \right] \quad (6.40)$$

The second term is independent of π so can be dropped. The first term can be minimized by setting

$$\pi^*(y|x) = \frac{1}{Z(x)} \pi_{\text{ref}}(y|x) \exp(\frac{1}{\beta} R(x, y)) \quad (6.41)$$

from which we can derive the optimal reward function as

$$R^*(x, y) = \beta \log \frac{\pi^*(y|x)}{\pi_{\text{ref}}(y|x)} + \beta \log Z(x) \quad (6.42)$$

Plugging this into the Bradley Terry model of Equation (6.1) we get

$$p^*(y_w \succ y_l|x) = \frac{\exp \left(\beta \log \frac{\pi^*(y_w|x)}{\pi_{\text{ref}}(y_w|x)} + \beta \log Z(x) \right)}{\exp \left(\beta \log \frac{\pi^*(y_w|x)}{\pi_{\text{ref}}(y_w|x)} + \beta \log Z(x) \right) + \exp \left(\beta \log \frac{\pi^*(y_l|x)}{\pi_{\text{ref}}(y_l|x)} + \beta \log Z(x) \right)} \quad (6.43)$$

$$= \frac{1}{1 + \exp \left(\beta \log \frac{\pi^*(y_l|x)}{\pi_{\text{ref}}(y_l|x)} - \beta \log \frac{\pi^*(y_w|x)}{\pi_{\text{ref}}(y_w|x)} \right)} \quad (6.44)$$

Thus we can fit the policy by minimizing

$$\mathcal{L}(\theta) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[\log \sigma \left(\beta \log \frac{\pi_{\theta}(y_w|x)}{\pi_{\text{ref}}(y_w|x)} - \beta \log \frac{\pi_{\theta}(y_l|x)}{\pi_{\text{ref}}(y_l|x)} \right) \right] \quad (6.45)$$

The main downside of DPO is that it is limited to learning from preference data, whereas the other policy gradient methods can work with any reward function, including verifiable (non-learned) rewards.

6.2.4.11 Inference-time scaling using posterior sampling

We can view the problem of generating samples from an LLM that maximizes some reward as equivalent to posterior sampling from a **tilted distribution** [Li+23a], that combines the prior $\pi_{\text{ref}}(y|x)$ with a likelihood $p(O = 1|x, y)$, where O is known as an optimality variable, to get the posterior $p(y|x, O = 1) \propto \pi_{\text{ref}}(y|x)p(O = 1|x, y)$. We define $p(O = 1|x, y) = \exp(R(x, y)/\beta)$, as in the “control as inference” paradigm (see Section 3.6), where $\beta > 0$ is a temperature parameter. (For example, $R(x, y) = \log p(O = 1|x, y)$ could be the log-probability that the response y to prompt x is non-offensive, or that y is the correct answer to the question x .) Thus we want to sample from

$$p(y|x, O = 1) = \frac{1}{Z_x} p(O = 1|x, y) \pi_{\text{ref}}(y|x) = \frac{1}{Z_x} \exp(\beta^{-1} R(x, y)) \pi_{\text{ref}}(y|x) \quad (6.46)$$

where $Z_x = p_{\text{ref}}(O = 1|x) = \sum_y p(O = 1|x, y) \pi_{\text{ref}}(y|x)$ is the normalization constant. (Henceforth we omit the conditioning prompt x for notational simplicity.)

An alternative (but equivalent) formulation is to write the target posterior as

$$\pi^*(y|x) \propto e^{\frac{1}{\beta} R(x, y)} \pi_{\text{ref}}(y|x) \quad (6.47)$$

This is the optimal solution to the following KL-regularized RL problem (see e.g., [Bei+25]):

$$\max_{\pi(\cdot|x)} \mathbb{E}_{\pi(y|x)} R(x, y) - \beta D_{\text{KL}}(\pi(\cdot|x) \parallel \pi_{\text{ref}}(\cdot|x)) \quad (6.48)$$

Note that, if we set $\beta = 1$ and define $R(x, y) = \log(\pi_{\text{ref}}(y|x)^{\alpha-1})$, the optimal solution becomes the **tempered distribution** or **power distribution**

$$\pi^*(y|x) \propto \pi_{\text{ref}}(y|x)^{\alpha-1} \pi_{\text{ref}}(y|x)^1 = \pi_{\text{ref}}(y|x)^{\alpha} \quad (6.49)$$

which can be flatter if $\alpha < 1$ or sharper if $\alpha > 1$ than the reference distribution.⁹

Sampling from such distributions can be done using various methods. A simple method is known as **best-of- N** sampling, which just generates N trajectories, and picks the best. This is equivalent to ancestral sampling from the forwards model, and then weighting by the final likelihood (a soft version of rejection sampling). Despite its simplicity, it was shown in [GGV24; Bei+25] that this method achieves an optimal reward vs KL tradeoff. However, performance can decrease when N increases, due to deviating too far from the base model (see e.g., [Hua+25a; FS25]).

A potentially more computationally efficient method is to use **twisted SMC**, which combines particle filtering (a form of sequential Monte Carlo) with a “twist” function, which predicts the future reward given the current state¹⁰, analogous to a value function (see e.g. [NLS19; CP20; Law+18; Law+22]). This is sometimes called **SMC steering**, and has been used in several papers (see e.g., [Lew+23; Zha+24e; Fen+25; Lou+25; Pur+25]). (See also Section 4.2.3 for closely related work which uses twisted SMC to sample action sequences, where the twist is approximated by a learned advantage function.)

The posterior sampling approach discussed above is an example of using more compute at test time to improve the generation process of an LLM. This is known as **test time compute** (see e.g., [Ji+25] for survey). This provides another kind of scaling law, known as **inference time scaling**, besides just improving the size (and training time) of the base model [Sne+24].

6.2.4.12 RLFT as amortized posterior sampling

The disadvantage of decision-time planning (online posterior sampling) is that it can be slow. Hence it can be desirable to “amortize” the cost by fine-tuning the base LLM so it matches the tilted posterior [KPB22]. We can do this by minimizing $J(\theta) = D_{\text{KL}}(\pi_{\theta} \parallel \pi^*)$, where π^* is the tilted distribution in Equation (6.47), and then sampling from π^* .

Following Equation (3.131), we have

$$D_{\text{KL}}(\pi_{\theta}(y) \parallel \pi^*(y)) = \mathbb{E}_{\pi_{\theta}(y)} [\log \pi_{\theta}(y) - \log \Phi(y) - \log \pi_{\text{ref}}(y)] + \log p_{\text{ref}}(O = 1) \quad (6.50)$$

where the first term is the negative ELBO (evidence lower bound), and the second $\log p_{\text{ref}}(O = 1)$ term is independent of the function being optimized (namely π^*). Hence we can minimize the KL by maximizing the ELBO:

$$J'(\theta) = \mathbb{E}_{\pi_{\theta}(y)} [\log \Phi(y)] - D_{\text{KL}}(\pi_{\theta}(y) \parallel \pi_{\text{ref}}(y)) \propto \mathbb{E}_{\pi_{\theta}(y)} [R(y)] - \beta D_{\text{KL}}(\pi_{\theta}(y) \parallel \pi_{\text{ref}}(y)) \quad (6.51)$$

which is exactly the KL-regularized RL objective used in DPO Equation (6.36).

⁹In [KD25] they show that the sharpened distribution upweights tokens with few but high likelihood future paths, while low-temperature sampling (on a per-token basis) upweights tokens with several but low likelihood completions. See also [SSE23; Zha+24e].

¹⁰Formally, the optimal twist function on a partial sequence y is defined as $\Phi^*(y) = \mathbb{E}_{\pi_{\text{ref}}(y')} [\Phi(y') | y \text{ is a prefix of } y']$.

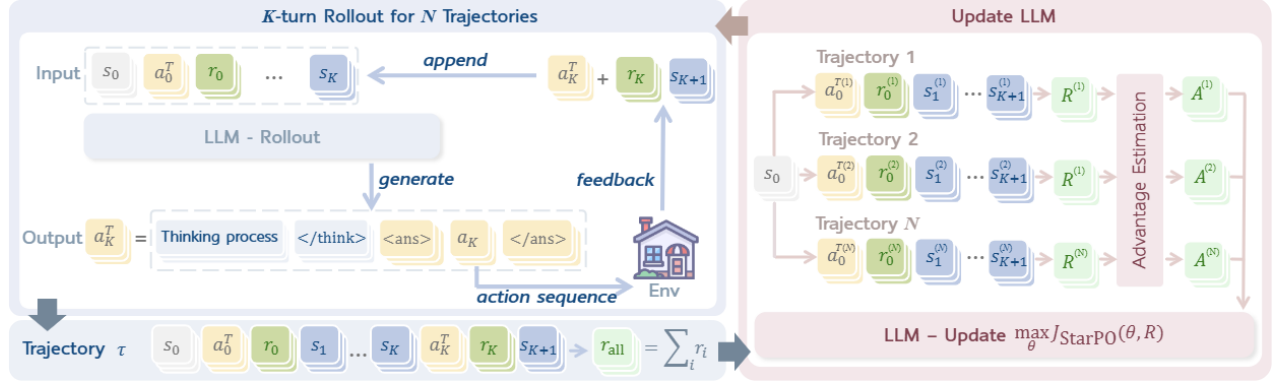


Figure 6.2: Illustration of how to train an LLM to both “think” internally and “act” externally. From the initial state s_0 , we roll out N trajectories, each of length K , and then use the empirical average across all N as a baseline for computing the advantage function. From [Wan+25c]. Used with kind permission of Zihan Wang.

The advantage of this probabilistic (distribution matching) perspective over the RL (reward maximizing) perspective is that it suggests natural alternatives, such as optimizing the inclusive or forwards KL, $D_{\text{KL}}(\pi^* \parallel \pi_\theta)$, which is “mode covering” rather than “mode seeking”. This can prevent “catastrophic forgetting”, in which the tuned policy loses diversity as well as some of its original capabilities [Kor+22]. (This is similar to the advantage of **reweighted wake sleep** training (see e.g., [Le+20; MLR24; Zho+25a]) compared to (amortized) VI training for latent variable models.)

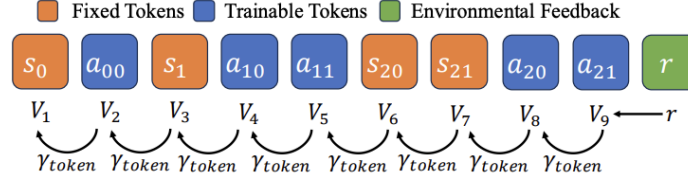
Note, however, that these offline approaches to LLM finetuning (whether using forwards or reverse KL penalty) have the disadvantage, compared to the online (decision-time) approach to inference, that they cannot easily handle hard constraints, since they only train policies that respect the constraints on average. (This is why MPC, which is a similar form of decision-time inference (see Section 4.2.4), is so widely used in the robotics community, where hard constraints are prevalent.)

6.2.5 Algorithms for multi-turn RL

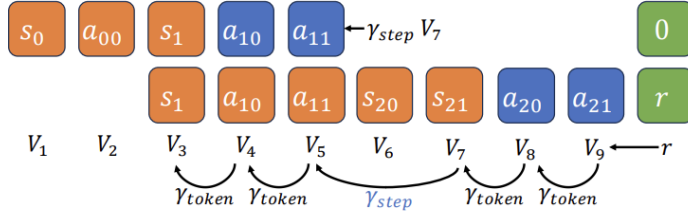
So far, we have considered the bandit setting, in which a single prompt, s , is presented, and then the agent optionally generates some thinking tokens, followed by the answer tokens, and then immediately gets a reward. In this section, we discuss how to train agents that can interact with an external environment, which enables **tool use** (see e.g., [Ngu+25; Mai+25]), or **dialog agents** (see e.g., [Lab+25; Abd+25; Zha+25g]). The difference from the standard contextual bandit LLM reasoning setting is that the effect of an action on the external environment is typically unknown, and may be stochastic, and the reward may be delayed. In addition, the external environment is often stateful, so actions may be irreversible. Training agents in this setting requires true **multi-turn RL** methods, as we have discussed throughout the rest of this document. There are currently only a few papers on using LLMs for this multi-turn setting (see e.g., [Car+23a; WA25; Li+25; Yab+25; Che+25; Liu+25d; Jin+25; Son+25]), but this area will likely grow in the future.

6.2.5.1 Example: RAGEN

As an example of a multi-turn RL system, Figure 6.2 illustrates the **RAGEN** system of [Wan+25c]. Here the policy is prompted to generate a set of thinking tokens, followed by a set of action tokens, until the STOP token is generated (representing the end of action). The policy is used to generate N trajectories, from which the advantage can be estimated using a Monte Carlo estimate. They call this algorithm **StarPO** (State-Thinking-Action-Reward Policy Optimization).



(a) Multi-turn PPO with full memory.



(b) Multi-turn PPO with memory length 1.

Figure 6.3: (a) In standard multi-turn RL, the context grows longer at each step. (b) In turn-level training, we truncate the past context to a fixed window (here, the previous state, previous action, and current state). In addition, we learn a value function, so we can truncate trajectories early and still get a learning signal. State tokens are skipped over when computing the loss, since they are not generated by the model. Finally, we use different discounting factors for steps corresponding to action token generation and environment tokens. From [Che+25]. Used with kind permission of Wen-Tse Chen.

6.2.5.2 Dealing with invalid actions

When interacting with the external environment, we often find that some trajectories result in a an **invalid action** during the generation of an action sequence. There are several ways to deal with such **partial rollouts**. For example, we could just truncate the sequence to length t (the step where the error occurred) to give $\tau_n = (s_{n1}, a_{n1}, \dots, s_{nt}, a_{nt})$. Or we could start a new rollout from the state where the failure occurred, s_{nt} (c.f. [Ngu+25]). Or we can replace the action that caused the error with some new action (e.g., chosen uniformly from the set of legal actions), and then continue the rollout (c.f., [Che+25]).

6.2.5.3 Turn-level training

Methods like RAGEN (Section 6.2.5.1) do not work well in long horizon tasks, for at least two reasons: (1) The context length (which is input to the policy) can grow very large; and (2) the REINFORCE estimator can be very high variance once the number of steps needed to reach a terminal reward becomes large. We can tackle (1) by truncating (or summarizing, using an LLM) the history of previous states and actions, as illustrated in Figure 6.3. We can tackle (2) by learning a value function or critic, and then using the Generalized Advantage Estimator of Section 3.2.2 to compute A_t^n , instead of using the Monte Carlo estimate of the return to go, $G_t^n = \sum_{k=t}^T \gamma^{k-t} r_k^n$, as in REINFORCE. The loss (for the on-policy case) becomes

$$\mathcal{L}_{A2C}(\theta) = \frac{1}{N} \sum_{n=1}^N \sum_{t=1}^{T_n} A_t^n \log \pi_\theta(a_t^n | s_t^n) \quad (6.52)$$

To improve performance, we can warm start the actor (policy) from a prompted, instruction-tuned LLM, and warm start the critic by fitting it to the MC value estimates computed from rollouts using the initial frozen policy, before updating both actor and critic, as suggested in [Che+25].

In [Liu+25d], they propose to replace the GAE estimate (which requires learning a critic) with return

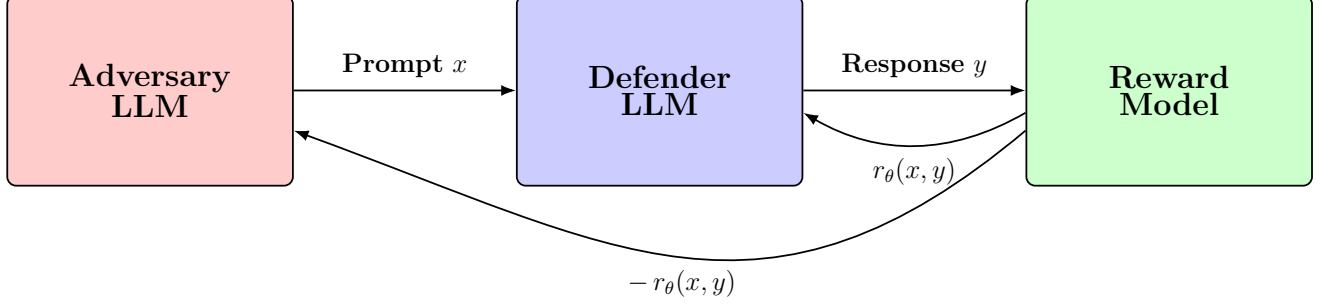


Figure 6.4: Illustration of zero-sum two-player game for training a robust and safe LLM. Based on a slide from Natasha Jacques.

batch normalization (**ReBN**), which is similar to GRPO but defined for the multi-turn setting:

$$A_t^n = \frac{G_t^n - \mu}{\sigma}, \mu = \text{mean}(G_{1:T}^{1:N}), \sigma = \text{std}(G_{1:T}^{1:N}) \quad (6.53)$$

6.2.5.4 Self-play for LLM training

Although RL is mostly used to train LLMs to act in single agent environments (such as reasoning tasks, or tool use), it is also possible to train them to act in multi-agent environments. We discuss multi-agent RL in detail in Chapter 5, but here we briefly mention some LLM-related papers.

In [Liu+25a] they use a two-player zero-sum adversarial game setup to improve LLM safety, by training an LLM adversary to attack an LLM defender. This approach is called **red-teaming**. The basic setup is illustrated in Figure 6.4. The goal is to find the Nash equilibrium of this game:

$$\min_{x \sim \pi_{\text{Adv}}} \max_{y \sim \pi_{\text{Def}}} r_{\theta}(x, y) \quad (6.54)$$

where $r_{\theta}(x, y)$ is a frozen LLM reward model. We assume the reward lies in $[-1, 1]$, where positive scores correspond to safe responses. At the Nash equilibrium, (π_A^*, π_D^*) , one can show that for any prompt $x \sim \pi_A^*$, we have $r_{\theta}(x, \pi_D^*(x)) \geq 0$, meaning that the response is safe.

In practice they implement this by using the same LLM for both adversary and defender, but just giving it a different system prompt. Thus this is an example of self-play (see Section 5.3.5). Training is done using the REINFORCE++ algorithm of Section 6.2.4.7.

In [Bo+25], they introduce the **SPiRAL** algorithm, which trains a single LLM to play various zero-sum two-player games (Tic-Tac-Toe, Kuhn Poker, and Simple Negotiation) using self play. They show that the resulting system is better able to solve math and reasoning problems, which share some commonalities with game play. Training is done using REINFORCE, where the gradient has the form

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{G \sim p(G)} \mathbb{E}_{\tau \sim \pi_{\theta} \times \pi_{\theta} | G} \left[\sum_{p=0}^1 \sum_{t \in \tau_p} A_{G,p}(\tau) \cdot \nabla_{\theta} \log \pi_{\theta}(a_t^p | s_t, p, G) \right] \quad (6.55)$$

where $\pi_{\theta}(a_t^p | s_t, p, G)$ is the policy for player p in game G , and $A_{G,p}$ is a game and role-specific advantage function. The latter is estimated using Role-conditioned Advantage Estimation (**RAE**), which is simply an EMA estimator of the following form

$$b_{G,p} = \alpha b_{G,p} + (1 - \alpha) R_p(\tau) \quad (\text{Update baseline}) \quad (6.56)$$

$$A_{G,p}(\tau) = R_p(\tau) - b_{G,p} \quad (\text{Compute advantage}) \quad (6.57)$$

6.2.6 Alignment and the assistance game

Encouraging an agent to behave in a way that satisfies one or more human preferences is called **alignment**. We can use RL for this, by creating suitable reward functions. However, any objective-maximizing agent may engage in reward hacking (Section 1.3.6.3), in which it finds ways to maximize the specified reward but which humans consider undesirable. This is due to **reward misspecification**, or simply, the law of unintended consequences.

A classic example of this is the poem known as *The Sorcerer’s Apprentice*, written by the German poet Goethe in 1797. This was later made famous in Disney’s cartoon “Fantasia” from 1940. In the cartoon version, Mickey Mouse is an apprentice to a powerful sorcerer. Mickey is tasked with fetching water from the well. Feeling lazy, Mickey puts on the sorcerer’s magic hat, and enchants the broom to carry the buckets of water for him. However, Mickey forgot to ask the magic broom to stop after the first bucket of water has been carried, so soon the room is filled with an army of tireless, water-carrying brooms, until the room floods, at which point Mickey asks the wizard to intervene.

Another example is known as the **cobra effect**. This is named after a policy adopted by a government in India to attempt to reduce the number of cobras in a certain town. To do this, they offered a reward for every cobra head that was turned in. Not surprisingly, this just resulted in a market for raising cobras and then decapitating them, thus increasing the overall number of cobras in the town.

The above parables are typical of many problems that arise when trying to define a reward function that fails to capture all the edge cases we might not have thought of. This is summarized in **Goodhart’s Law**, which states “When a measure becomes a target, it ceases to be a good measure”.

In [Rus19], Stuart Russell proposed a clever solution to this fundamental problem. Specifically, the human and machine are both treated as agents in a two-player, partially observed cooperative game (an instance of a Dec-POMDP, see Section 5.1.3), called an **assistance game**, where the machine’s goal is to maximize the user’s utility (reward) function, which is inferred based on the human’s behavior using inverse RL. That is, instead of trying to learn the reward function using RLHF, and then optimizing that, we treat the reward function as an unknown part of the environment. If we adopt a Bayesian perspective on this, we can maintain a posterior belief over the model parameters, which will incentivize the agent to perform information gathering actions (see Section 7.2.1.2). For example, if the machine is uncertain about whether something is a good idea or not, it will proceed cautiously (e.g., by asking the user for their preference), rather than blindly solving the wrong problem. For more details on this framework, see [Sha+20; Lai+25]. For a more general discussion of (mis)alignment and risks posed by AI agents and “AGI”, see e.g., [Ham+25].

6.3 LLMs for RL

In this section, we discuss how to use LLMs to help create agents that themselves may or may not use language. The LLMs can be used for their prior knowledge, their ability to generate code, their “reasoning” ability, and their ability to perform **in-context learning**. The survey in [Cao+24] groups the literature into four main categories: LLMs for pre-processing the inputs, LLMs for rewards, LLMs for world models, and LLMs for decision making or policies. In our brief presentation below, we follow this categorization. See also e.g., [Spi+24; Hu+24; Pte+24] for more information.

6.3.1 LLMs for pre-processing the input

If the input observations \mathbf{o}_t sent to the agent are in natural language (or some other textual representation, such as JSON), it is natural to use an LLM to process them, in order to compute a more compact representation, $\mathbf{s}_t = \phi(\mathbf{o}_t)$, where ϕ can be the hidden state of the last layer of an LLM. This encoder can either be frozen, or fine-tuned with the policy network. Note that we can also pass in the entire past observation history, $\mathbf{o}_{1:t}$, as well as static “side information”, such as instruction manuals or human hints; these can all be concatenated to form the LLM prompt.

6.3.1.1 Example: AlphaProof

The **AlphaProof** system¹¹ uses an LLM (called the “formalizer network”) to translate an informal specification of a math problem into the formal Lean representation, which is then passed to an agent (called the “solver network”) which is trained, using the AlphaZero method (see Section 4.2.2.1), to generate proofs inside the Lean theorem proving environment. In this environment, the reward is 0 or 1 (proof is correct or not), the state space is a structured set of previously proved facts and the current goal, and the action space is a set of proof tactics. The agent itself is a separate transformer policy network (distinct from the formalizer network) that is a pre-trained LLM, that is fine-tuned on math, Lean and code, and then further trained using RL.

6.3.1.2 VLMs for parsing images into structured data

If the observations are images, it is traditional to use a CNN to process the input, so $\mathbf{s}_t \in \mathbb{R}^N$ would be an embedding vector. However, we could alternatively use a VLM to compute a structured representation, where \mathbf{s}_t might be a set of tokens describing the scene at a high level, or potentially a JSON dictionary. We can then pass this symbolic representation to the policy function. We can also fine tune the VLM with RL (see e.g., [Zha+24a]).

6.3.1.3 Active control of LLM sensor/preprocessor

Note that the information that is extracted will heavily depend on the prompt that is used. Thus we should think of an LLM/VLM as an **active sensor** that we can control via prompts. Choosing how to control this sensor requires expanding the action space of the agent to include computational actions [Che+24d]. Note also that these kinds of “sensors” are very expensive to invoke, so an agent with some limits on its time and compute (which is all practical agents) will need to reason about the value of information and the cost of computation. This is called **metareasoning** [RW91; AT17; Wan+25a]. Devising good ways to train agents to perform both computational actions (e.g., invoking an LLM or VLM) and environment actions (e.g., taking a step in the environment or calling a tool) is an open research problem.

6.3.2 LLMs for rewards

It is difficult to design a reward function to cause an agent to exhibit some desired behavior, as we discussed in Section 1.3.6. Fortunately LLMs can often help with this task, especially when using goal-conditioned RL (Section 1.2.3). We discuss a few approaches below.

In [Kli+24], they present the **Motif** system, that uses an LLM in lieu of a human to provide preference judgements to an RLHF system. In more detail, a pre-trained policy is used to collect trajectories, from which pairs of states, $(\mathbf{o}, \mathbf{o}')$, are selected at random. The LLM is then asked which state is preferable, thus generating $(\mathbf{o}, \mathbf{o}', y)$ tuples, which can be used to train a binary classifier from which a reward model is extracted, a technique known as “AI feedback” (see Section 6.2.2.3). In [Kli+24], the observations \mathbf{o} are text captions generated by the NetHack game. The learned reward model is then used in lieu of the environment reward, or as a shaping function (Section 1.3.6.5). They applied their method to train an agent in the NetHack environment, which has very sparse reward. In [Kli+25b], they extend this method to use a VLM, instead of just an LLM, in order to generate the preference dataset for visual domains, such as MetaWorld. They also show that using the LLM to generate AI feedback works better than using it as a policy (i.e., to directly predict actions). In [Zhe+24a] they extend this work to the online setting, avoiding the need for an informative offline trajectory dataset.

In [Ma+24], they present the **Eureka** system, that learns the reward using bilevel optimization, with RL on the inner loop and LLM-powered evolutionary search on the outer loop. In particular, in the inner loop, given a candidate reward function R_i , we use PPO to train a policy, and then return a scalar quality score $S_i = S(R_i)$. In the outer loop, we ask an LLM to generate a new set of reward functions, R'_i , given a population of old reward functions and their scores, (R_i, S_i) , which have been trained and evaluated in parallel on a fleet of GPUs. The prompt also includes the source code of the environment simulator. Each

¹¹See <https://deepmind.google/discover/blog/ai-solves-imo-problems-at-silver-medal-level/>.

generated reward function R_i is represented as a Python function, that has access to the ground truth state of the underlying robot simulator. The resulting system is able to learn a complex reward function that is sufficient to train a policy (using PPO) that can control a simulated robot hand to perform various dexterous manipulation tasks, including spinning a pen with its finger tips. In [Li+24a], they present a somewhat related approach and apply it to Minecraft.

In [Ven+24], they propose **code as reward**, in which they prompt a VLM with an initial and goal image, and ask it to describe the corresponding sequence of tasks needed to reach the goal. They then ask the LLM to synthesize code that checks for completion of each subtask (based on processing of object properties, such as relative location, derived from the image). These reward functions are then “verified” by applying them to an offline set of expert and random trajectories; a good reward function should allocate high reward to the expert trajectories and low reward to the random ones. Finally, the reward functions are used as auxiliary rewards inside an RL agent.

It is also common to use a VLM to define a reward function (see e.g., [Du+23a; Bau+23; Hua+23; Son+23; Roc+24]).

6.3.3 LLMs for world models

In this section, we discuss how to use LLMs to create world models of the form $p(s_t|s_{1:t-1}, a_{1:t})$, which we will denote by $p(s'|s, a)$ for brevity. We can either do this by treating the LLM itself as a WM (which is then updated using in-context learning), or asking the LLM to generate another artefact, such as some python code, that represents the WM. The advantage of the latter approach is that the resulting WM will be much faster to run, and may be more interpretable. We discuss both versions below.

6.3.3.1 LLMs as world models

In principle it is possible to treat a pre-trained LLM (or other kind of foundation model) as an implicit model of the form $p(s'|s, a)$ by sampling responses to a suitable prompt, which encodes s and a . This rarely works out of the box. However it can be made to work by suitable pre-training.

For example, [Yan+24] presents **UniSim**, which is an action-conditioned video diffusion model trained on large amounts of robotics and visual navigation data. Combined with a VLM reward model, this can be used for decision-time planning as follows: sample a candidate action sequence, generate the corresponding images, feed them to the reward model, score the rollouts, and then pick the best action from this set. This is just standard model-predictive control (Section 4.2.4) in image space with a diffusion WM and a random shooting planning algorithm. Unfortunately, the method can be quite slow, since it needs to call the diffusion model HM times at each planning step, where H is the lookahead horizon and M is the number of samples.

6.3.3.2 LLMs for generating code world models

Calling the LLM at every step to sample from the WM $p(s'|s, a)$ is very slow (see Section 6.3.5), so an alternative is to use the LLM to generate code that represents the world model. This is called a **code world model** (CWM).

One approach is to rely on zero-shot prompting of the LLM to generate the CWM just from a text description of the environment (see e.g., [Sun+24; Won+23]), possibly combined with feedback that checks the validity of the generated model (see e.g., the prompt-based PDDL model learning method of [Gua+23], or the text-code consistency method of [Min+24]). However, below we focus on methods that use feedback from trajectory data, generated either offline or online, as is standard in model-based learning.

In [Dai+24], they present **GIF-MCTS** (Generate, Improve and Fix with Monte Carlo Tree Search) for learning CWMs given a natural language description of the task, and a fixed offline dataset of trajectories (about 10 per task). These trajectories are collected using a behavior policy, which should demonstrate at least some successful trials. The method maintains a representation of the posterior over the WM, $M = p(s'|s, a)$, as a tree of partial programs. At each step, a node is chosen from the tree using the UCT formula (see Section 4.2.2). This node can then be expanded in one of three ways: (G) the LLM is asked to generate code to solve the task, using the current node as a seed program (by adding new lines to it); (I) the LLM is asked

to improve the current code so it passes more unit tests, evaluated on the offline trajectories; (F) the LLM is asked to fix execution bugs in the current code. Each of these tree mutation operations involve passing a custom-written prompt to the LLM, in addition to the program stored in the relevant tree node. They apply their method to fully observed, deterministic environments, with both discrete and continuous actions. The quality of the CWM is measured both in terms of offline prediction accuracy, and also the reward that can be obtained (relative to that of a random policy) when using it inside of a run-time planning algorithm (MCTS for discrete actions, CEM for continuous actions).

In [TKE24], they present **WorldCoder**, which learns a CWM in an online fashion by interacting with the environment, and prompting an LLM. More precisely, it maintains a sample-based representation of the posterior of $p(M|\mathcal{D}(1:t))$, where M is the WM, where the weight for each sampled program ρ are represented by a Beta distribution, $B(\alpha, \beta)$, with initial parameters $\alpha = C + Cr(\rho)$, and $\beta = C + C(1 - r(\rho))$, where $r(\rho)$ is the fraction of unit tests that pass, and C is a constant. This representation about the quality of each program is similar to the one used to represent the reward for the arms in a Bernoulli bandit, and is based on the **REx** (Refine, Explore, Exploit) algorithm of [Tan+24]. At each step, it samples one of these models (programs) from this weighted posterior, and then uses it inside of a planning algorithm, similar to Thompson sampling (Section 7.2.2.2). The agent then executes this in the environment, and passes back failed predictions to the LLM, asking it to improve the WM, or to fix bugs if it does not run. (This refinement step is similar to the I and F steps of GIF-MCTS.) To encourage exploration, they introduce a new learning objective that prefers world models which a planner thinks lead to rewarding states, particularly when the agent is uncertain as to where the rewards are.

We can also use the LLM as a mutation operator inside of an evolutionary search algorithm, as in the **FunSearch** system [RP+24] (recently rebranded as **AlphaEvolve** [Dee25]), where the goal is to search over program space to find code that minimizes some objective, such as prediction errors on a given dataset.

6.3.3.3 LLMs for generating partial code world models

In [Pir+25; Kha+25a] they extend this approach to create a **PoE World Model**, defined in terms of a product of experts. Each term in the product is a distribution over a single element of the state space, and the terms themselves are defined as a product of deterministic experts, each of which is learned using code synthesis. This approach allows parts of the world model to be learned independently, although it does not capture constraints or correlations between the parts.

6.3.4 LLMs for policies

Finally we turn to using LLMs for creating policies. We can either do this by treating the LLM itself as a policy (which is then updated using in-context learning), or asking the LLM to generate some code that represents the policy. We discuss both versions below.

6.3.4.1 LLMs for generating actions

We can sample an action from a policy $\pi(a_t|o_t, h_{t-1})$ by using an LLM, where the input context contains the past data (o_t, h_{t-1}) , and then output token is interpreted as the action. For this to work, the model must be pre-trained on state-action sequences using behavior cloning. See e.g., **Gato** model [Ree+22] **RT-2** [Zit+23], and **RoboCat** [Bou+23].

An alternative approach is to enumerate all possible discrete actions, and use the LLM to score them in terms of their likelihoods given the goal, and their suitability given a learned value function applied to the current state, i.e. $\pi(a_t = k|g, p_t, o_t, h_t) \propto \text{LLM}(w_k|g_t, p_t, h_t)V_k(o_t)$, where g_t is the current goal, w_k is a text description of action k , and V_k is the value function for action k . This is the approach used in the robotics **SayCan** approach [Ich+23], where the primitive actions a_k are separately trained goal-conditioned policies.

Alternatively, we can use a general purpose pre-trained LLM, combined with a suitable prompt chosen by the human, to request the LLM to generate the right kind of output. This approach is used by the **ReAct** paper [Yao+22] which works by prompting the LLM to do some Chain of Thought reasoning (see Section 6.2.3.1) before acting. This approach can be extended by giving feedback on earlier actions, a

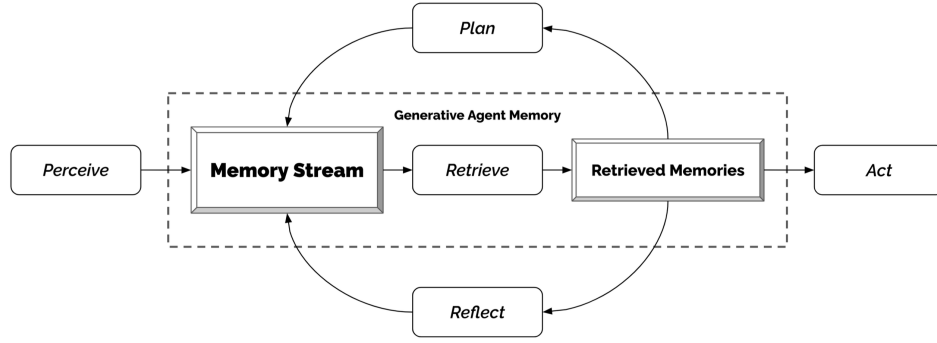


Figure 6.5: Illustration of how to use a pretrained LLM (combined with RAG) as a policy. From Figure 5 of [Par+23]. Used with kind permission of Joon Park.

technique called **Reflexion** [SLG23]. We can also prompt the LLM to first retrieve relevant past examples from an external “memory”, rather than explicitly storing the entire history h_t in the context (this is called **retrieval augmented generation** or **RAG**); see Figure 6.5 for an illustration.

Note that no explicit learning (in the form of parametric updates) is performed in any of these systems; instead they rely entirely on in-context learning and **prompt engineering** / **context engineering** (see e.g., [Zha+25d]). This can offer an alternative to RL finetuning.

6.3.4.2 LLMs for generating code policies

Calling the LLM at every step is very slow (see Section 6.3.5), so an alternative is to use the LLM to generate code that represents (parts of) the policy. This is called a **code policy**.

For example, the **Voyager** system in [Wan+24a] builds up a reusable skill library (represented as Python functions), by alternating between environment exploration and prompting the (frozen) LLM to generate new tasks and skills, given the feedback (environment trajectories) collected so far.

We can also use the LLM as a mutation operator inside of an evolutionary search algorithm, as in the **FunSearch** system [RP+24], where the objective is to maximize performance of the generated policy when deployed in one or more environments. A related method was used in [Ebe+24] to generate code for game playing agent policies.

6.3.4.3 LLMs for generating code actions

An alternative to asking the LLM to generate a code policy, that will then be run without the LLM in the loop, is to ask the LLM to generate code for the current action (see e.g., the **CodeAct** system of [Wan+24c]). This is often better than asking it to call a tool multiple times, since the generated code can represent this “action chunk” with a for-loop, and add extra logic (see e.g., [Ta25]). This is different to a policy, since it is a not a function mapping all states to actions. Instead it is a generating a (potentially closed-loop) plan to be executed from the current state.

6.3.4.4 In-context RL

Large LLMs have shown a surprising property, known as **In-Context Learning** or **ICL** (see e.g., [Don+24]), in which they can be “taught” to do function approximation just by being given (x, y) pairs in their context (prompt), and then being asked to predict the output for a novel x . This can be used to train LLMs without needing to do any gradient updates to the underlying parameters. For a review of methods that apply ICL to RL, see [Moe+25].

6.3.5 Speeding up LLMs

Most LLMs are based on transformers. These work well, but can be very slow, particularly for long trajectories, and/or settings in which each observation takes many tokens to encode (e.g., an image). To explain the problem, we focus on the case of LLM policies, which are non-Markovian models of the form $p(a_t|o_{1:t}, a_{1:t-1})$, which condition on the entire history of past observations and actions. As we discuss below, this model (when implemented with standard transformers) needs $O(t)$ time at step t to generate an action, so generating a trajectory of length T takes $O(T^2)$ time and $O(T)$ memory. This is problematic for **lifelong learning agents**, which continually learn many new tasks (see e.g., [Zhe+25c] for a survey of this topic), as well as problems with real-time constraints [Wen+25].

6.3.5.1 Computational complexity of transformer models

To understand the time and space complexity of autoregressive transformer models in more detail, consider a generic transformer which represents $p(y|x)$; this can be used to generate a response y given a prompt x . The **time-to-first-token** (TTFT) latency — which is the time to generate the first output token y_1 , given $x_{1:N}$, (a process known as the **prefill phase**) — grows as $O(N^2)$, due to the use of full cross attention on x .

To generate M subsequent tokens, the **time-to-iterative-token** generation — which is the time to generate each output token of y during the **decode phase** — is $O(NM + M^2)$ [LCZ25]. To see why, note that we can use **KV caching** to avoid recomputing previous keys and values. Then, when we generate a new token at step t , the model only needs to compute the Query for the single token it just produced. It then computes a $1 \times N_t$ attention vector, where $N_t = N + t$, by comparing to all the previously computed keys, and returns the corresponding weighted sum of previously computed values. Thus the first decoding step takes $O(N + 1)$ time, the second takes $O(N + 2)$ time, all the way up to the M 'th step, which takes $O(N + M)$ time, for a total decoding cost of $O((N + 1) + (N + 2) + \dots + (N + M)) = O(NM + M^2)$.

When using an LLM as a policy, we are basically always in the decode phase since there is no initial prompt. That is, at each step, we add the previously generated action to the context, as well as the external observation, and then we grow the KV-cache. Thus the first action takes $O(1)$ time, the second action takes $O(2)$ time, the t 'th action takes $O(t)$ time.

6.3.5.2 Modern RNNs

To reduce the per-step complexity to $O(1)$, while still maintaining non-Markovianity, it is possible to use RNN-type policies. Traditional RNNs are slow to train, since they have a serial dependency between the latent states. More “modern” RNNs — such as **xLSTMs** [Bec+24; Bec+25] and linear **state space models** (SSMs) such as **Mamba** [GD23] — enjoy fast parallel offline training, like transformers, but still support constant time inference. (See [Som+25] for a review.)

Although these linear RNNs are fast, they are not as accurate as transformers. However, it is possible to create hybrid RNN/transformer models, such as **Griffin** [De+24], which combines recurrent connections with local attention, in alternating layers, to get the best of both worlds. One can also develop RNN-type methods based on **test-time training** (TTT), such as **ATLAS** [Beh+25a; Beh+25b] and **MesaNet** [Osw+25], which perform an inner iterative (or closed form) optimization of the state after each step. (See [Beh+25b; WSF25] for a review.)

These “modern” RNNs can be used for world models (e.g., [Wan+25b] uses Mamba), for policies (e.g., [Sch+25] uses xLSTM and [Hua+24b] uses Mamba), or for representing any other kind of non-Markovian function.

6.4 Implementation details

In addition to the algorithmic issues we have discussed, RL for LLMs (which are very large models) requires a lot of engineering effort, to ensure things run efficiently and stably. In this section, we briefly discuss some of these details.

6.4.1 Policy gradient using Tinker

```

1  async def main(config: Config):
2      wandb_name = (config.wandb_name or config.model_name.split("/")[-1] + f"_{config.env_id}")
3      wandb_name += "_" + datetime.now().strftime("%m%dT%H:%M:%S")
4      save_path = os.path.join("./tinker_output", wandb_name)
5      os.makedirs(save_path, exist_ok=True)
6
7      tokenizer = get_tokenizer(config.model_name)
8      envs = make_envs(config, tokenizer)
9      service_client = tinker.ServiceClient()
10     training_client = await service_client.create_lora_training_client_async(
11         base_model=config.model_name, rank=config.lora_rank)
12     sampling_params = tinker.types.SamplingParams(max_tokens=config.max_tokens)
13     adam_params = types.AdamParams(
14         learning_rate=config.learning_rate, beta1=0.9, beta2=0.95, eps=1e-8)
15
16     for policy_iteration_step in range(config.max_steps):
17         sampling_path = (training_client.save_weights_for_sampler(
18             name=f"{policy_iteration_step:06d}").result().path)
19         sampling_client = service_client.create_sampling_client(model_path=sampling_path)
20
21         episodes_buffer = await collect_episodes_buffer_async(
22             sampling_client, sampling_params, envs, tokenizer, config)
23         transitions = augment_transitions_with_advantages(episodes_buffer, config)
24         training_datums = make_training_data(transitions)
25
26         #fwd_bwd_future = training_client.forward_backward(training_datums, loss_fn="ppo")
27         fwd_bwd_future = training_client.forward_backward(training_datums, loss_fn="importance_sampling")
28         optim_step_future = training_client.optim_step(adam_params)
29         fwd_bwd_result = fwd_bwd_future.result()
30         _ = optim_step_future.result()
31
32     await save_checkpoint_async(training_client, f"{policy_iteration_step:06d}",
33                               log_path=save_path, kind="state", loop_state={"policy_iteration_step": policy_iteration_step})
34

```

Figure 6.6: Main Tinker training loop.

In this section, we show some (slightly abbreviated, but runnable) Python code for implementing a multi-step LLM agent using policy gradient descent. We use the Tinker library¹² to make working with large LLMs easy. In particular, Tinker performs asynchronous rollouts (sampling) and gradient based computation in the cloud, providing a simple low-level interface that can be called from a Python script running on your laptop. Our code is based on a Tinker-based script included as part of the GEM library [Liu+25d].¹³

The main loop is shown in Figure 6.6. At each step of policy updating, it rollouts out some episodes, computes the advantages, converts the data into a format suitable for Tinker, computes the gradient of the loss, and then updates the parameters.

As explained in <https://tinker-docs.thinkingmachines.ai/losses>, Tinker supports several different loss functions, including an importance-weighted version of advantage actor critic (A2C), as discussed in Section 3.4.2.2, and PPO, as discussed in Section 3.3.3. (The importance sampling correction is needed because the distribution q used to collect data may differ from the distribution π_θ which is being trained.)

¹²<https://github.com/thinking-machines-lab/tinker>. We do not use the tinker-cookbook extension.

¹³The full version of the code shown in this chapter is available at https://github.com/probml/pyprobml/blob/master/scripts/tinker_train_gem.py. It is a modified version of https://github.com/axon-rl/gem/blob/main/examples/train_tinker/tinker_train.py.

In more detail, suppose we sample E episodes (trajectories), each with T steps. (In reality, each episode may have a different length, but we omit this from the notation.) Each step of the trajectory contains a sequence of N_s tokens representing the state, \mathbf{s}_t^e , and a sequence of N_a tokens representing the action, \mathbf{a}_t^e . (Again, in reality, the number of tokens per step can differ. Also, note that the action tokens can contain both thinking tokens and “real” action tokens that are sent to the environment, but we omit this distinction from the notation.) Thus a single trajectory τ^e has the form

$$\tau^e = [\underbrace{(s_{1,1}^e, \dots, s_{1,N_s}^e)}_{\mathbf{s}_1^e}, \underbrace{(a_{1,1}^e, \dots, a_{1,N_a}^e)}_{\mathbf{a}_1^e}, r_1, \dots, \underbrace{(s_{T,1}^e, \dots, s_{T,N_s}^e)}_{\mathbf{s}_T^e}, \underbrace{(a_{T,1}^e, \dots, a_{T,N_a}^e)}_{\mathbf{a}_T^e}, r_T] \quad (6.58)$$

To simplify notation, we replace the double index $(\mathbf{s}_t^e, \mathbf{a}_t^e)$ with a single index $(\mathbf{s}^n, \mathbf{a}^n)$, where $n = 1 : N$ indexes the transition, and $N = T \times E$ is the total number of transitions. Then the importance sampling loss is given by

$$\mathcal{L}(\theta) = -\mathbb{E}_{p_{\text{env}}(\mathbf{s})q(\mathbf{a}|\mathbf{s})} \left[\frac{\pi_{\theta}(\mathbf{a}|\mathbf{s})}{q(\mathbf{a}|\mathbf{s})} A(\mathbf{s}, \mathbf{a}) \right] \quad (6.59)$$

$$\approx -\frac{1}{N} \sum_{n=1}^N \frac{\pi_{\theta}(\mathbf{a}^n|\mathbf{s}^n)}{q(\mathbf{a}^n|\mathbf{s}^n)} A(\mathbf{s}^n, \mathbf{a}^n) \quad (6.60)$$

where p_{env} is the environment’s distribution used to sample states, q is the distribution used to sample the actions, and $A(\mathbf{s}^n, \mathbf{a}^n)$ is the advantage. (We discuss how to compute the advantages in Section 6.4.3.)

To simplify notation even further, let us define $\mathbf{x}^n = (\mathbf{s}^n, \mathbf{a}^n)$ as the sequence of $K = N_s + N_a$ tokens representing both state and action, so x_k^n is the k ’th token in the n ’th transition. Let us also define a mask variable $m_k^n = 0$ if token k corresponds to a state token (which is observed by the LLM, but not generated by it), and $m_k^n = 1$ if token k corresponds to an action token. Then the token-level loss is given by

$$\mathcal{L}(\theta) = -\frac{1}{N} \sum_{n=1}^N \sum_{k=1}^K m_k^n \frac{\pi_{\theta}(x_k^n|\mathbf{x}_{1:k-1}^n)}{q(x_k^n|\mathbf{x}_{1:k-1}^n)} A(x_k^n) \quad (6.61)$$

In practice, rather than defining the mask variable, the user can simply set the loss terms corresponding to state tokens to 0, so the loss only comes from the action tokens. Thus the loss per transition (or “datum”), denoted by $\mathcal{L}(\theta|\mathbf{x}^n)$, can be implemented as follows:

```
prob_ratio = torch.exp(target_logprobs - sampling_logprobs)
loss = -(prob_ratio * advantages).sum()
```

Note that we do not divide by the number of total tokens $|\mathbf{x}^n|$, or action tokens $|\mathbf{a}^n|$, since we average over transitions, not over elements of the action sequence. This also avoids a bias towards longer thinking traces, as discussed in [Liu+25b].¹⁴

Tinker also supports a clipped version of the above loss, as in PPO. This is defined by

$$\mathcal{L}(\theta) = -\mathbb{E}_{q(\mathbf{x})} \left[\min \left(\frac{\pi_{\theta}(\mathbf{x})}{q(\mathbf{x})} \cdot A(\mathbf{x}), \text{clip} \left(\frac{\pi_{\theta}(\mathbf{x})}{q(\mathbf{x})}, 1 - \epsilon_{\text{low}}, 1 + \epsilon_{\text{high}} \right) \cdot A(\mathbf{x}) \right) \right] \quad (6.62)$$

where we abuse notation somewhat by assuming $\mathbf{x} = \mathbf{a}$ and ignoring the input prompt \mathbf{s} .

Note that Tinker computes gradients of the above losses with respect to the LoRA (low-rank adaptation) parameters \mathbf{A} and \mathbf{B} of each layer of the transformer, where $\mathbf{W}' = \mathbf{W} + \mathbf{A}\mathbf{B}^T$ are the weights of the fine-tuned model, and \mathbf{W} are the weights of the frozen base model. However, this detail is hidden from users, except for the fact that they must specify the rank of the approximation to use. In [SL25] they show experimentally that a rank one approximation often suffices, which is much more efficient than optimizing all the weights, of which there may be billions.

¹⁴Recall that in standard “reasoning” benchmarks, each trajectory consists of a single state-action pair, so $\tau^e = (\mathbf{s}^e, \mathbf{a}^e)$, where \mathbf{s}^e is the question (prompt) and \mathbf{a}^e is the answer (including thinking tokens). In this case, each episode τ^e is the same as a single transition \mathbf{x}^n , so averaging over transitions is the same as averaging over problems in the dataset (i.e., episodes drawn from the stateless (bandit) distribution of problems).

6.4.2 Rolling out episodes

```

1  async def collect_episodes_buffer_async(sampling_client, sampling_params, envs, tokenizer, config):
2      episodes_buffer = []
3      while True:
4          batch_episodes = await asyncio.gather(
5              *[collect_episode_async(sampling_client, sampling_params, env, tokenizer, config)
6                for env in envs])
7          batch_episodes = [x for x in batch_episodes if x != []]
8          episodes_buffer.extend(batch_episodes)
9          if sum([len(ep) for ep in episodes_buffer]) >= config.batch_size:
10             break
11     return episodes_buffer
12
13  async def collect_episode_async(sampling_client, sampling_params, env, tokenizer, config):
14      transitions = []
15      obs, _ = env.reset()
16      while True:
17          obs = TEMPLATE_FACTORY[config.template](obs) # add system prompt
18          obs_tokens = tokenizer.encode(obs, add_special_tokens=False)
19          sample_result = await sampling_client.sample_async(
20              prompt=types.ModelInput.from_ints(tokens=obs_tokens), num_samples=1,
21              sampling_params=sampling_params)
22
23          sampled_tokens = sample_result.sequences[0].tokens
24          sampled_logprobs = sample_result.sequences[0].logprobs
25          action = tokenizer.decode(sampled_tokens)
26
27          next_obs, reward, terminated, truncated, _ = env.step(action)
28          done = terminated | truncated
29          obs = next_obs
30
31          transitions.append(
32              { "obs_tokens": obs_tokens, "act_tokens": sampled_tokens,
33                "act_logprobs": sampled_logprobs, "reward": reward, "done": done})
34
35          if done:
36             break
37     return transitions

```

Figure 6.7: Code to collect trajectories.

The code to perform the rollouts is shown in Figure 6.7. We perform E rollouts in parallel, where E is the number of copies of the environment. If each environment e returns an observation \mathbf{o}_t^e at each step t pertaining to the current state, we can optionally use a wrapper to concatenate this with the past history of actions and observations, $\mathbf{h}_t^e = (\mathbf{o}_{1:t-1}^e, \mathbf{a}_{1:t-1}^e)$, to get a Markovian state $\mathbf{s}_t^e = (\mathbf{o}_t^e, \mathbf{h}_t^e)$. (We can also compress (or “reflect on”) the history before concatenation, to save space, as suggested in [Che+25].) In addition, we add the system prompt to \mathbf{s}_t^e . This then becomes the context that is passed to the LLM policy to generate \mathbf{a}_t^e , which corresponds to the thinking and action sequence. In GEM, the actual action is extracted from the sequence of tokens \mathbf{a}_t^e by the environment, but we could instead make the agent parse its own generated output before passing it to the environment.

6.4.3 Computing the advantages

In Figure 6.8 we show code to augment each transition with an estimate of its advantage. To do this, we first compute the return-to-go

$$G_t^e = r_t^e + \gamma G_{t+1}^e \quad (6.63)$$

where $G_T^e = r_T^e$ is the initial condition, where T is the first step for which $\text{done}=1$. We then use the return-batch normalization (ReBN) method of [Liu+25d] — which is a multi-turn variant of GRPO — to

```

1  def augment_transitions_with_advantages(epochs_buffer, config):
2      transitions = []
3      for episode in epochs_buffer:
4          # Augment each (s, a, r) transition with MC estimate of return to go.
5          rewards = [transition["reward"] for transition in episode]
6          cur = 0.0
7          for i in reversed(range(len(rewards))):
8              cur = rewards[i] + config.gamma * cur
9              episode[i]["return"] = cur
10         transitions.extend(episode)
11
12     # return batch normalization
13     if config.use_rebn:
14         returns = torch.tensor([transition["return"] for transition in transitions]).float()
15         returns = (returns - returns.mean()) / (returns.std() + 1e-9)
16         for i, transition in enumerate(transitions):
17             transition["return"] = returns[i].item()
18
19     # subsample to make a constant batch size
20     if len(transitions) > config.batch_size:
21         transitions = np.random.choice(transitions, config.batch_size, replace=False)
22
23     return transitions

```

Figure 6.8: Code to compute advantages.

convert these into advantages, using

$$A_t^e = \frac{G_t^e - \mu}{\sigma} \quad (6.64)$$

where $\mu = \text{mean}(\{G_t^e : t = 1 : T_e, e = 1 : E\})$, and $\sigma = \text{std}(\{G_t^e : t = 1 : T_e, e = 1 : E\})$. We then flatten the data to get a list of $N = E \times T$ transitions $(\mathbf{s}^n, \mathbf{a}^n, r^n, A^n)$.

6.4.4 Computing token level loss

Finally, Figure 6.9 shows the code to convert the set of transitions into a set of “training datums”, each of which contains the input tokens $\mathbf{x}^n = (\mathbf{s}^n, \mathbf{a}^n[: -1])$, the target tokens $\mathbf{y}^n = (\mathbf{s}^n[1 :], \mathbf{a}^n)$, the masked sampling log probabilities $\log q(x_k^n)$ and the masked token advantages A_k^n , which we set to be equal to the turn-level advantages A^n , derived from the turn-level rewards R^n . (The masking ensures the state tokens contribute 0 to the loss.)

We can optionally add a KL penalty to the policy to ensure it does not deviate too far from the base model. However, it is not correct to add this to the advantage itself, as pointed out in [Zha+25e; TM25]. Instead, we modify the token-level reward to be

$$R_k^n = R^n + \beta D_{\text{KL}}(\pi_{\theta}(a_k^n | \mathbf{a}_{1:k-1}^n, \mathbf{s}^n) \parallel \pi_0(a_k^n | \mathbf{a}_{1:k-1}^n, \mathbf{s}^n)) \quad (6.65)$$

where π_0 is the reference prior, and R^n is the reward from the environment after transition n .

6.4.5 Computing metrics related to training stability

In practice it is important to log various metrics, to monitor the training process. For example, we might want to compute $D_{\text{KL}}(q \parallel p)$, where q is the sampling distribution (used to rollout the episodes), and p is the training distribution (used to compute gradients). We can estimate the KL using k_1 estimator from Equation (6.33) or k_2 estimator from Equation (6.34). If the KL exceeds 0.01, it means that learning is very off-policy, and results might be unstable.¹⁵ In addition, we can compute the entropy of the policy. We want to ensure this is initially not too small, to enable exploration, but that it does not blow up over time. See Figure 6.10 for some code to compute these metrics.

¹⁵This heuristic is from <https://tinker-docs.thinkingmachines.ai/rl/rl-hyperparams>.

```

1  def make_training_data(transitions):
2      training_datums = []
3      for transition in transitions:
4          ob_len_m1 = len(transition["obs_tokens"]) - 1 # -1 due to shifting
5          tokens = transition["obs_tokens"] + transition["act_tokens"]
6
7          input_tokens = tokens[:-1]
8          target_tokens = tokens[1:]
9          all_logprobs = [0.0] * ob_len_m1 + transition["act_logprobs"]
10         all_advantages = [0.0] * ob_len_m1 + [transition["return"]] * (len(input_tokens) - ob_len_m1)
11
12         datum = types.Datum(
13             model_input=types.ModelInput.from_ints(tokens=input_tokens),
14             loss_fn_inputs={
15                 "target_tokens": TensorData.from_torch(torch.tensor(target_tokens)),
16                 "logprobs": TensorData.from_torch(torch.tensor(all_logprobs)),
17                 "advantages": TensorData.from_torch(torch.tensor(all_advantages)),
18             },
19         )
20         training_datums.append(datum)
21     return training_datums
22

```

Figure 6.9: Code to compute data for token level loss computation.

6.4.6 Example

In this section, we give a brief example. We consider the **GuessTheNumber-v0-easy** game, in which the agent must guess the secret number n_e randomly chosen by the environment at the start of each episode e . In the easy version, the number is between $n_{\min} = 1$ and $n_{\max} = 10$, and the max number of turns is $T = 4$. If the agent guesses correctly in 4 or fewer turns, it gets a success reward of 1. If it runs out of turns, it gets a reward of $1 - d/(n_{\max} - n_{\min})$, where $d = |n^e - a_T^e|$ is the distance between the true answer and the final guess. If it generates an action sequence \mathbf{a}_t that cannot be parsed in to a valid numeric action, the agent gets a format penalty of -0.1. For details, see the code.¹⁶

The optimal strategy is to perform binary search. So for the easy version of the game, the expected number of turns is $\log_2(10) = 3.32$. We fine-tune a Qwen3-4B-Instruct-2507 model using 20 steps of ReBN, rolling out $E = 32$ episodes per step. We use LoRA of rank 8, a learning rate of 4×10^{-5} , and a discount factor of $\gamma = 0.9$. Further details on the hyper-parameters can be found in the code. Training takes about 20 minutes on a laptop, using the standard Tinker client.

In Figure 6.11(a), we plot the average reward vs training steps, and in Figure 6.11(c), we plot the average number of action tokens per turn vs training steps. We see that initially the model gets an average reward of 0.8, but that the resulting action sequences are quite verbose (about 900 tokens). At the end of training, we see that the average reward has improved to 1.0, and the answer length has reduced to about 500 tokens per step. In Figure 6.11(b), we see that the average number of turns is always around 3, since the initial base model already “knows” that the optimal strategy is binary search. (But by using $\gamma < 1$, we discourage taking too many turns; see Figure 6.12 for an example.) Finally in Figure 6.11(d), we see the entropy of the policy reduces over time, as the model specializes to this task.

¹⁶https://github.com/axon-rl/gem/blob/main/gem/envs/game_env/guess_the_number.py.

```

1  def compute_policy_metrics(config, transitions, fwd_bwd_result):
2      # compute policy entropy and sampler-learner difference
3      act_token_logprobs = []
4      act_token_diffs = []
5      for i in range(config.batch_size):
6          transition = transitions[i]
7          train_output = fwd_bwd_result.loss_fn_outputs[i]
8          nact = len(transition["act_logprobs"])
9          act_token_logprobs.extend(transition["act_logprobs"])
10         sampling_token_logprobs = torch.tensor(transition["act_logprobs"])
11         policy_token_logprobs = train_output["logprobs"].to_torch()[-nact:]
12         # kl = E_{qsample(a)} [log q_sample(a) - log p_policy(a)]
13         act_token_diffs.append(sampling_token_logprobs - policy_token_logprobs)
14
15     act_token_diffs = torch.cat(act_token_diffs)
16     kl_sample_train_v1 = act_token_diffs.mean().item()
17     kl_sample_train_v2 = 0.5 * (act_token_diffs**2).mean().item()
18     return {
19         "token_entropy": -torch.tensor(act_token_logprobs).mean().item(),
20         "kl_sample_train_v1": kl_sample_train_v1,
21         "kl_sample_train_v2": kl_sample_train_v2,
22     }

```

Figure 6.10: Code to compute metrics.

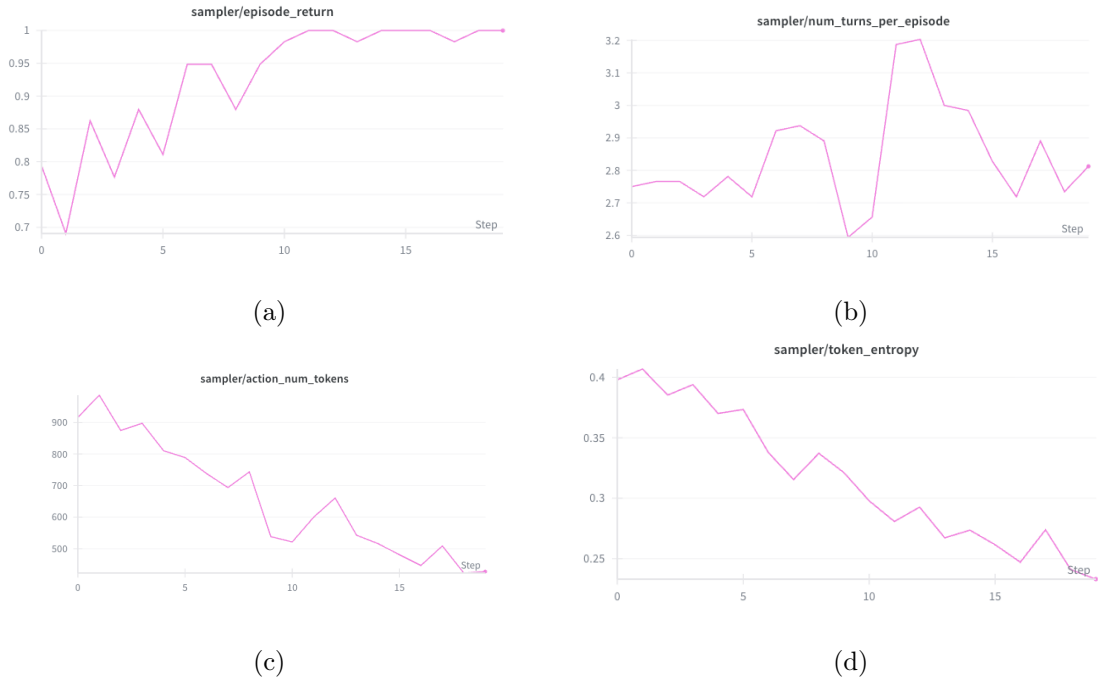


Figure 6.11: Results of RL on the guess-the-number game from GEM.

STATE: You are playing Guess The Number. You have to guess the number between 1 and 50 (inclusive) within 50 turns. As you play, the history of your guesses will be appended below. Use the information to complete the game before you run out of guesses. Enter your first guess to start the game.

At turn 1, you guessed 25, and the target number is lower than 25.
 At turn 2, you guessed 12, and the target number is higher than 12.
 At turn 3, you guessed 18, and the target number is higher than 18.
 At turn 4, you guessed 21, and the target number is higher than 21.
 At turn 5, you guessed 23, and the target number is lower than 23.
 Enter your next guess.

ACTION: To determine the target number between 1 and 50, we can use a binary search strategy, which is effective for narrowing down a range when given hints like higher or lower.

Step-by-Step Reasoning:

1. **Initial Guess (Turn 1):** - The target number is lower than 25. - This means the target number is in the range [1, 24].
2. **Second Guess (Turn 2):** - The target number is higher than 12. - This means the target number is in the range [13, 24].
3. **Third Guess (Turn 3):** - The target number is higher than 18. - This means the target number is in the range [19, 24].
4. **Fourth Guess (Turn 4):** - The target number is higher than 21. - This means the target number is in the range [22, 24].
5. **Fifth Guess (Turn 5):** - The target number is lower than 23. - This means the target number is in the range [22, 22].

Final Guess:
 The target number is in the range [22, 22], which means the target number is 22. Final Answer:
`\boxed{22}`

Figure 6.12: Example of a state and action at the end of turn 5 of guess-the-number after training the policy with $\gamma = 0.9$.

Chapter 7

Other topics in RL

In this section, we briefly mention some other important topics in RL.

7.1 Regret minimization

In this book, we have assumed the environment is a stochastic process (e.g., an MDP), and the agent is trying to maximize its expected future utility (i.e., minimize its risk) given the data it has seen so far in the past, as is standard in the Bayesian approach to decision making. (If there are multiple agents, we must add them to the model, as in Chapter 5.)

An alternative approach is to assume the agent encounters an arbitrary stream of observations, and it must choose actions that are as close as possible to what an optimal policy would have achieved, even if the sequence is chosen by an unknown **adversary**. Algorithms that achieve this goal are said to be **regret minimization** algorithms. We give more details below.

7.1.1 Regret for static MDPs

In online learning, we assume the decision at each step is to choose the parameters θ_t such that we minimize the loss at that step, such as $\ell_t(\theta) = (y_t - f(\mathbf{x}_t, \theta))^2$. Let $\theta_t = \pi(\mathbf{x}_{1:t}, y_{1:t-1})$ be the policy for choosing parameters. We then define the **regret** to be

$$\text{Regret}_T(\pi) = \left[\sum_{t=1}^T \ell_t(\theta_t) - \ell_t(\theta^*) \right] \quad (7.1)$$

where

$$\theta^* = \arg \min_{\theta \in \Theta} \sum_{t=1}^T \ell_t(\theta) \quad (7.2)$$

See e.g., [Haz19] for details.

In the case of the MDP, the regret of a policy π as the difference between its expected return and the expected return of an optimal policy π^* .¹

$$\text{Regret}_T(\pi; M, \Pi) = \mathbb{E}_{s_t \sim M(\cdot | s_{t-1}, a_t), a_t \sim \pi(\cdot | s_t), a_t^* \sim \pi^*(\cdot | s_t)} \left[\sum_{t=1}^T (r_t(s_t, a_t^*) - r_t(s_t, a_t)) \right] \quad (7.3)$$

where

$$\pi^* = \underset{\pi \in \Pi}{\operatorname{argmax}} \mathbb{E}_{s_0 \sim M} [V^\pi(s_0 | M)] \quad (7.4)$$

¹We can also define the regret in terms of value functions: $\text{Regret}_T(\pi | M, \Pi) = \mathbb{E}_{s_0 \sim M} [V_T^{\pi^*}(s_0 | M) - V_T^\pi(s_0 | M)]$, where $V^\pi(s | M)$ refers to the value of policy π starting from state s in MDP M .

is the optimal policy from some policy class Π which has access to the true MDP M . This is often referred to as the **best policy in hindsight**, since if we average over enough sequences (or over enough time steps), the policy we wish we had chosen will of course be the optimal policy that knows the true environment.

Since the true MDP is usually unknown, we can define the **maximum regret** of a policy as its worst case regret wrt some class of models \mathcal{M} :

$$\text{MaxRegret}_T(\pi; \mathcal{M}, \Pi) = \max_{M \in \mathcal{M}} \text{Regret}_T(\pi | M, \Pi) \quad (7.5)$$

We can then define the **minimax optimal policy** as the one that minimizes the maximum regret:²

$$\pi_{MM}^*(\mathcal{M}, \Pi) = \operatorname{argmin}_{\pi \in \Pi} \max_{M \in \mathcal{M}} \text{Regret}_T(\pi | \pi^*(M), M) \quad (7.6)$$

The main quantity of interest in the theoretical RL literature is how fast the regret grows as a function of time T . In the case of a tabular episodic MDP, the optimal minimax regret is $O(\sqrt{HSAT})$ (ignoring logarithmic factors), where H is the horizon length (number of steps per episode), S is the number of states, A is the number of actions, and T is the total number of steps. When using parametric functions to define the MDP, the bounds depend on the complexity of the function class. For details, see e.g., [AJO08; JOA10; LS19].

7.1.2 Regret for non-stationary MDPs

When the world can change, there may be no single optimal policy π^* we can compare to. Instead, the **dynamic regret** (aka **adaptive regret**) compares to a sequence of optimal policies:

$$\text{DynamicRegret}_T(\pi_{1:T} | M_{1:T}, \Pi) = \mathbb{E}_{s_t \sim M_t, a_t \sim \pi_t, a_t^* \sim \pi_t^*} \left[\sum_{t=1}^T (r(s_t, a_t^*) - r(s_t, a_t)) \right] \quad (7.7)$$

where

$$\pi_t^* = \operatorname{argmax}_{\pi \in \Pi} V^\pi(s_t | M_t) \quad (7.8)$$

is the optimal policy at that moment in time.

To compute bounds on the optimal dynamic regret, we need to make assumptions about how often the world changes, and by how much. This is called a **variational budget**. This is defined as

$$\text{VB}_T = \sum_{t=2}^T \text{dist}(\mathcal{M}_t, \mathcal{M}_{t-1}) \quad (7.9)$$

where the distance function measures the similarity of the MDPs at adjacent episodes (e.g., ℓ_1 distance of the reward and transition functions). The optimal dynamic regret can then be bounded in terms of VB (see e.g., [CSLZ23; Aue+19]).

7.1.3 Minimizing regret vs maximizing expected utility

The Bayes optimal agent is the one that maximizes its expected utility (minimizes its risk), where we take expectations not only over the sequence of observations and rewards, but also over the unknown environment M itself, rather than assuming it is known [Gha+15]. That is, the optimal learning algorithm \mathcal{A} is given by

$$\mathcal{A}_{\text{Bayes}}^*(P_0) = \operatorname{argmax}_{\mathcal{A}} U_T(\mathcal{A} | P_0) \quad (7.10)$$

$$U_T(\mathcal{A} | P_0) = \mathbb{E}_{\mathcal{M} \sim P_0(\mathcal{M})} \left[\mathbb{E}_{s_t \sim \mathcal{M}, a_t \sim \pi_t, \pi_t = \mathcal{A}(s_{1:t-1}, a_{1:t-1}, r_{1:t-1})} \left[\sum_{t=1}^T r(s_t, a_t) \right] \right] \quad (7.11)$$

²We can also consider a non-stochastic setting, in which we allow an adversary to choose the state sequence, rather than taking expectations over them. In this case, we take the maximum over individual sequences rather than models. For details, see [HS22].

Aspect	Bayes-Optimal (BAMDP)	Regret-Minimizing (Minimax)
Knowledge	Requires a known prior over MDPs	No prior; judged against best policy in hindsight
Objective	Maximize expected return under the prior	Minimize regret w.r.t. optimal policy in true MDP
Exploration	Performs optimal Bayesian exploration	Often uses optimism or randomness (e.g., UCB, TS)
Adaptation	Fully adaptive via posterior updates	May use confidence bounds, resets, or pessimism
Setting	Bayesian RL	Frequentist or adversarial RL

Table 7.1: Key differences between Bayes-optimal and regret-minimizing policies in RL.

where $P_0(\mathcal{M})$ is our prior over models, and \mathcal{A} is our learning algorithm that generates the policy (decision procedure) to use at each step, as discussed in Section 1.1.3. Note that the uncertainty over models automatically encourages the optimal amount of exploration, as we discuss in Section 7.2.1.2. Note also that if we can do exact inference, the optimal algorithm is uniquely determined by the prior P_0 .

By contrast, the regret minimizing policy is the one that minimizes the maximum regret

$$\mathcal{A}_{MM}^*(\mathcal{M}, \Pi) = \operatorname{argmin}_{\pi \in \Pi} \max_{M \in \mathcal{M}} \operatorname{Regret}_T(\mathcal{A}|M, \Pi) \quad (7.12)$$

where we define the regret of a learning algorithm as

$$\operatorname{Regret}_T(\mathcal{A}|M, \Pi) = \mathbb{E}_{s_t \sim M, a_t \sim \pi_t, \pi_t = \mathcal{A}(s_{1:t-1}, a_{1:t-1}, r_{1:t-1}), a_t^* \sim \pi^*} \left[\sum_{t=1}^T (r(s_t, a_t^*) - r(s_t, a_t)) \right] \quad (7.13)$$

where π^* is given in Equation (7.4). Unlike the Bayesian case, we must now manually design the algorithm to solve the exploration-exploitation problem (e.g., using the Thompson sampling method of Section 7.2.2 or the UCB method of Section 7.2.3), i.e., there is no automatic solution to the problem.

This distinction between minimizing risk and minimizing regret is equivalent to the standard difference between Bayesian and frequentist approaches to decision making (see e.g., [Mur23, Sec 34.1]). The advantage of the Bayesian approach is that it can use prior knowledge (e.g., based on experience with other tasks, or knowledge from an LLM) to adapt quickly to changes, and to make predictions about the future, allowing for optimal long-range planning. The advantage of the regret-minimizing approach is that it avoids the need for a specifying prior over models, it can be robustly adapt to unmodeled changes, and it can handle adversarial (non-stochastic) noise. See Table 7.1 for a summary, and [HT15] for more discussion.

7.2 Exploration-exploitation tradeoff

In this section, we discuss solutions to the exploration-exploitation tradeoff that go beyond the simple heuristics introduced in Section 1.3.5.

7.2.1 Optimal (Bayesian) approach

We can compute an optimal solution to the exploration-exploitation tradeoff by adopting a Bayesian approach to the problem, where we augment the state space with our beliefs about the underlying model, as discussed in Section 1.2.6.

7.2.1.1 Bandit case (Gittins indices)

In the special case of context-free bandits with a finite number of arms, the optimal policy of this belief state MDP can be computed using dynamic programming. To explain this, we follow the presentation of [KWW22, Sec 15.5], and consider a Bernoulli bandit with n arms. Let the belief state be denoted by $b = (w_1, l_1, \dots, w_n, l_n)$, where w_a is the number of times arm a has won (given reward 1) and l_a is the number of times arm a has lost (given reward 0). Using Bellman’s equation, and the expression for the probability of winning under a beta-Bernoulli distribution with a uniform prior, we have

$$V^*(b) = \max_a Q^*(b, a) \quad (7.14)$$

$$Q^*(b, a) = \frac{w_a + 1}{w_a + l_a + 2} (1 + V^*(\dots, w_a + 1, l_a, \dots)) \quad (7.15)$$

$$+ \left(1 - \frac{w_a + 1}{w_a + l_a + 2}\right) V^*(\dots, w_a, l_a + 1, \dots) \quad (7.16)$$

In the finite horizon case, with h steps, We can compute Q^* using dynamic programming. We start with terminal belief states b with $\sum_a (w_a + l_a) = h$, where $V^*(b) = 0$. We then work backwards to states b satisfying $\sum_a (w_a + l_a) = h - 1$, and then applying the above equation recursively until time step 0.

Unfortunately, although this process is optimal, the number of belief states is $O(h^{2n})$, rendering it intractable. Fortunately, for the infinite horizon discounted case, the problem can be solved efficiently using **Gittins indices** [Git89] (see [PR12; Pow22] for details). However, these optimal methods do not extend to contextual bandits, where the problem is provably intractable [PT87].

7.2.1.2 MDP case (Bayes Adaptive MDPs)

We can extend the above techniques to the MDP case by constructing a **BAMDP**, which stands for “Bayes-Adaptive MDP” [Duf02]. The basic idea is quite simple. We define an MDP with an augmented state-space, consisting of the original state s and the belief state b , representing a distribution over the model parameters. The transition function is given by

$$\mathcal{T}'(s', b' | s, b, a) = \delta(b' = BU(s, b, a, s')) P(s' | s, b, a) \quad (7.17)$$

where BU is the (deterministic) Bayes updating procedure (e.g., incrementing the pseudo counts of the Dirichlet distribution, in the case of a discrete MDP), and the second term is the posterior predictive distribution over states:

$$P(s' | s, b, a) = \int b(\theta) \mathcal{T}(s' | s, a; \theta) d\theta \quad (7.18)$$

The rewards and actions of the augmented MDP are the same as the base MDP. Thus Bellman’s equation gives us

$$V^*(s, b) = \max_a \left(R(s, a) + \gamma \sum_{s'} P(s' | s, b, a) V^*(s', BU(s, b, a, s')) \right) \quad (7.19)$$

Unfortunately, this is computationally intractable to solve. Fortunately, various approximations have been proposed (see e.g., [Zin+21; AS22; Mik+20]).

7.2.2 Thompson sampling

The fully Bayesian approach is computationally intractable. A common approximation is to use **Thompson sampling** [Tho33], also called **probability matching** [Sco10]. We start by describing this in the bandit case, then extend to the MDP case. For more details, see [Rus+18]. (See also [Ger18] for some evidence that humans use Thompson-sampling like mechanisms.)

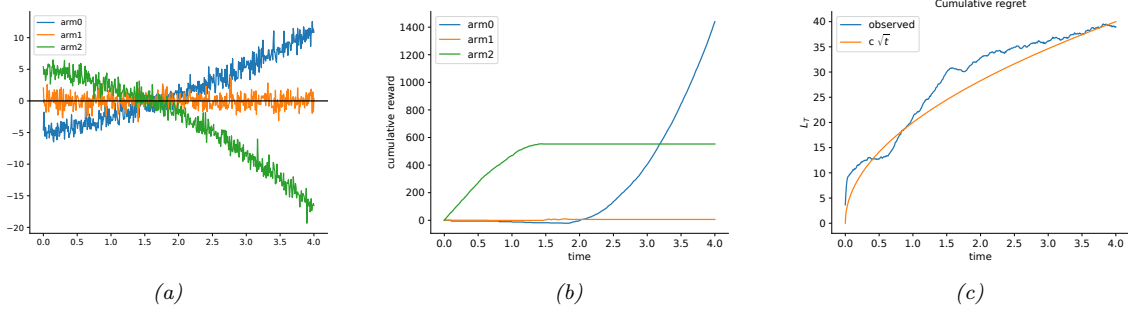


Figure 7.1: Illustration of Thompson sampling applied to a linear-Gaussian contextual bandit. The context has the form $\mathbf{s}_t = (1, t, t^2)$. (a) True reward for each arm vs time. (b) Cumulative reward per arm vs time. (c) Cumulative regret vs time. Generated by `thompson_sampling_linear_gaussian.ipynb`.

7.2.2.1 Bandit case

In Thompson sampling, we define the policy at step t to be $\pi_t(a|\mathbf{s}_t, \mathbf{h}_t) = p_a$, where p_a is the probability that a is the optimal action. This can be computed using

$$p_a = \Pr(a = a_* | \mathbf{s}_t, \mathbf{h}_t) = \int \mathbb{I}\left(a = \operatorname{argmax}_{a'} R(\mathbf{s}_t, a'; \boldsymbol{\theta})\right) p(\boldsymbol{\theta} | \mathbf{h}_t) d\boldsymbol{\theta} \quad (7.20)$$

If the posterior is uncertain, the agent will sample many different actions, automatically resulting in exploration. As the uncertainty decreases, it will start to exploit its knowledge.

To see how we can implement this method, note that we can compute the expression in Equation (7.20) by using a single Monte Carlo sample $\tilde{\boldsymbol{\theta}}_t \sim p(\boldsymbol{\theta} | \mathbf{h}_t)$. We then plug in this parameter into our reward model, and greedily pick the best action:

$$a_t = \operatorname{argmax}_{a'} R(\mathbf{s}_t, a'; \tilde{\boldsymbol{\theta}}_t) \quad (7.21)$$

This sample-then-exploit approach will choose actions with exactly the desired probability, since

$$p_a = \int \mathbb{I}\left(a = \operatorname{argmax}_{a'} R(\mathbf{s}_t, a'; \tilde{\boldsymbol{\theta}}_t)\right) p(\tilde{\boldsymbol{\theta}}_t | \mathbf{h}_t) = \Pr_{\tilde{\boldsymbol{\theta}}_t \sim p(\boldsymbol{\theta} | \mathbf{h}_t)}\left(a = \operatorname{argmax}_{a'} R(\mathbf{s}_t, a'; \tilde{\boldsymbol{\theta}}_t)\right) \quad (7.22)$$

Despite its simplicity, this approach can be shown to achieve optimal regret (see e.g., [Rus+18] for a survey). In addition, it is very easy to implement, and hence is widely used in practice [Gra+10; Sco10; CL11].

In Figure 7.1, we give a simple example of Thompson sampling applied to a linear regression bandit. The context has the form $\mathbf{s}_t = (1, t, t^2)$. The true reward function has the form $R(\mathbf{s}_t, a) = \mathbf{w}_a^\top \mathbf{s}_t$. The weights per arm are chosen as follows: $\mathbf{w}_0 = (-5, 2, 0.5)$, $\mathbf{w}_1 = (0, 0, 0)$, $\mathbf{w}_2 = (5, -1.5, -1)$. Thus we see that arm 0 is initially worse (large negative bias) but gets better over time (positive slope), arm 1 is useless, and arm 2 is initially better (large positive bias) but gets worse over time. The observation noise is the same for all arms, $\sigma^2 = 1$. See Figure 7.1(a) for a plot of the reward function. We use a conjugate Gaussian-gamma prior and perform exact Bayesian updating. Thompson sampling quickly discovers that arm 1 is useless. Initially it pulls arm 2 more, but it adapts to the non-stationary nature of the problem and switches over to arm 0, as shown in Figure 7.1(b). In Figure 7.1(c), we show that the empirical cumulative regret in blue is close to the optimal lower bound in red.

7.2.2.2 MDP case (posterior sampling RL)

We can generalize Thompson sampling to the (episodic) MDP case by maintaining a posterior over all the model parameters (reward function and transition model), sampling an MDP from this belief state at the start

Algorithm 21: Posterior sampling RL. We define the history at step k to be the set of previous trajectories, $H_k = \{\tau_1, \dots, \tau_{k-1}\}$, each of length H , where $\tau_k = (s_1^k, a_1^k, r_1^k, \dots, s_H^k, a_H^k, r_H^k, s_{H+1}^k)$.

```

1 Input: Prior over models  $P(M)$ 
2 History  $H_1 = \emptyset$ 
3 for Episode  $k = 1 : K$  do
4   Sample model from posterior,  $M_k \sim P(M|H_k)$ 
5   Compute optimal policy  $\pi_k^* = \text{solve}(M_k)$ 
6   Execute  $\pi_k^*$  for  $H$  steps to get  $\tau_k$ 
7   Update history  $H_{k+1} = H_k \cup \tau_k$ 
8   Update posterior  $p(M|H_{k+1})$ 
9 Return  $\pi_K^*$ 

```

of each episode, solving for the optimal policy corresponding to the sampled MDP, using the resulting policy to collect new data, and then updating the belief state at the end of the episode. This is called **posterior sampling RL** [Str00; ORVR13; RR14; OVR17; WCM24]. See Algorithm 21 for the pseudocode.³

As a more computationally efficient alternative, it is also possible to maintain a posterior over policies or Q functions instead of over world models; see e.g., [Osb+16] for a simple implementation of this idea using **bootstreap DQN**, [Osb+23a] for an implementation based on **epistemic neural networks** [Osb+23b], and **epistemic value estimation** [SSTVH23] for an implementation based on Laplace approximation. Another approach is to use successor features (Section 4.5.3), where the Q function is assumed to have the form $Q^\pi(s, a) = \psi^\pi(s, a)^\top \mathbf{w}$. In particular, [Jan+19b] proposes **Successor Uncertainties**, in which they model the uncertainty over \mathbf{w} as a Gaussian, $p(\mathbf{w}) = \mathcal{N}(\mu_{\mathbf{w}}, \Sigma_{\mathbf{w}})$. From this they can derive the posterior distribution over Q values as $p(Q(s, a)) = \mathcal{N}(\Psi^\pi \mu_{\mathbf{w}}, \Psi^\pi \Sigma_{\mathbf{w}} (\Psi^\pi)^\top)$, where $\Psi^\pi = [\psi^\pi(s, a)]^\top$ is a matrix of features, one per state-action pair.

7.2.3 Upper confidence bounds (UCBs)

The optimal solution to explore-exploit is intractable. However, an intuitively sensible approach is based on the principle known as “**optimism in the face of uncertainty**” (OFU). The principle selects actions greedily, but based on optimistic estimates of their rewards. This approach is optimal in the regret minimization sense, as proved in the **R-Max** paper of [Ten02], which builds on the earlier **E3** paper of [KS02].

The most common implementation of this principle is based on the notion of an **upper confidence bound** or **UCB**. We will initially explain this for the bandit case, then extend to the MDP case.

7.2.3.1 Basic idea

To use a UCB strategy, the agent maintains an optimistic reward function estimate \tilde{R}_t , so that $\tilde{R}_t(s_t, a) \geq R(s_t, a)$ for all a with high probability, and then chooses the greedy action accordingly:

$$a_t = \underset{a}{\operatorname{argmax}} \tilde{R}_t(s_t, a) \quad (7.23)$$

UCB can be viewed a form of **exploration bonus**, where the optimistic estimate encourages exploration. Typically, the amount of optimism, $\tilde{R}_t - R$, decreases over time so that the agent gradually reduces exploration. With properly constructed optimistic reward estimates, the UCB strategy has been shown to achieve near-optimal regret in many variants of bandits [LS19]. (We discuss regret in Section 7.1.)

The optimistic function \tilde{R} can be obtained in different ways, sometimes in closed forms, as we discuss below.

³In [AG25], they used prompted LLMs to implement a crude approximation to the sampling, planning, and posterior updating steps of the PSRL algorithm for some simple tabular problems. Although the method worked surprisingly well (in the sense of having low Bayesian regret) on very small problems (e.g., 3 state RiverSwim), it failed on larger problems (e.g., 4 state) and more stochastic problems.

7.2.3.2 Bandit case: Frequentist approach

A frequentist approach to computing a confidence bound can be based on a **concentration inequality** [BLM16] to derive a high-probability upper bound of the estimation error: $|\hat{R}_t(s, a) - R_t(s, a)| \leq \delta_t(s, a)$, where \hat{R}_t is a usual estimate of R (often the MLE), and δ_t is a properly selected function. An optimistic reward is then obtained by setting $\tilde{R}_t(s, a) = \hat{R}_t(s, a) + \delta_t(s, a)$.

As an example, consider again the context-free Bernoulli bandit, $R(a) \sim \text{Ber}(\mu(a))$. The MLE $\hat{R}_t(a) = \hat{\mu}_t(a)$ is given by the empirical average of observed rewards whenever action a was taken:

$$\hat{\mu}_t(a) = \frac{N_t^1(a)}{N_t(a)} = \frac{N_t^1(a)}{N_t^0(a) + N_t^1(a)} \quad (7.24)$$

where $N_t^r(a)$ is the number of times (up to step $t - 1$) that action a has been tried and the observed reward was r , and $N_t(a)$ is the total number of times action a has been tried:

$$N_t(a) = \sum_{s=1}^{t-1} \mathbb{I}(a_s = a) \quad (7.25)$$

Then the **Chernoff-Hoeffding inequality** [BLM16] leads to $\delta_t(a) = c/\sqrt{N_t(a)}$ for some constant c , so

$$\tilde{R}_t(a) = \hat{\mu}_t(a) + \frac{c}{\sqrt{N_t(a)}} \quad (7.26)$$

7.2.3.3 Bandit case: Bayesian approach

We can also derive an upper confidence about using Bayesian inference. If we use a beta prior, we can compute the posterior in closed form, as shown in Equation (1.26). The posterior mean is $\hat{\mu}_t(a) = \mathbb{E}[\mu(a)|\mathbf{h}_t] = \frac{\alpha_t^a}{\alpha_t^a + \beta_t^a}$, and the posterior standard deviation is approximately

$$\hat{\sigma}_t(a) = \sqrt{\mathbb{V}[\mu(a)|\mathbf{h}_t]} \approx \sqrt{\frac{\hat{\mu}_t(a)(1 - \hat{\mu}_t(a))}{N_t(a)}} \quad (7.27)$$

We can use similar techniques for a Gaussian bandit, where $p_R(R|a, \boldsymbol{\theta}) = \mathcal{N}(R|\mu_a, \sigma_a^2)$, μ_a is the expected reward, and σ_a^2 the variance. If we use a conjugate prior, we can compute $p(\mu_a, \sigma_a|\mathcal{D}_t)$ in closed form. Using an uninformative version of the conjugate prior, we find $\mathbb{E}[\mu_a|\mathbf{h}_t] = \hat{\mu}_t(a)$, which is just the empirical mean of rewards for action a . The uncertainty in this estimate is the standard error of the mean, i.e., $\sqrt{\mathbb{V}[\mu_a|\mathbf{h}_t]} = \hat{\sigma}_t(a)/\sqrt{N_t(a)}$, where $\hat{\sigma}_t(a)$ is the empirical standard deviation of the rewards for action a .

Once we have computed the mean and posterior standard deviation, we define the optimistic reward estimate as

$$\tilde{R}_t(a) = \hat{\mu}_t(a) + c\hat{\sigma}_t(a) \quad (7.28)$$

for some constant c that controls how greedy the policy is. See Figure 7.2 for an illustration. We see that this is similar to the frequentist method based on concentration inequalities, but is more general.

7.2.3.4 MDP case

The UCB idea (especially in its frequentist form) has been extended to the MDP case in several works. (The Bayesian version is discussed in Section 7.2.2.) For example, [ACBF02] proposes to combine UCB with Q learning, by defining the policy as

$$\pi(a|s) = \mathbb{I}\left(a = \underset{a'}{\operatorname{argmax}} Q(s, a') + c\sqrt{\log(t)/N_t(s, a')}\right) \quad (7.29)$$

[AJO08] presents the more sophisticated **UCRL2** algorithm, which computes confidence intervals on all the MDP model parameters at the start of each episode; it then computes the resulting **optimistic MDP** and solves for the optimal policy, which it uses to collect more data.

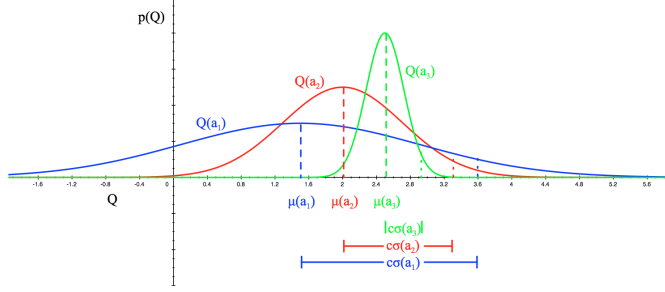


Figure 7.2: Illustration of the reward distribution $Q(a)$ for a Gaussian bandit with 3 different actions, and the corresponding lower and upper confidence bounds. We show the posterior means $Q(a) = \mu(a)$ with a vertical dotted line, and the scaled posterior standard deviations $co(a)$ as a horizontal solid line. From [Sil18]. Used with kind permission of David Silver.

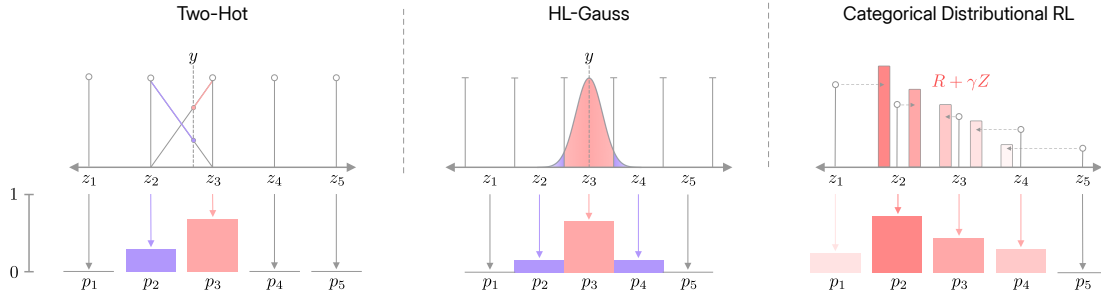


Figure 7.3: Illustration of how to encode a scalar target y or distributional target Z using a categorical distribution. From Figure 1 of [Far+24]. Used with kind permission of Jesse Farebrother.

7.3 Distributional RL

The **distributional RL** approach of [BDM17; BDR23], predicts the distribution of (discounted) returns, not just the expected return. More precisely, let $Z_t^\pi = \sum_{k=0}^{T-t} \gamma^k R(s_{t+k}, a_{t+k})$ be a random variable representing the (discounted) reward-to-go from step t . The standard value function is defined to compute the expectation of this variable: $V^\pi(s) = \mathbb{E}[Z_0^\pi | s_0 = s]$. In DRL, we instead attempt to learn the full distribution, $p(Z_0^\pi | s_0 = s)$ when training the critic. We then compute the expectation of this distribution when training the actor. For a general review of distributional regression, see [KSS23]. Below we briefly mention a few algorithms in this class that have been explored in the context of RL.

7.3.1 Quantile regression methods

An alternative to predicting a full distribution is to predict a fixed set of quantiles. This is called quantile regression, and has been used with DQN in [Dab+17] to get **QR-DQN**, and with SAC in [Wur+22] to get **QR-SAC**. (The latter was used in Sony’s **GTSophy** Gran Turismo AI racing agent.)

7.3.2 Replacing regression with classification

An alternative to quantile regression is to approximate the distribution over returns using a histogram, and then fit it using cross entropy loss (see Figure 7.3). This approach was first suggested in [BDM17], who called it **categorical DQN**. (In their paper, they use 51 discrete categories (atoms), giving rise to the name **C51**.)

An even simpler approach is to replace the distributional target with the standard scalar target (representing the mean), and then discretize this target and use cross entropy loss instead of squared error.⁴ Unfortunately,

⁴Technically speaking, this is no longer a distributional RL method, since the prediction target is the mean, but the mechanism

this encoding is lossy. In [Sch+20], they proposed the **two-hot** transform, that is a lossless encoding of the target based on putting appropriate weight on the nearest two bins (see Figure 7.3). In [IW18], they proposed the **HL-Gauss** histogram loss, that convolves the target value y with a Gaussian, and then discretizes the resulting continuous distribution. This is more symmetric than two-hot encoding, as shown in Figure 7.3. Regardless of how the discrete target is chosen, predictions are made using $\hat{y}(s; \theta) = \sum_k p_k(s) b_k$, where $p_k(s)$ is the probability of bin k , and b_k is the bin center.

In [Far+24], they show that the HL-Gauss trick works much better than MSE, two-hot and C51 across a variety of problems (both offline and online), especially when they scale to large networks. They conjecture that the reason it beats MSE is that cross entropy is more robust to noisy targets (e.g., due to stochasticity) and nonstationary targets. They also conjecture that the reason HL works better than two-hot is that HL is closer to ordinal regression, and reduces overfitting by having a softer (more entropic) target distribution (similar to label smoothing in classification problems).

Recently, [Gil+25] proposed an alternative called **Fourier head**, in which the linear output layer is replaced by a fourier transform, before discretizing. This is compatible with standard transformer training, and gives improved results when generating continuous outputs.

7.4 Intrinsic motivation for reward-free RL

When the extrinsic reward is sparse, or does not exist at all (“unsupervised RL”), it can be useful to reward the agent for solving “generally useful” tasks, such as learning about the world, or developing a set of skills. This is called **intrinsically motivated RL** [AMH19; Lin+19; Ami+21; Lad+22; Yua22; Col+22; AMH23; Yua+25].

We can classify these methods into two main types: (1) **knowledge-based intrinsic motivation**, or **artificial curiosity**, where the agent is rewarded for learning about its environment (this is focused on reducing prediction error); and (2) **competence-based intrinsic motivation**, where the agent is rewarded for achieving novel goals or mastering new skills (this is focused on control).

7.4.1 Knowledge-based intrinsic motivation

In this section, we discuss some approaches to knowledge-based intrinsic motivation.

7.4.1.1 Exploration bonuses

One simple approach is to create an intrinsic **exploration bonus** $R_t^i(s_t)$ which is high when the agent visits novel states. For tabular environments, we can just count the number of visits to each state, $N_t(s)$, and define $R_t^i(s) = 1/N_t(s)$ or $R_t^i(s) = 1/\sqrt{N_t(s)}$, which is similar to the UCB heuristic used in bandits (see Section 7.2.3). We can extend exploration bonuses to high dimensional states (e.g. images) using density models [Bel+16]. Alternatively, [MBB20] propose to use the ℓ_1 norm of the successor feature (Section 4.5.3) representation as an alternative to the visitation count, giving rise to an intrinsic reward of the form $R^i(s) = 1/||\psi^\pi(s)||_1$. Recently [Yu+23] extended this to combine SFs with *predecessor* representations, which encode retrospective information about the previous state (c.f., inverse dynamics models, mentioned below). This encourages exploration towards bottleneck states.

7.4.1.2 Random Network Distillation (RND)

Another approach is the **Random Network Distillation** or **RND** method of [Bur+18]. This uses a fixed random neural network feature extractor $z_t = f(s_t; \theta^*)$ to define a target, and then trains a predictor $\hat{z}_t = f(s_t; \hat{\theta}_t)$ to predict these targets. If s_t is similar to previously seen states, then the trained model will have low prediction error. We can thus define the intrinsic reward as proportional to the squared error $||\hat{z}_t - z_t||_2^2$. The **BYOL-Explore** method of [Guo+22] goes beyond RND by learning the target representation (for the next state), rather than using a fixed random projection, but is still based on prediction error.

for predicting the mean leverages a distribution, for robustness and ease of optimization.

7.4.1.3 Information-theoretic measures

We can also define an intrinsic reward in terms of the information theoretic **surprise** of the next state given the current one:

$$R(\mathbf{s}, \mathbf{a}, \mathbf{s}') = -\log q(\mathbf{s}'|\mathbf{s}, \mathbf{a}) \quad (7.30)$$

This is the same as methods based on rewarding states for prediction error. Unfortunately such methods can suffer from the **noisy TV problem** (also called a **stochastic trap**), in which an agent is attracted to states which are intrinsically to predict. To see this, note that by averaging over future states we see that the above reward reduces to

$$R(\mathbf{s}, \mathbf{a}) = -\mathbb{E}_{p^*(\mathbf{s}'|\mathbf{s}, \mathbf{a})} [\log q(\mathbf{s}'|\mathbf{s}, \mathbf{a})] = \mathbb{H}_{ce}(p^*, q) \quad (7.31)$$

where p^* is the true model and q is the learned dynamics model, and \mathbb{H}_{ce} is the cross -entropy. As we learn the optimal model, $q = p^*$, this reduces to the conditional entropy of the predictive distribution, which can be non-zero for inherently unpredictable states.

To help filter out such random noise, [Pat+17] proposes an **Intrinsic Curiosity Module**. This first learns an **inverse dynamics model** of the form $a = f(\mathbf{s}, \mathbf{s}')$, which tries to predict which action was used, given that the agent was in \mathbf{s} and is now in \mathbf{s}' . The classifier has the form $\text{softmax}(g(\phi(\mathbf{s}), \phi(\mathbf{s}'), a))$, where $\mathbf{z} = \phi(\mathbf{s})$ is a representation function that focuses on parts of the state that the agent can control. Then the agent learns a forwards dynamics model in \mathbf{z} -space. Finally it defines the intrinsic reward as

$$R(\mathbf{s}, \mathbf{a}, \mathbf{s}') = -\log q(\phi(\mathbf{s}')|\phi(\mathbf{s}), a) \quad (7.32)$$

Thus the agent is rewarded for visiting states that lead to unpredictable consequences, where the difference in outcomes is measured in a (hopefully more meaningful) latent space.

Another solution is to replace the cross entropy with the KL divergence, $R(\mathbf{s}, \mathbf{a}) = D_{\text{KL}}(p||q) = \mathbb{H}_{ce}(p, q) - \mathbb{H}(p)$, which goes to zero once the learned model matches the true model, even for unpredictable states. This has the desired effect of encouraging exploration towards states which have epistemic uncertainty (reducible noise) but not aleatoric uncertainty (irreducible noise) [MP+22]. The **BYOL-Hindsight** method of [Jar+23] is one recent approach that attempts to use the $R(\mathbf{s}, \mathbf{a}) = D_{\text{KL}}(p||q)$ objective. Unfortunately, computing the $D_{\text{KL}}(p||q)$ term is much harder than the usual variational objective of $D_{\text{KL}}(q||p)$. A related idea, proposed in the RL context by [Sch10], is to use the **information gain** as a reward. This is defined as $R_t(\mathbf{s}_t, \mathbf{a}_t) = D_{\text{KL}}(q(\mathbf{s}_t|\mathbf{h}_t, \mathbf{a}_t, \boldsymbol{\theta}_t)||q(\mathbf{s}_t|\mathbf{h}_t, \mathbf{a}_t, \boldsymbol{\theta}_{t-1}))$, where \mathbf{h}_t is the history of past observations, and $\boldsymbol{\theta}_t = \text{update}(\boldsymbol{\theta}_{t-1}, \mathbf{h}_t, \mathbf{a}_t, \mathbf{s}_t)$ are the new model parameters. This is closely related to the BALD (Bayesian Active Learning by Disagreement) criterion [Hou+11; KAG19], and has the advantage of being easier to compute, since it does not reference the true distribution p .

7.4.2 Competence-based intrinsic motivation

Another way to explore the environment is to use goal-conditioned RL (see Section 1.2.3), where the agent creates its own goals; this is known as an **autotelic agent** [Col+22]. Intuitively it is desirable to choose goals that cover the set of states. We discuss some ways of choosing these goals below.

7.4.2.1 Empowerment

One approach to choosing goals is to use the concept of **empowerment** [KPN05], which is defined as the mutual information between the goal G (or past action A) and the future state S . This is given by

$$I(G, S) = H(S) - H(S|G) \quad (7.33)$$

Thus we can maximize empowerment by maximizing the entropy of the states (a form of diversity) while minimizing the conditional entropy of the states given the goal (ensuring that the goal is predictable in its effects). For more details, see e.g., [GRW17; Eys+19; Cho+21; Tio+24], and for a cognitive science perspective, see [Sch12; Du+23b; GG24; Yiu+25].

7.4.2.2 Curriculum design

Since the space of possible goals is usually too vast to explore, it is important to choose useful goals for the agent to learn from. A good goal is often defined as one that is not too hard or too easy to learn, since this maximizes **learning progress**, also called the “**zone of proximal development**”. Choosing the best order in which to tackle various goals is an example of **automatic curriculum** design; similar methods can also be used to automatically design new environments (see e.g., [Por+20; For+22] for a review), a process which is sometimes called **open-ended learning**. It is also possible to train one agent to design an environment that another other agent finds challenging to solve; this is known as **asymmetric self-play** (see Section 7.5.3.1). The convergence of such setter-solver systems can be studied using the tools of Stackelberg equilibria (Section 5.2.13).

7.4.2.3 Using an LLM to choose goals

Another approach to goal generation is to use suitably prompted LLMs (see e.g., [Du+23c; Col+23]). This can leverage the prior knowledge of LLMs to not only propose novel goals, but also ones that are plausibly useful to humans.

7.4.2.4 Go-Explore

The **Go-Explore** algorithm from [Eco+19; Eco+21] proposes to first follow a goal-conditioned policy to reach (or to reset the environment state to) a goal state, which is chosen from an archive of “interestingly new” previously visited states (e.g., rarely visited ones), and then switch to an exploration policy (using random actions) to expand the coverage of the state space.

In [LHC25], they present **Intelligent Go-Explore**, which uses an LLM to decide what is an interesting goal to return to (leveraging the LLM’s ability to determine interestingness [Zha+24b]), to decide what exploratory actions to take after reaching the goal, and to decide whether to add any newly explored states to the archive.

In [Hu+23b], they choose a goal g for Go-Explore based on the expected value of the state they would end up in, if they followed the goal-conditioned policy towards g and then switched to exploration mode. This expected value is computed using a learned world model, following the LEXA method of [Men+21] (which in turn is based on Dreamer, discussed in Section 4.4.1.3).

7.5 Hierarchical RL

So far we have focused on MDPs that work at a single time scale. However, this is very limiting. For example, imagine planning a trip from San Francisco to New York: we need to choose high level actions first, such as which airline to fly, and then medium level actions, such as how to get to the airport, followed by low level actions, such as motor commands. Thus we need to consider actions that operate multiple levels of **temporal abstraction**. Fitting policies that use temporal (and optionally also state) abstraction is called **hierarchical RL** or **HRL**. This is a big and important topic, and we only brief mention a few key ideas and methods. Our summary is based on [Pat+22], but see [Kli+25a] for a more recent survey. (See also Section 4.5 where we discuss multi-step predictive models; by contrast, in this section we focus on model-free methods.)

7.5.1 HRL using Options

7.5.1.1 Introduction

Before we define the HRL problem, we define the concept of an “option” [SPS99], which is a form of temporally extended action. Formally, an **option** $\omega = (I, \pi, \beta)$ is a tuple consisting of the following:

- The **subpolicy** (aka **intra-option policy**, or **action policy**) $\pi_\omega(a|s) \in [0, 1]$.

- The **termination probability** $\beta_\omega(s) \in [0, 1]$, which gives the probability of finishing in state s . This induces a geometric distribution over option durations, which we denote by $\tau \sim \beta_\omega$.
- The **initiation set** $I_\omega \subset S$, which is the set of states this option can start from. Alternatively, we can define $I_\omega(s) \in [0, 1]$ as the probability that ω can be started from s and achieve its goal. (Note that the **affordances** of a state, $A(s) = \{\omega : I_\omega(s) > \epsilon\}$, is the set of options that can be initiated in a given state in order to achieve a specific intent or goal [Khe+20].)

We may also optionally associate a reward function with each option, denoted $R_\omega(s)$. The set of all options is denoted Ω , and the distribution over options (i.e., the **high level policy**) is denoted by $\mu(o|s)$.

Executing an option at step t entails choosing an action using $a_t = \pi_\omega(s_t)$ and then deciding whether to terminate at step $t + 1$ with probability $1 - \beta_\omega(s_{t+1})$ or to continue following the option at step $t + 1$. (This is an example of a **semi-Markov decision process** [Put94].) If we define $\pi_\omega(s) = a$ and $\beta_\omega(s) = 0$ for all s , then this option corresponds to primitive action a that terminates in one step. But with options we can expand the repertoire of actions to include those that take many steps to finish.

Note that goal-conditioned RL (Section 1.2.3) can be considered a special case of options where each option corresponds to a different goal. The reward function for each option has the form $R_\omega(s) = \text{sim}(s, \omega)$, the termination function is $\beta_\omega(s) = \text{sim}(s, \omega) > \text{thresh}$, and the initiation set is the entire state space, so $I_\omega(s) = 1$ for all s . (An option that can be initiated from any state is called a **global option**.)

To create a semi-MDP with options, we need to define the reward function and dynamics model. The reward is defined as follows:

$$R(s, \omega) = \mathbb{E} [R_1 + \gamma R^2 + \dots + \gamma^{\tau-1} R_\tau | S_0 = s, A_{0:\tau-1} \sim \pi_\omega, \tau \sim \beta_\omega] \quad (7.34)$$

The dynamics model is defined as follows:

$$T_\gamma(s'|s, \omega) = \sum_{k=1}^{\infty} \gamma^k \Pr(S_k = s', \beta_\omega(s_k) | S_0 = s, A_{0:k-1} \sim \pi_\omega) \quad (7.35)$$

Note that $p_\gamma(s'|s, \omega)$ is not a conditional probability distribution, because of the γ^k term, but we can usually treat it like one. Note also that a dynamics model that can predict multiple steps ahead is sometimes called a **jumpy model** (see also Section 4.5.4.2).

We can use these definitions to define the value function for a hierarchical policy using a generalized Bellman equation, as follows:

$$V_\pi(s) = \sum_{\omega \in \Omega(s)} \pi(\omega|s) \left[R(s, \omega) + \sum_{s'} T_\gamma(s'|s, \omega) V_\pi(s') \right] \quad (7.36)$$

We can compute this using value iteration. We can then learn a policy using policy iteration, or a policy gradient method. In other words, once we have defined the options, we can use all the standard RL machinery.

If we have a set of options, learning the corresponding semi-MDP is harder than learning a standard MDP. We can sidestep this by using the **Expected-Length Model** (ELM) of options proposed in [Abe+19], as opposed to the standard Multi-Time Model (MLM). Let $\tilde{\tau}_\omega$ be the average number of timesteps needed to execute option ω . Then we define the dynamics model as

$$T_\gamma(s'|s, \omega) = \gamma^{\tilde{\tau}_\omega} T(s'|s, \omega) \quad (7.37)$$

where $T(s'|s, \omega)$ is the distribution over next states s' after ω has finished running, when started in s . In [Abe+19], they show that this is a (biased) estimator of the MTM transition model, but which is easier to learn. In particular, we can learn this model from trajectory data which records (s, s', τ, r) tuples, by maximizing the log likelihood of transitions $T(s'|s, \omega)$ and fitting a regression model to $\tilde{\tau}$ and R .

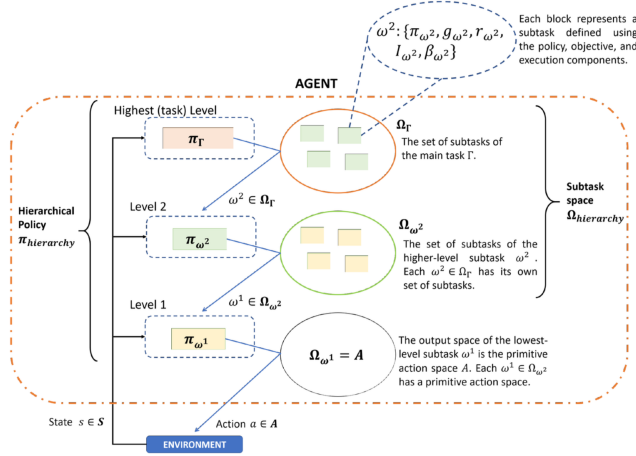


Figure 7.4: Illustration of a 3 level hierarchical RL system. From Fig 2 of [Pat+22].

7.5.1.2 Option hierarchies

We can now define the HRL problem in terms of a nested set of options. Let Ω_l be the set of options or subtasks at level l of the hierarchy, where $\Omega_1 = \mathcal{A}$ is the set of primitive actions. Let Γ denote the top of the hierarchy, corresponding to the main task. The policy at level l is denoted $\pi_l(\omega^l|s)$, where $\omega^l \in \Omega_l$ is the chosen subtask. The **subtask space** is $\Omega_{\text{hier}} = \{\Omega_2, \Omega_3, \dots, \Omega_\Gamma\}$. The **hierarchical policy** is $\pi_{\text{hier}} = \pi_1 \odot \dots \odot \pi_\Gamma$, which maps states to primitive actions, by successively invoking policies at descending levels of the hierarchy. See Figure 7.4 for an illustration.

Finally we can define HRL to be the following optimization problem:

$$\Omega_{\text{hier}}^*, \pi_{\text{hier}}^* = \underset{\Omega_{\text{hier}}}{\operatorname{argmax}} \underset{\pi_{\text{hier}}|\Omega_{\text{hier}}}{\operatorname{argmax}} Q^{\text{hier}}(s, a) \quad (7.38)$$

$$Q^{\text{hier}}(s_t, a_t) = \mathbb{E}_{a \sim \pi_{\text{hier}}|\Omega_{\text{hier}}} \left[\sum_{i=0}^{\infty} \gamma^{t+i} R(s_{t+1}, a_{t+i}) | s_t, a_t \right] \quad (7.39)$$

We see the problem breaks down into two parts: learning the hierarchical policy π_{hier} given a fixed hierarchy of tasks, Ω_{hier} , and learning the hierarchy itself. The latter is called the **subtask discovery** problem, and can be tackled jointly with policy learning, or can be done in a pre-training phase, as we discuss below.

7.5.1.3 Hierarchical Q learning

Suppose we have a set of options \mathcal{O} , with fixed initiation and termination functions. We can learn the corresponding option policy using Q-learning, as follows. Suppose we have a trajectory of $(s_t, o_t, r_{t:t+\tau}, s_{t+\tau})$ tuples, where τ is the duration of the option o_t . We then fit $Q(s, o)$ by regressing towards the following target:

$$y_t = \sum_{t'=t}^{\tau} \gamma^{t'-t} r_{t'} + \gamma^{\tau-t} Q(s_{t+\tau}, o_{t+\tau}^*) \quad (7.40)$$

$$o_{t+\tau}^* = \underset{o' \in \mathcal{O}'(s_{t+\tau})}{\operatorname{argmax}} Q(s_{t+\tau}, o') \quad (7.41)$$

$$\mathcal{O}'(s) = \{o_i : I_{o_i}(s) = 1 \cap \beta_{o_i}(s) = 0, \forall o_i \in \mathcal{O}\} \quad (7.42)$$

The individual option policies can be trained by 1-step Q learning in the usual way.

7.5.1.4 MAXQ

The Q-value function of any an option represents the expected cumulative reward until the end of the main task, rather than only until the horizon (or termination) of the option itself. Hence, an option is not a standalone subtask unit; instead it blends into the **core MDP**. This limits the transferability of the learned options to other tasks, but theoretically guarantees the optimality of the learned hierarchical policy.

In the **MAXQ** approach of [Die00], the core MDP is decomposed into smaller sub-MDP components. Each sub-MDP is associated to a subtask whose policy can be learned separately from other subtasks. This is achieved by decomposing the main Q-value function into the separate Q-value functions of the subtasks. Hence the Q-value of any subtask represents the expected cumulative reward only until the horizon (or termination) of that subtask. This allows each subtask policy to be learned as a standalone unit. However, the resulting policy will only be recursively optimal, rather than globally optimal.

7.5.1.5 Option learning using EM

If the set of options is unknown, we can learn them by segmenting the trajectories into sub-trajectories, which correspond to the latent options. This can be done using the EM algorithm [Dan+16].

7.5.1.6 Skill chaining

In [KB09; BK20], they propose the **skill chaining** method. The process begins by establishing an initial skill, let's call it ω , with the primary goal of the overall task as its objective. The agent then learns the necessary actions (the policy, denoted as π_ω) to achieve this main goal.

Next, a system is trained to identify the starting points, or “initiation states”, for this skill. States from which the main goal can be successfully reached within a set timeframe are marked as “positive”. Conversely, states where the goal is not reached within this time limit are labeled “negative.”. These successful starting points are then grouped together to form the initiation set for skill ω , designated as I_ω .

Following this, a new skill, ω' , is created. The objective of this new skill is to reach the initiation set (I_ω) of the first skill. The same process of learning the policy and identifying the initiation set is then repeated for ω' . This creates a chain, where each new skill's purpose is to get the agent to a state where the previous skill can be activated. This method allows the agent to autonomously discover and build a variable and adaptable number of interconnected skills. However, it only works for tasks that have explicit goal states, and requires strong exploration methods, without using options, to find out how to reach the goal states, which is particularly problematic for the initial goal. (One practical approach, used in the **DemoStart** system of [Bau+24], is to rely on a human demonstration that reaches the goal, and then initialize the agent in a state that is near the end of this trajectory.)

7.5.1.7 Option critic

We can also learn options using end-to-end training, i.e., the options and their policies are jointly learned online. For example, [BHP17] propose the **option-critic** architecture. The number of options is manually specified, and all policies are randomly initialized. Then they are jointly trained using policy gradient methods designed for semi-MDPs. (See also [RLT18] for a hierarchical extension of option-critic to support options calling options.) However, since the learning signal is just the main task reward, the method can work poorly in problems with sparse reward compared to subgoal methods (see discussion in [Vez+17; Nac+19]).

7.5.1.8 Double actor critic (DAC)

One problem with option-critic is that it requires specialized methods that are designed for optimizing semi-MDPs. In [ZW19], they propose **double actor critic**, which allows the use of standard policy gradient methods. This works by defining two parallel **augmented MDPs**, where the state space of each MDP is the cross-product of the original state space and the set of options. For a two level hierarchy, the manager learns a policy over options, and the worker learns a policy over states for each option. Both MDPs just use task rewards, without subgoals or subtask rewards.

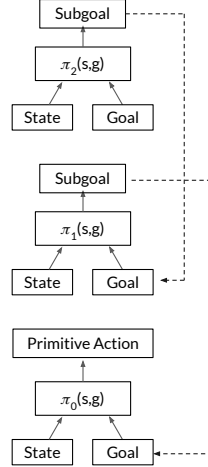


Figure 7.5: Illustration of a 3 level hierarchical goal-conditioned controller. From <http://bigai.cs.brown.edu/2019/09/03/hac.html>. Used with kind permission of Andrew Levy.

7.5.1.9 Avoiding excessive (or insufficient) option switching

It has been observed that option learning using option-critic or double actor-critic can fail, in the sense that the top level controller may learn to switch from one option to the next at almost every time step [ZW19; Har+18]. The reason is that the optimal policy does not require the use of temporally extended options, but instead can be defined in terms of primitive actions (as in standard RL). Therefore in [Har+18] they propose to add a regularizer called the **deliberation cost**, in which the higher level policy is penalized whenever it switches options. This can speed up learning, at the cost of a potentially suboptimal policy.

Another possible failure mode in option learning is if the higher level policy selects a single option for the entire task duration. To combat this, [KP19] propose the **Interest Option Critic**, which learns the initiation condition I_ω so that the option is selected only in certain states of interest, rather than the entire state space.

7.5.1.10 MBRL using options

Most work on options is based on model-free RL. In [Sut+23], they discuss the model-based RL case.

7.5.2 HRL using feudal hierarchies

7.5.2.1 Introduction

We now discuss the other main framework for HRL, known as **feudal RL** [DH92]. In this approach, the policy at level l (known as a **manager**) chooses a goal from some goal space (equal to the state space, or some abstraction thereof), and passes that down to the level below (known as a **worker**), as shown in Figure 7.5. Thus rather than having a finite number of options to choose from, we can have a nested set of parameterized (universal) policies, $\pi_l(g_{l-1}|s, g_l)$, for each manager level l , and $\pi_1(a|s, g_1)$ for the worker level. The value for a policy at a given level is the expected reward until the policy finishes, where the reward is defined intrinsically in terms of reaching the specified goal. Thus only the top level manager gets to see the external (environment) reward, a principle known as **reward hiding**.

7.5.2.2 Comparison with options

Although the feudal approach is somewhat easier to learn (due to the locality/modularity of subgoals, as we will see), the resulting hierarchical policy may be suboptimal compared to the optimal flat policy, since

learning is performed at each level wrt local goals, rather than optimizing a single global objective. In addition, not all subroutines or skills can be defined in terms of reaching a goal state (e.g., consider the skill of “driving in a circle”). By contrast, option-based HRL can match the optimality of a flat policy, and can have options with richer termination conditions. See Table 7.2 for a further comparison.

Feature	Options Framework	Feudal Reinforcement Learning
Hierarchical Structure	Flat controller + multiple sub-policies (options)	Manager \rightarrow Worker (two-level policy hierarchy)
Control Flow	Top-level policy chooses discrete option	Manager emits subgoals; worker acts toward them
Subpolicy Input	Option policies take full state or goal	Worker is conditioned on a goal from the manager
Termination Handling	Explicit termination function $\beta(s)$	Often fixed horizon or implicit (e.g., N steps)
Initiation Sets	Explicit $\mathcal{I}_o \subseteq \mathcal{S}$	Often ignored or assumed full initiation
Goal Communication	Top policy selects an option (index)	Manager gives a vectorial subgoal (e.g., in state space)
Training Paradigm	Semi-MDP; option-critic and variants	Two-agent structure; manager learns to guide worker
Interpretability	Each option has a nameable skill	Subgoals may be more abstract (e.g., latent vectors)

Table 7.2: Comparison between the Options Framework and Feudal Reinforcement Learning

7.5.2.3 Feudal Q learning

In [Kul+16], they tackle the feudal learning problem using two layers of DQN networks. The manager has to choose from a finite number of predefined subgoals, and is trained to maximize the task reward. The worker receives the goal, and has to choose from a finite number of discrete primitive actions, and is trained to maximize the (possibly sparse) intrinsic reward, based on reaching the goal.

7.5.2.4 Dealing with nonstationarity using hindsight relabeling (HIRO, HAC)

A major difficulty in HRL, which arises when training multiple levels of policies simultaneously, is that the resulting data distribution is **non-stationary**. To see this, note that the worker generates a distribution over trajectories in response to a (state, subgoal) pair that, but this distribution changes as the worker’s policy improves. Thus the data that the manager sees is also coming from a non-stationary distribution, making it hard to disentangle any changes due to the managers choice of actions from changes due to the worker’s policy changing.

In [Nac+18] they propose **HIRO** (Hierarchical Reinforcement Learning with Off-policy Correction), which is a way to tackle this non-stationarity by using hindsight relabeling (see Section 2.5.5). The data for the manager are transition tuples of the form $(s_t, g_t, \sum r_{t:t+c}, s_{t+c})$, where c is the time taken for the worker to reach the goal (or some maximum time), and r_t is the main task reward function at step t . The data for the worker are transition tuples of the form $(s_{t+i}, g_t, a_{t+i}, r_{t+i}^{g_t}, s_{t+i+1})$ for $i = 0 : c$, where r_t^g is the reward wrt reaching goal g . This data can be used to train the two policies. However, if the worker fails to achieve the goal in the given time limit, all the rewards will be 0, and no learning will take place. To combat this, if the worker does not achieve g_t after c timesteps, the subgoal is relabeled in the transition data with another

subgoal g'_t which is sampled from $p(g|\tau)$, where τ is the observed trajectory. Thus both policies treat g'_t as the goal in hindsight, so they can use the actually collected data for training

The **hierarchical actor critic** (HAC) method of [Lev+18] is a simpler version of HIRO that can be extended to multiple levels of hierarchy, where the lowest level corresponds to primitive actions (see Figure 7.4). In the HAC approach, the output subgoal in the higher level data, and the input subgoal in the lower-level data, are replaced with the actual state that was achieved in hindsight. This allows the training of each level of the hierarchy independently of the lower levels, by assuming the lower level policies are already optimal (since they achieved the specified goal). As a result, the distribution of (s, a, s') tuples experienced by a higher level will be stable, providing a stationary learning target.

7.5.2.5 Learning the goal space and policy

In the previous approaches, the subgoals are defined in terms of the states that were achieved at the end of each trajectory, $g' = s_T$. This can be generalized by using a state abstraction function to get $g' = \phi(s_T)$. The methods in Section 7.5.2.4 assumed that ϕ was manually specified. We now mention some ways to learn ϕ .

In [Vez+17], they present **Feudal Networks** for learning a two level hierarchy. The manager samples subgoals in a learned latent subgoal space. The worker uses distance to this subgoal as a reward, and is trained in the usual way. The manager uses the “transition gradient” as a reward, which is derived from the task reward as well as the distance between the subgoal and the actual state transition made by the worker. This reward signal is used to learn the manager policy and the latent subgoal space.

Feudal networks do not guarantee that the learned subgoal space will result in optimal behavior. In [Nac+19], they present a method to optimize the policy and ϕ function so as to minimize a bound on the suboptimality of the hierarchical policy. This approach is combined with HIRO (Section 7.5.2.4) to tackle the non-stationarity issue.

7.5.3 Subtask discovery

In this section, we discuss ways of learning hierarchical structure in an environment (given a dataset), independent of the specific task that needs to be solved. Thus this can be thought of as a “pre-training” phase. The resulting subtasks can then be used to define a hierarchical structure, for which a hierarchical policy can be trained, using the methods described above.

7.5.3.1 Discovery of subgoals

One way to define subtasks is in terms of subgoals that might be worth achieving. These are often chosen to be **bottleneck states**, through which many paths (in state space) must pass when going from different starting states to different goal states. These states are often identified by performing various graph-theoretic analyses of the graph G derived from the state transition diagram by defining an edge $E_{s,s'}$ iff $\sum_a p(s'|s, a) > 0$. For example, [MB01] proposes the **diverse density** metric, and [SB08] proposes a metric on betweenness centrality. (This quantifies how important a node is in a network by counting how many times it appears on the shortest path between other nodes.)

A notable drawback of existing subgoal discovery methods is their reliance on a discrete subgoal space. This discreteness prevents the interpolation of new subgoals, which in turn restricts the variety of tasks the agent can learn. To overcome this, the development of continuous subgoal spaces is crucial. To tackle this, [Suk+18] introduced **Hierarchical Self Play** (HSP), a method for learning a continuous embedding of subgoals through an unsupervised pre-training technique called **asymmetric self-play** (see Section 5.3.5). This process begins by initializing two reinforcement learning policies, designated as Alice (π_A) and Bob (π_B), also called a **setter-solver** pair. The asymmetric self-play unfolds as follows:

1. Starting from an initial state s_0 , Alice’s policy π_A acts for T_A steps, arriving at a final state denoted as $s^* = s_A^{T_A}$.
2. The environment is then reset to s_0 , and this state s^* is assigned as a target for Bob’s policy, π_B .

3. As Bob executes its policy, a learned encoder, E , generates a low-dimensional subgoal embedding at each timestep t : $g_t = E(s_t^B, s^*)$, where s_t^B is Bob’s state at that time.
4. Bob’s policy then selects actions based on its current state and this subgoal embedding: $a_t^B = \pi_B(s_t^B, g_t)$.
5. Bob has T_B steps (where T_B is typically close to T_A) to reach the target state s^* . A reward of $R_B = 1$ is given for success and 0 for failure.
6. Alice receives an opposing reward, $R_A = 1 - R_B$.

This setup creates a dynamic where Bob learns to reach goals set by Alice, while Alice is incentivized to discover novel states that are currently challenging for Bob. This encourages effective, unsupervised exploration of the environment.

The encoder, $E(s_t^B, s^*)$, effectively maps various target states into a continuous, low-dimensional subgoal space. Following this pre-training phase, this learned space is utilized as the continuous action space for the high-level policy for an HRL agent. The low-level policy of this agent is initialized with Bob’s learned policy and is subsequently fine-tuned for specific downstream tasks.

In [Ope+20], researchers at OpenAI used this method to train a robot by learning its own curriculum. And in [Zho+25b], researchers at FAIR use a related method to train an LLM to perform tasks.

Note that this kind of goal generation is related to goal-based intrinsic motivation (see Section 7.4.2).

7.5.3.2 Discovery of skills

Approaches centered on subgoal discovery are inadequate for identifying subtasks that do not have a specific, concrete objective. For instance, a subtask like "navigating through traffic", within the broader goal of reaching a location, involves complex maneuvers but lacks a distinct subgoal state, making it unsuitable for subgoal discovery techniques.

In this section, we focus on methods for learning a varied collection of **skills**, where we define a skill to be a policy for subtask, encapsulating the agent’s ability to perform a certain action. Once discovered, these skills can be integrated as the lower-level policies within an HRL agent and subsequently fine-tuned for a given task.

A prominent strategy for discovering diverse skills involves maximizing the Mutual Information (MI) between a given skill, ω , and the resulting states or trajectories produced when that skill is executed. This is typically implemented by conditioning a universal policy on a latent skill vector, z_ω , creating a skill-specific policy $\pi_\omega(s) = \pi(s, z_\omega)$. In many such methods, z_ω is a one-hot vector. By maximizing the MI between the skill vector z_ω and the outcomes of its corresponding policy $\pi(s, z_\omega)$, the agent is encouraged to learn a set of distinct skills that lead to different states or trajectories.

For example, the **Variational Intrinsic Control** (VIC) method of [GRW17] focuses on maximizing the mutual information between the skill vector z_ω and the terminal state s_T , given an initial state s_0 . The **Diversity Is All You Need** (DIAYN) method of [Eys+19] aims to discover skills by maximizing the mutual information between the skill vector z_ω and every state visited within the trajectory generated by $\pi(s, z_\omega)$, without considering the order of the states. (This approach is also called **Mutual Information Skill Learning** or MISL [Zhe+25a].)

The **VALOR** (Variational Autoencoding Learning of Options by Reinforcement) method of [Ach+18] uses a VAE-like method to learn skills from trajectories, respecting the order of the states. The **SeCTAR** (Self-Consistent Trajectory Autoencoder) method of [CR+18] also uses an autoencoder, but based on an LSTM encoder and decoder. They combine this with an exploration mechanism to generate diverse trajectories, so that the continuous latent space represents a diverse set of skills.

7.6 Imitation learning

In previous sections, an RL agent is to learn an optimal sequential decision making policy so that the total reward is maximized. **Imitation learning** (IL) (see e.g., [CC25]), also known as **apprenticeship**

learning and **learning from demonstration** (LfD), is a different setting, in which the agent does not observe rewards, but has access to a collection \mathcal{D}_{exp} of trajectories generated by an expert policy π_{exp} ; that is, $\tau = (s_0, a_0, s_1, a_1, \dots, s_T)$ and $a_t \sim \pi_{\text{exp}}(s_t)$ for $\tau \in \mathcal{D}_{\text{exp}}$. The goal is to learn a good policy by imitating the expert, in the absence of reward signals. IL finds many applications in scenarios where we have demonstrations of experts (often humans) but designing a good reward function is not easy, such as car driving and conversational systems. (See also Section 7.7, where we discuss the closely related topic of offline RL, where we also learn from a collection of trajectories, but no longer assume they are generated by an optimal policy.)

7.6.1 Imitation learning by behavior cloning

A natural method is **behavior cloning**, which reduces IL to supervised learning; see [Pom89] for an early application to autonomous driving. It interprets a policy as a classifier that maps states (inputs) to actions (labels), and finds a policy by minimizing the imitation error, such as

$$\min_{\pi} \mathbb{E}_{p_{\pi_{\text{exp}}}^{\gamma}(s)} [D_{\text{KL}}(\pi_{\text{exp}}(s) \parallel \pi(s))] \quad (7.43)$$

where the expectation wrt $p_{\pi_{\text{exp}}}^{\gamma}$ may be approximated by averaging over states in \mathcal{D}_{exp} . A challenge with this method is that the loss does not consider the sequential nature of IL: future state distribution is not fixed but instead depends on earlier actions. Therefore, if we learn a policy $\hat{\pi}$ that has a low imitation error under distribution $p_{\pi_{\text{exp}}}^{\gamma}$, as defined in Equation (7.43), it may still incur a large error under distribution $p_{\hat{\pi}}^{\gamma}$ (when the policy $\hat{\pi}$ is actually run). This problem has been tackled by the offline RL literature, which we discuss in Section 7.7.

7.6.2 Imitation learning by inverse reinforcement learning

An effective approach to IL is **inverse reinforcement learning** (IRL) or **inverse optimal control** (IOC). Here, we first infer a reward function that “explains” the observed expert trajectories, and then compute a (near-)optimal policy against this learned reward using any standard RL algorithms studied in earlier sections. The key step of reward learning (from expert trajectories) is the opposite of standard RL, thus called inverse RL [NR00].

It is clear that there are infinitely many reward functions for which the expert policy is optimal, for example by several optimality-preserving transformations [NHR99]. To address this challenge, we can follow the maximum entropy principle, and use an energy-based probability model to capture how expert trajectories are generated [Zie+08]:

$$p(\tau) \propto \exp\left(\sum_{t=0}^{T-1} R_{\theta}(s_t, a_t)\right) \quad (7.44)$$

where R_{θ} is an unknown reward function with parameter θ . Abusing notation slightly, we denote by $R_{\theta}(\tau) = \sum_{t=0}^{T-1} R_{\theta}(s_t, a_t)$ the cumulative reward along the trajectory τ . This model assigns exponentially small probabilities to trajectories with lower cumulative rewards. The partition function, $Z_{\theta} \triangleq \int_{\tau} \exp(R_{\theta}(\tau))$, is in general intractable to compute, and must be approximated. Here, we can take a sample-based approach. Let \mathcal{D}_{exp} and \mathcal{D} be the sets of trajectories generated by an expert, and by some known distribution q , respectively. We may infer θ by maximizing the likelihood, $p(\mathcal{D}_{\text{exp}}|\theta)$, or equivalently, minimizing the negative log-likelihood loss

$$\mathcal{L}(\theta) = -\frac{1}{|\mathcal{D}_{\text{exp}}|} \sum_{\tau \in \mathcal{D}_{\text{exp}}} R_{\theta}(\tau) + \log \frac{1}{|\mathcal{D}|} \sum_{\tau \in \mathcal{D}} \frac{\exp(R_{\theta}(\tau))}{q(\tau)} \quad (7.45)$$

The term inside the log of the loss is an importance sampling estimate of Z that is unbiased as long as $q(\tau) > 0$ for all τ . However, in order to reduce the variance, we can choose q adaptively as θ is being updated.

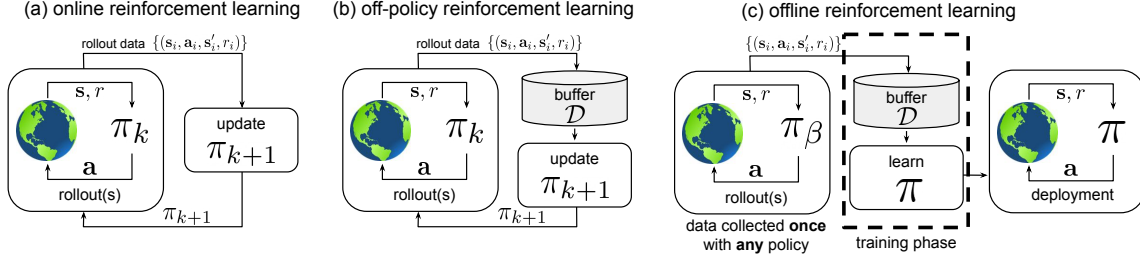


Figure 7.6: Comparison of online on-policy RL, online off-policy RL, and offline RL. From Figure 1 of [Lev+20a]. Used with kind permission of Sergey Levine.

The optimal sampling distribution, $q_*(\tau) \propto \exp(R_\theta(\tau))$, is hard to obtain. Instead, we may find a policy $\hat{\pi}$ which induces a distribution that is close to q_* , for instance, using methods of maximum entropy RL discussed in Section 3.6.4. Interestingly, the process above produces the inferred reward R_θ as well as an approximate optimal policy $\hat{\pi}$. This approach is used by **guided cost learning** [FLA16], and found effective in robotics applications.

7.6.3 Imitation learning by divergence minimization

We now discuss a different, but related, approach to IL. Recall that the reward function depends only on the state and action in an MDP. It implies that if we can find a policy π , so that $p_\pi^\gamma(s, a)$ and $p_{\pi_{\text{exp}}}^\gamma(s, a)$ are close, then π receives similar long-term reward as π_{exp} , and is a good imitation of π_{exp} in this regard. A number of IL algorithms find π by minimizing the divergence between p_π^γ and $p_{\pi_{\text{exp}}}^\gamma$. We will largely follow the exposition of [GZG19]; see [Ke+19] for a similar derivation.

Let f be a convex function, and D_f be the corresponding f -divergence [Mor63; AS66; Csi67; LV06; CS04]. From the above intuition, we want to minimize $D_f(p_{\pi_{\text{exp}}}^\gamma \| p_\pi^\gamma)$. Then, using a variational approximation of D_f [NWJ10], we can solve the following optimization problem for π :

$$\min_{\pi} \max_{\mathbf{w}} \mathbb{E}_{p_{\pi_{\text{exp}}}^\gamma(s, a)} [T_{\mathbf{w}}(s, a)] - \mathbb{E}_{p_\pi^\gamma(s, a)} [f^*(T_{\mathbf{w}}(s, a))] \quad (7.46)$$

where f^* is the convex conjugate of f , and $T_{\mathbf{w}} : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is some function parameterized by \mathbf{w} . We can think of π as a generator (of actions) and $T_{\mathbf{w}}$ as an adversarial critic that is used to compare the generated (s, a) pairs to the real ones. Thus the first expectation can be estimated using \mathcal{D}_{exp} , as in behavior cloning, and the second can be estimated using trajectories generated by policy π . Furthermore, to implement this algorithm, we often use a parametric policy representation π_θ , and then perform stochastic gradient updates to find a saddle-point to Equation (7.46). With different choices of the convex function f , we can obtain many existing IL algorithms, such as **generative adversarial imitation learning (GAIL)** [HE16] and **adversarial inverse RL (AIRL)** [FLL18], etc.

7.7 Offline RL

Offline reinforcement learning (also called **batch reinforcement learning** [LGR12]) is concerned with learning a reward maximizing policy from a fixed, static dataset, collected by some existing policy, known as the **behavior policy**. Thus no interaction with the environment is allowed (see Figure 7.6). This makes policy learning harder than the online case, since we do not know the consequences of actions that were not taken in a given state, and cannot test any such “counterfactual” predictions by trying them. (This is the same problem as in off-policy RL, which we discussed in Section 3.4.) In addition, the policy will be deployed on new states that it may not have seen, requiring that the policy generalize out-of-distribution, which is the main bottleneck for current offline RL methods [Par+24b].

A very simple and widely used offline RL method is known as behavior cloning or BC. This amounts to training a policy to predict the observed output action a_t associated with each observed state s_t , so we aim to ensure $\pi(s_t) \approx a_t$, as in supervised learning. This assumes the offline dataset was created by an expert, and so falls under the umbrella of imitation learning (see Section 7.6.1 for details). By contrast, offline RL methods can leverage suboptimal data. We give a brief summary of some of these methods below. For more details, see e.g., [Lev+20b; Che+24b; Cet+24; YWW25; Jac+25] and the list of papers at <https://github.com/hanjuku-kaso/awesome-offline-rl>. For some offline RL benchmarks, see DR4L [Fu+20], RL Unplugged [Gul+20], OGBench (Offline Goal-Conditioned benchmark) [Par+24a], and D5RL [Raf+24].

7.7.1 Offline model-free RL

In principle, we can tackle offline RL using the off-policy methods that we discussed in Section 3.4. These use some form of importance sampling, based on $\pi(a|s)/\pi_b(a|s)$, to reweight the data in the replay buffer \mathcal{D} , which was collected by the behavior policy, towards the current policy (the one being evaluated/learned). Unfortunately, such methods only work well if the behavior policy is close to the new policy. In the online RL case, this can be ensured by gradually updating the new policy away from the behavior policy, and then sampling new data from the updated policy (which becomes the new behavior policy). Unfortunately, this is not an option in the offline case. Thus we need to use other strategies to control the discrepancy between the behavior policy and learned policy, as we discuss below. (Besides the algorithmic techniques we discuss, another reliable way to get better offline RL performance is to train on larger, more diverse datasets, as shown in [Kum+23].)

7.7.1.1 Policy constraint methods

In the **policy constraint** method, we use a modified form of actor-critic, which, at iteration k , uses an update of the form

$$Q_{k+1}^\pi \leftarrow \operatorname{argmin}_Q \mathbb{E}_{(s,a,s') \sim \mathcal{D}} \left[\left(Q(s,a) - (R(s,a) + \gamma \mathbb{E}_{\pi_k(a'|s')} [Q_k^\pi(s',a')]) \right)^2 \right] \quad (7.47)$$

$$\pi_{k+1} \leftarrow \operatorname{argmax}_\pi \mathbb{E}_{s \sim \mathcal{D}} [\mathbb{E}_{\pi(a|s)} [Q_{k+1}^\pi(s,a)]] \quad \text{s.t.} \quad D(\pi, \pi_b) \leq \epsilon \quad (7.48)$$

where $D(\pi(\cdot|s), \pi_b(\cdot|s))$ is a divergence measure on distributions, such as KL divergence or another f -divergence. This ensures that we do not try to evaluate the Q function on actions a' that are too dissimilar from those seen in the data buffer (for each sampled state s), which might otherwise result in artefacts similar to an adversarial attack.

As an alternative to adding a constraint, we can add a penalty of $\alpha D(\pi(\cdot|s), \pi_b(\cdot|s))$ to the target Q value and the actor objective, resulting in the following update:

$$Q_{k+1}^\pi \leftarrow \operatorname{argmin}_Q \mathbb{E}_{(s,a,s') \sim \mathcal{D}} \left[\left(Q(s,a) - (R(s,a) + \gamma \mathbb{E}_{\pi_k(a'|s')} [Q_k^\pi(s',a') - \alpha D(\pi_k(\cdot|s'), \pi_b(\cdot|s'))]) \right)^2 \right] \quad (7.49)$$

$$\pi_{k+1} \leftarrow \operatorname{argmax}_\pi \mathbb{E}_{s \sim \mathcal{D}} [\mathbb{E}_{\pi(a|s)} [Q_{k+1}^\pi(s,a)] - \alpha D(\pi(\cdot|s'), \pi_b(\cdot|s'))] \quad (7.50)$$

One problem with the above method is that we have to fit a parametric model to $\pi_b(a|s)$ in order to evaluate the divergence term. Fortunately, in the case of KL, the divergence can be enforced implicitly, as in the **advantage weighted regression** or **AWR** method of [Pen+19], the **reward weighted regression** method of [PS07], the **advantage weighted actor critic** or **AWAC** method of [Nai+20], the **advantage weighted behavior model** or **ABM** method of [Sie+20]. In this approach, we first solve (nonparametrically) for the new policy under the KL divergence constraint to get $\bar{\pi}_{k+1}$, and then we project this into the required

policy function class via supervised regression, as follows:

$$\bar{\pi}_{k+1}(a|s) \leftarrow \frac{1}{Z} \pi_b(a|s) \exp \left(\frac{1}{\alpha} Q_k^\pi(s, a) \right) \quad (7.51)$$

$$\pi_{k+1} \leftarrow \underset{\pi}{\operatorname{argmin}} D_{\mathbb{KL}}(\bar{\pi}_{k+1} \parallel \pi) \quad (7.52)$$

In practice the first step can be implemented by weighting samples from $\pi_b(a|s)$ (i.e., from the data buffer) using importance weights given by $\exp(\frac{1}{\alpha} Q_k^\pi(s, a))$, and the second step can be implemented via supervised learning (i.e., maximum likelihood estimation) using these weights.

It is also possible to replace the KL divergence with an integral probability metric (IPM), such as the maximum mean discrepancy (MMD) distance, which can be computed from samples, without needing to fit a distribution $\pi_b(a|s)$. This approach is used in [Kum+19]. This has the advantage that it can constrain the support of the learned policy to be a subset of the behavior policy, rather than just remaining close to it. To see why this can be advantageous, consider the case where the behavior policy is uniform. In this case, constraining the learned policy to remain close (in KL divergence) to this distribution could result in suboptimal behavior, since the optimal policy may just want to put all its mass on a single action (for each state).

7.7.1.2 Behavior-constrained policy gradient methods

Recently a class of methods has been developed that is simple and effective: we first learn a baseline policy $\pi(a|s)$ (using BC) and a Q function (using Bellman minimization) on the offline data, and then update the policy parameters to pick actions that have high expected value according to Q and which are also likely under the BC prior. An early example of this is the Q^\dagger algorithm of [Fuj+19]. In [FG21], they present the **DDPG+BC** method, which optimizes

$$\max_{\pi} J(\pi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [Q(s, \mu^\pi(s)) + \alpha \log \pi(a|s)] \quad (7.53)$$

where $\mu^\pi(s) = \mathbb{E}_{\pi(a|s)}[a]$ is the mean of the predicted action, and α is a hyper-parameter. As another example, the **DQL** method of [WHZ23] optimizes a diffusion policy using

$$\min_{\pi} \mathcal{L}(\pi) = \mathcal{L}_{\text{diffusion}}(\pi) + \mathcal{L}_q(\pi) = \mathcal{L}_{\text{diffusion}}(\pi) - \alpha \mathbb{E}_{s \sim \mathcal{D}, a \sim \pi(\cdot|s)} [Q(s, a)] \quad (7.54)$$

where the second term is a penalty derived from Conservative Q Learning (Section 7.7.1.4), that ensures the Q values do not get too small. Finally, [Aga+22b] discusses how to transfer the policy from a previous agent to a new agent by combining BC with Q learning.

7.7.1.3 Uncertainty penalties

An alternative way to avoid picking out-of-distribution actions, where the Q function might be unreliable, is to add a penalty term to the Q function based on the estimated epistemic uncertainty, given the dataset \mathcal{D} , which we denote by $\text{Unc}(P_D(Q^\pi))$, where $P_D(Q^\pi)$ is the distribution over Q functions, and Unc is some metric on distributions. For example, we can use a deep ensemble to represent the distribution, and use the variance of $Q(s, a)$ across ensemble members as a measure of uncertainty. This gives rise to the following policy improvement update:

$$\pi_{k+1} \leftarrow \underset{\pi}{\operatorname{argmax}} \mathbb{E}_{s \sim \mathcal{D}} \left[\mathbb{E}_{\pi(a|s)} \left[\mathbb{E}_{P_D(Q_{k+1}^\pi)} [Q_{k+1}^\pi(s, a)] \right] - \alpha \text{Unc}(P_D(Q_{k+1}^\pi)) \right] \quad (7.55)$$

For examples of this approach, see e.g., [An+21; Wu+21; GGN22].

7.7.1.4 Conservative Q-learning

An alternative to explicitly estimating uncertainty is to add a **conservative penalty** directly to the Q -learning error term. That is, we minimize the following wrt \mathbf{w} using each batch of data \mathcal{B} :

$$\bar{\mathcal{E}}(\mathcal{B}, \mathbf{w}) = \alpha \mathcal{C}(\mathcal{B}, \mathbf{w}) + \mathcal{E}(\mathcal{B}, \mathbf{w}) \quad (7.56)$$

where $\mathcal{E}(\mathcal{B}, \mathbf{w}) = \mathbb{E}_{(s,a,s') \in \mathcal{B}} [(Q_{\mathbf{w}}(s,a) - (r + \gamma \max_{a'} Q_{\mathbf{w}}(s',a'))^2]$ is the usual loss for Q -learning, and $\mathcal{C}(\mathcal{B}, \mathbf{w})$ is some conservative penalty.

In the **conservative Q learning** or **CQL** method of [Kum+20], we use the following penalty term:

$$\mathcal{C}(\mathcal{B}, \mathbf{w}) = \mathbb{E}_{s \sim \mathcal{D}} [\mathbb{E}_{a \sim \mu(\cdot|s)} [Q_{\mathbf{w}}(s,a)] - \mathbb{E}_{a \sim \pi_b(\cdot|s)} [Q_{\mathbf{w}}(s,a)]] + R(\mu) \quad (7.57)$$

where μ is the new policy derived from Q , and $R(\mu) = -D_{\text{KL}}(\mu \parallel \rho)$ is a regularizer, and ρ is the action prior, which we discuss below. Since we are minimizing $\mathcal{C}(\mathcal{B}, \mathbf{w})$ (in addition to $\mathcal{E}(\mathcal{B}, \mathbf{w})$), we see that we are simultaneously maximizing the Q values for actions that are drawn from the behavior policy while minimizing the Q values for actions sampled from μ . This is to combat the optimism bias of Q -learning (hence the term “conservative”).

Now we derive the expression for μ . From Section 3.6.4 we know that the optimal solution has the form $\mu(a|s) = \frac{1}{Z} \rho(a|s) \exp(Q(s,a))$, where $Z = \sum_{a'} \exp(Q(s,a'))$ is the normalizer, and $\rho(a|s)$ is the prior. (For example, we can set $\rho(a|s)$ to be the previous policy.) We can then approximate the first term in the penalty using importance sampling, with $\rho(a|s)$ as the proposal:

$$\mathbb{E}_{a \sim \mu(\cdot|s)} [Q(s,a)] = \mathbb{E}_{\rho(a|s)} \left[\frac{\mu(a|s)}{\rho(a|s)} Q(s,a) \right] = \mathbb{E}_{\rho(a|s)} \left[\frac{\exp(Q(s,a))}{\sum_{a'} \exp(Q(s,a'))} Q(s,a) \right] \quad (7.58)$$

Alternatively, suppose we set $\rho(a|s)$ to be uniform, as in maxent RL (Section 3.6.4). In this case, we should replace the value function with the soft value function. From Equation (3.161), using a penalty coefficient of $\alpha = 1$, we have

$$\mathbb{E}_a [Q_{\text{soft}}(s,a)] = V_{\text{soft}}(s) = \log \sum_a \exp(Q(s,a)) \quad (7.59)$$

Note, however, this can be intractable for high-dimensional actions.

7.7.2 Offline model-based RL

In Chapter 4, we discussed model-based RL, which can train a dynamics model given a fixed dataset, and then use this to generate synthetic data to evaluate and then optimize different possible policies. However, if the model is wrong, the method may learn a suboptimal policy, as we discussed in Section 4.4.4. This problem is particularly severe in the offline RL case, since we cannot recover from any errors by collecting more data. Therefore various conservative MBRL algorithms have been developed, to avoid exploiting model errors. For example, [Kid+20] present the **MOREL** algorithm, and [Yu+20] present the **MOPO** algorithm. Unlike the value function uncertainty method of Section 7.7.1.3, or the conservative value function method of Section 7.7.1.4, these model-based methods add a penalty for visiting states where the model is likely to be incorrect.

In more detail, let $u(s,a)$ be an estimate of the uncertainty of the model’s predictions given input (s,a) . In MOPO, they define a conservative reward using $\bar{R}(s,a) = R(s,a) - \lambda u(s,a)$, and in MOREL, they modify the MDP so that the agent enters an absorbing state with a low reward when $u(s,a)$ is sufficiently large. In both cases, it is possible to prove that the model-based estimate of the policy’s performance under the modified reward or dynamics is a lower bound of the performance of the policy’s true performance in the real MDP, provided that the uncertainty function u is an error oracle, which means that it satisfies $D(M_{\theta}(s'|s,a), M^*(s'|s,a)) \leq u(s,a)$, where M^* is the true dynamics, and M_{θ} is the estimated dynamics.

For more information on offline MBRL methods, see [Che+24c].

7.7.3 Offline RL using reward-conditioned sequence modeling

Recently an approach to offline RL based on sequence modeling has become very popular. The basic idea — known as **upside down RL** [Sch19] or **RvS** (RL via Supervised learning) [KPL19; Emm+21] — is to train a generative model over future states and/or actions conditioned on the observed reward, rather than predicting the reward given a state-action trajectory. At test time, the conditioning is changed to represent the desired reward, and futures are sampled from the model. The implementation of this idea then depends on what kind of generative model to use, as we discuss below.

The **trajectory transformer** method of [JLL21] learns a joint model of the form $p(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}, \mathbf{r}_{1:T})$ using a transformer, and then samples from this using beam search, selecting the ones with high reward (similar to MPC, Section 4.2.4). The **decision transformer** [Che+21b] is related, but just generates action sequences, and conditions on the past observations and the future reward-to-go. That is, it fits

$$\operatorname{argmax}_{\theta} \mathbb{E}_{p_{\mathcal{D}}} [\log \pi_{\theta}(a_t | s_{0:t}, a_{0:t-1}, \text{RTG}_{0:t})] \quad (7.60)$$

where $\text{RTG}_t = \sum_{k=t}^T r_k$ is the return to go. This is just like BC policy learning, except we also condition on the RTG. At run time, RTG_0 is set to some desired high value (e.g., the maximum RTG observed during training), and is then updated online using $\text{RTG}_{t+1} = \text{RTG}_t - r_t$. To set more plausible RTG values, [Lee+22] propose to learn the distribution $p(\text{RTG}_t | s_{\leq t}, a_{\leq t}, \text{RTG}_{< t})$, which we can view as a critic, in addition to training the actor $p(a_t | s_{\leq t}, a_{< t}, \text{RTG}_{\leq t})$. [YKS23] propose the Q-learning Decision Transformer (**QDT**), which conditions on a Q value (learned using Q learning) instead of RTG. This combines the benefits of dynamic programming methods (that can “stitch” suboptimal trajectories together) with the stability of supervised learning using by DT.

The **diffuser** method of [Jan+22] is a diffusion version of trajectory transformer, so it fits $p(\mathbf{s}_{1:T}, \mathbf{a}_{1:T}, \mathbf{r}_{1:T})$ using diffusion, where the action space is assumed to be continuous. They also replace beam search with classifier guidance. The **decision diffuser** method of [Aja+23] extends diffuser by using classifier-free guidance, where the conditioning signal is the reward-to-go, similar to decision transformer. However, unlike diffuser, the decision diffuser just models the future state trajectories (rather than learning a joint distribution over states and actions), and infers the actions using an **inverse dynamics model** $a_t = \pi(s_t, s_{t+1})$, which is trained using supervised learning.

One problem with the above approaches is that conditioning on a desired return and taking the predicted action can fail dramatically in stochastic environments, since trajectories that result in a return may have only achieved that return due to chance [PMB22; Yan+23; Bra+22; Vil+22]. (This is the same as the optimism bias problem in the control-as-inference approach discussed in Section 3.6.2.)

In [Kon+24], they propose the **latent plan transformer**, that replaces conditioning on the reward-to-go with conditioning on a latent “plan”, $\mathbf{z} \in \mathbb{R}^D$. In more detail, they fit the following latent variable sequence model using MC-EM, where the E step is implemented using Langevin dynamics:

$$p(\mathbf{z})p(\tau|\mathbf{z})p(y|\mathbf{z}) \quad (7.61)$$

where $\tau = (s_1, a_1, \dots, s_T, a_T)$ is the state-action trajectory and y is the observed (trajectory-level) reward. The model for $p(\tau|\mathbf{z})$ is a causal transformer, which generates the action at each step given the previous states and actions and the latent plan. The model for $p(y|\mathbf{z})$ is just a Gaussian. The model for $p(\mathbf{z})$ is a Gaussian passed through a U-net style CNN, thus providing a richer prior. The latent variables provide a way to “stitch together” individual (high performing) trajectories, so that the learned policy can predict $p(a_t | s_t, \mathbf{z})$ even if the current state s_t is not on the training manifold (thus requiring generalization, a problem that behavior cloning faces [Ghu+24]). During decision time, they infer $\hat{\mathbf{z}} = \operatorname{argmax}_{\mathbf{z}} p(\mathbf{z} | y = y_{\max})$ using gradient ascent, and then autoregressively generate actions from $p(a_t | s_{1:t}, a_{1:t-1}, \hat{\mathbf{z}})$.

7.7.4 Offline-to-online methods

Despite the progress in offline RL, it is fundamentally more limited in what it can learn compared to online RL [OCD21], because agent cannot explore the consequences of its own actions (c.f., results in psychology

[MG14]). Therefore, there is a lot of interest in pre-training offline, and then using online finetuning. This is called the **offline-to-online** (O2O) paradigm. Unfortunately, due to the significant distribution shift between online experiences and offline data, most offline RL algorithms suffer from performance drops when they are finetuned online. Many different methods have been proposed to tackle this, a few of which we mention below. See <https://github.com/linhlpv/awesome-offline-to-online-RL-papers> for a more extensive list.

7.7.4.1 Calibrated Q learning

[Nak+23] suggest pre-training with CQL followed by online finetuning. Naively this does not work that well, because CQL can be too conservative, requiring the online learning to waste some time at the beginning fixing the pessimism. So they propose a small modification to Conservative Q learning, discussed in Section 7.7.1.4, known as **Calibrated Q learning**. This simply prevents CQL from being too conservative, by replacing the CQL regularizer in Equation (7.57) with a slightly modified expression. Then online finetuning is performed in the usual way.

7.7.4.2 Dagger

An alternative approach is the **Dagger** algorithm of [RGB11]. (Dagger is short for Dataset Aggregation.) This iteratively trains the policy on expert provided data. We start with an initial dataset \mathcal{D} (e.g., empty) and an initial policy π_1 (e.g., random). At iteration t , we run the current policy π_t in the environment to collect states $\{s_i\}$. We then ask an expert policy for the correct actions $a_i^* = \pi^*(s_i)$. We then aggregate the data to compute $\mathcal{D} = \mathcal{D} \cup \{(s_i, a_i^*)\}$, and train the new policy π_{t+1} on \mathcal{D} . The key idea is to not train passively on expert trajectories as in BC, but to train on the states that the policy actually visits. This avoids overfitting to idealized data and improves robustness (avoids compounding error), since the policy is learning the effects of its own causal interventions [Ort+21].

7.8 General RL, AIXI and universal AGI

The term “**general RL**” (see e.g., [Hut05; LHS13; HQC24; Maj21]) refers to the setup in which an agent receives a stream of observations o_1, o_2, \dots and rewards r_1, r_2, \dots , and performs a sequence of actions in response, a_1, a_2, \dots , but where we do not make any Markovian (or even stationarity) assumptions about the environment that generates the observation stream. Instead, we assume that the environment is a computable function or program p^* , which generated the observations $o_{1:t}$ and $r_{1:t}$ seen so far in response to the actions taken, $a_{1:t-1}$. We denote this by $U(p^*, \mathbf{a}_{1:t}) = (o_1 r_1 \dots o_t r_t)$, where U is a universal Turing machine. If we use the receding horizon control strategy (see Section 4.2.4), the optimal action at each step is the one that maximizes the posterior expected reward-to-go (out to some horizon m steps into the future). If we assume the agent represents the unknown environment as a program $p \in \mathcal{M}$, then the optimal action is given by the following **expectimax** formula:

$$a_t = \operatorname{argmax}_{a_t} \sum_{o_t, r_t} \dots \max_{a_m} \sum_{o_m, r_m} [r_t + \dots + r_m] \sum_{p: U(p, \mathbf{a}_{1:m}) = (o_1 r_1 \dots o_m r_m)} \Pr(p) \quad (7.62)$$

where $\Pr(p)$ is the prior probability of p , and we assume the likelihood is 1 if p can generate the observations given the actions, and is 0 otherwise.

One important question is: what is a reasonable prior over programs? In [Hut05], Marcus Hutter proposed to apply the idea of **Solomonoff induction** [Sol64] to the case of an online decision making agent. This amounts to using the prior $\Pr(p) = 2^{-\ell(p)}$, where $\ell(p)$ is the length of program p . This prior favors shorter programs, and the likelihood filters out programs that cannot explain the data. The resulting agent is known as **AIXI**, where “AI” stands for “Artificial Intelligence” and “XI” referring to the Greek letter ξ used in Solomonoff induction. The AIXI agent has been called the “most intelligent general-purpose agent possible” [HQC24], and can be viewed as the theoretical foundation of (universal) **artificial general intelligence** or **AGI**.

Unfortunately, the AIXI agent is intractable to compute, for two main reasons: (1) it relies on Solomonoff induction and Kolmogorov complexity, both of which are intractable; and (2) the expectimax computation is intractable. Fortunately, various tractable approximations have been devised. In lieu of Kolmogorov complexity, we can use measures like MDL (minimum description length), and for Solomonoff induction, we can use various local search or optimization algorithms through suitable function classes. For the expectimax computation, we can use MCTS (see Section 4.2.2) to approximate it. Alternatively, [GM+24] showed that it is possible to use meta learning to train a generic sequence predictor, such as a transformer or LSTM, on data generated by random Turing machines, so that the transformer learns to approximate a universal predictor. Another approach is to learn a policy (to avoid searching over action sequences) using TD-learning (Section 2.3.2); the weighting term in the policy mixture requires that the agent predict its own future actions, so this approach is known as **self-AIXI** [Cat+23].

Note that AIXI is a normative theory for optimal agents, but is not very practical, since it does not take computational limitations into account. In [Aru+24a; Aru+24b], they describe an approach which extends the above Bayesian framework, while also taking into account the data budget (due to limited environment interactions) that real agents must contend with (which prohibits modeling the entire environment or finding the optimal action). This approach, known as **Capacity-Limited Bayesian RL** (CBRL), combines Bayesian inference, RL, and rate distortion theory, and can be seen as a normative theoretical foundation for computationally bounded rational agents.

Chapter 8

Acknowledgements

Parts of this monograph are borrowed from chapters 34 and 35 of my textbook [Mur23], some of which was written with Lihong Li. However, this text supercedes those chapters, and goes beyond it in many ways. Thanks to the following people for feedback on the current document: Pablo Samuel Castro, Elad Hazan, Tuan Ahn Le, Dieterich Lawson, Marc Lanctot, David Pfau, Theo Weber. And thanks to Xinghua Lou for help with some of the figures.

Bibliography

- [Abd+18] A. Abdolmaleki, J. T. Springenberg, Y. Tassa, R. Munos, N. Heess, and M. Riedmiller. “Maximum a Posteriori Policy Optimisation”. In: *International Conference on Learning Representations*. Feb. 2018. URL: <https://openreview.net/pdf?id=S1ANxQW0b>.
- [Abd+25] M. Abdulhai, I. White, C. V. Snell, C. Sun, J. Hong, Y. Zhai, K. Xu, and S. Levine. “LMRL Gym: Benchmarks for Multi-Turn Reinforcement Learning with Language Models”. In: *ICML*. June 2025. URL: <https://openreview.net/forum?id=hmGhP5D02W>.
- [Abe+19] D. Abel, J. Winder, M. desJardins, and M. Littman. “The Expected-Length Model of options”. In: *AISTATS*. Aug. 2019. URL: <http://dx.doi.org/10.5555/3367243.3367310>.
- [ABM10] J.-Y. Audibert, S. Bubeck, and R. Munos. “Best Arm Identification in Multi-Armed Bandits”. In: *COLT*. 2010, pp. 41–53.
- [ACBF02] P. Auer, N. Cesa-Bianchi, and P. Fischer. “Finite-time Analysis of the Multiarmed Bandit Problem”. In: *MLJ* 47.2 (May 2002), pp. 235–256. URL: <http://mercurio.srv.di.unimi.it/~cesabian/Pubblicazioni/ml-02.pdf>.
- [Ach+17] J. Achiam, D. Held, A. Tamar, and P. Abbeel. “Constrained Policy Optimization”. In: *ICML*. 2017. URL: <http://arxiv.org/abs/1705.10528>.
- [Ach+18] J. Achiam, H. Edwards, D. Amodei, and P. Abbeel. “Variational Option Discovery Algorithms”. In: *arXiv [cs.AI]* (July 2018). URL: <http://arxiv.org/abs/1807.10299>.
- [ACS24] S. V. Albrecht, F. Christianos, and L. Schäfer. *Multi-Agent Reinforcement Learning: Foundations and Modern Approaches*. MIT Press, 2024. URL: <https://www.marl-book.com>.
- [AG25] D. Arumugam and T. L. Griffiths. “Toward efficient exploration by large language model agents”. In: *arXiv [cs.LG]* (Apr. 2025). URL: <http://arxiv.org/abs/2504.20997>.
- [Aga+14] D. Agarwal, B. Long, J. Traupman, D. Xin, and L. Zhang. “LASER: a scalable response prediction platform for online advertising”. In: *WSDM*. 2014.
- [Aga+21a] A. Agarwal, S. M. Kakade, J. D. Lee, and G. Mahajan. “On the Theory of Policy Gradient Methods: Optimality, Approximation, and Distribution Shift”. In: *JMLR* 22.98 (2021), pp. 1–76. URL: <http://jmlr.org/papers/v22/19-736.html>.
- [Aga+21b] R. Agarwal, M. Schwarzer, P. S. Castro, A. Courville, and M. G. Bellemare. “Deep Reinforcement Learning at the Edge of the Statistical Precipice”. In: *NIPS*. Aug. 2021. URL: <http://arxiv.org/abs/2108.13264>.
- [Aga+22a] A. Agarwal, N. Jiang, S. Kakade, and W. Sun. *Reinforcement Learning: Theory and Algorithms*. 2022. URL: https://rltheorybook.github.io/rltheorybook_AJKS.pdf.
- [Aga+22b] R. Agarwal, M. Schwarzer, P. S. Castro, A. C. Courville, and M. Bellemare. “Reincarnating Reinforcement Learning: Reusing Prior Computation to Accelerate Progress”. In: *NIPS*. Vol. 35. 2022, pp. 28955–28971. URL: https://proceedings.neurips.cc/paper_files/paper/2022/hash/ba1c5356d9164bb64c446a4b690226b0-Abstract-Conference.html.
- [Agr+25] L. A. Agrawal et al. “GEPA: Reflective prompt evolution can outperform reinforcement learning”. In: *arXiv [cs.CL]* (July 2025). URL: <http://arxiv.org/abs/2507.19457>.

- [AH81] R. Axelrod and W. Hamilton. “The evolution of cooperation”. In: *Science* 4489 (1981), pp. 1390–1396.
- [Ahm+24] A. Ahmadian, C. Cremer, M. Gallé, M. Fadaee, J. Kreutzer, A. Üstün, and S. Hooker. “Back to basics: Revisiting REINFORCE style optimization for learning from Human Feedback in LLMs”. In: *arXiv [cs.LG]* (Feb. 2024). URL: <http://arxiv.org/abs/2402.14740>.
- [Aja+23] A. Ajay, Y. Du, A. Gupta, J. B. Tenenbaum, T. S. Jaakkola, and P. Agrawal. “Is Conditional Generative Modeling all you need for Decision Making?” In: *ICLR*. 2023. URL: <https://openreview.net/forum?id=sP1fo2K9DFG>.
- [AJO08] P. Auer, T. Jaksch, and R. Ortner. “Near-optimal Regret Bounds for Reinforcement Learning”. In: *NIPS*. Vol. 21. 2008. URL: https://proceedings.neurips.cc/paper_files/paper/2008/file/e4a6222cdb5b34375400904f03d8e6a5-Paper.pdf.
- [al24] J. P.-H. et al. “Genie 2: A large-scale foundation world model”. In: (2024). URL: <https://deepmind.google/discover/blog/genie-2-a-large-scale-foundation-world-model/>.
- [Ale+23] L. N. Alegre, A. L. C. Bazzan, A. Nowé, and B. C. da Silva. “Multi-step generalized policy improvement by leveraging approximate models”. In: *NIPS*. Vol. 36. Curran Associates, Inc., 2023, pp. 38181–38205. URL: https://proceedings.neurips.cc/paper_files/paper/2023/hash/77c7faab15002432ba1151e8d5cc389a-Abstract-Conference.html.
- [Alo+24] E. Alonso, A. Jelley, V. Micheli, A. Kanervisto, A. Storkey, T. Pearce, and F. Fleuret. “Diffusion for world modeling: Visual details matter in Atari”. In: *arXiv [cs.LG]* (May 2024). URL: <http://arxiv.org/abs/2405.12399>.
- [AM89] B. D. Anderson and J. B. Moore. *Optimal Control: Linear Quadratic Methods*. Prentice-Hall International, Inc., 1989.
- [Ama98] S. Amari. “Natural Gradient Works Efficiently in Learning”. In: *Neural Comput.* 10.2 (1998), pp. 251–276. URL: <http://dx.doi.org/10.1162/089976698300017746>.
- [AMH19] A. Aubret, L. Matignon, and S. Hassas. “A survey on intrinsic motivation in reinforcement learning”. In: *arXiv [cs.LG]* (Aug. 2019). URL: <http://arxiv.org/abs/1908.06976>.
- [AMH23] A. Aubret, L. Matignon, and S. Hassas. “An information-theoretic perspective on intrinsic motivation in reinforcement learning: A survey”. en. In: *Entropy* 25.2 (Feb. 2023), p. 327. URL: <https://www.mdpi.com/1099-4300/25/2/327>.
- [Ami+21] S. Amin, M. Gomrokchi, H. Satija, H. van Hoof, and D. Precup. “A survey of exploration methods in reinforcement learning”. In: *arXiv [cs.LG]* (Aug. 2021). URL: <http://arxiv.org/abs/2109.00157>.
- [An+21] G. An, S. Moon, J.-H. Kim, and H. O. Song. “Uncertainty-Based Offline Reinforcement Learning with Diversified Q-Ensemble”. In: *NIPS*. Vol. 34. Dec. 2021, pp. 7436–7447. URL: https://proceedings.neurips.cc/paper_files/paper/2021/file/3d3d286a8d153a4a58156d0e02d8570c-Paper.pdf.
- [And+17] M. Andrychowicz, F. Wolski, A. Ray, J. Schneider, R. Fong, P. Welinder, B. McGrew, J. Tobin, P. Abbeel, and W. Zaremba. “Hindsight Experience Replay”. In: *arXiv [cs.LG]* (July 2017). URL: <http://arxiv.org/abs/1707.01495>.
- [And+20] O. M. Andrychowicz et al. “Learning dexterous in-hand manipulation”. In: *Int. J. Rob. Res.* 39.1 (2020), pp. 3–20. URL: <https://doi.org/10.1177/0278364919887447>.
- [Ant22] Anthropic. “Constitutional AI: Harmlessness from AI Feedback”. In: *arXiv [cs.CL]* (Dec. 2022). URL: <http://arxiv.org/abs/2212.08073>.
- [Ant+22] I. Antonoglou, J. Schrittwieser, S. Ozair, T. K. Hubert, and D. Silver. “Planning in Stochastic Environments with a Learned Model”. In: *ICLR*. 2022. URL: <https://openreview.net/forum?id=X6D9bAHhBQ1>.

- [AP23] S. Alver and D. Precup. “Minimal Value-Equivalent Partial Models for Scalable and Robust Planning in Lifelong Reinforcement Learning”. en. In: *Conference on Lifelong Learning Agents*. PMLR, Nov. 2023, pp. 548–567. URL: <https://proceedings.mlr.press/v232/alver23a.html>.
- [AP24] S. Alver and D. Precup. “A Look at Value-Based Decision-Time vs. Background Planning Methods Across Different Settings”. In: *Seventeenth European Workshop on Reinforcement Learning*. Oct. 2024. URL: <https://openreview.net/pdf?id=Vx2ETvHId8>.
- [Arb+23] J. Arbel, K. Pitas, M. Vladimirova, and V. Fortuin. “A Primer on Bayesian Neural Networks: Review and Debates”. In: *arXiv [stat.ML]* (Sept. 2023). URL: <http://arxiv.org/abs/2309.16314>.
- [ARKP24] S. Alver, A. Rahimi-Kalahroudi, and D. Precup. “Partial models for building adaptive model-based reinforcement learning agents”. In: *COLLAS*. May 2024. URL: <https://arxiv.org/abs/2405.16899>.
- [Aru+17] K. Arulkumaran, M. P. Deisenroth, M. Brundage, and A. A. Bharath. “A Brief Survey of Deep Reinforcement Learning”. In: *IEEE Signal Processing Magazine, Special Issue on Deep Learning for Image Understanding* (2017). URL: <http://arxiv.org/abs/1708.05866>.
- [Aru+18] D. Arumugam, D. Abel, K. Asadi, N. Gopalan, C. Grimm, J. K. Lee, L. Lehnert, and M. L. Littman. “Mitigating planner overfitting in model-based reinforcement learning”. In: *arXiv [cs.LG]* (Dec. 2018). URL: <http://arxiv.org/abs/1812.01129>.
- [Aru+24a] D. Arumugam, M. K. Ho, N. D. Goodman, and B. Van Roy. “Bayesian Reinforcement Learning With Limited Cognitive Load”. en. In: *Open Mind* 8 (Apr. 2024), pp. 395–438. URL: https://direct.mit.edu/opmi/article-pdf/doi/10.1162/opmi_a_00132/2364075/opmi_a_00132.pdf.
- [Aru+24b] D. Arumugam, S. Kumar, R. Gummadi, and B. Van Roy. “Satisficing exploration for deep reinforcement learning”. In: *Finding the Frame Workshop at RLC*. July 2024. URL: <https://openreview.net/forum?id=tHCpsrزهb>.
- [AS18] S. V. Albrecht and P. Stone. “Autonomous agents modelling other agents: A comprehensive survey and open problems”. en. In: *Artif. Intell.* 258 (May 2018), pp. 66–95. URL: <http://dx.doi.org/10.1016/j.artint.2018.01.002>.
- [AS22] D. Arumugam and S. Singh. “Planning to the information horizon of BAMDPs via epistemic state abstraction”. In: *NIPS*. Oct. 2022.
- [AS66] S. M. Ali and S. D. Silvey. “A General Class of Coefficients of Divergence of One Distribution from Another”. In: *J. R. Stat. Soc. Series B Stat. Methodol.* 28.1 (1966), pp. 131–142. URL: <http://www.jstor.org/stable/2984279>.
- [ASN20] R. Agarwal, D. Schuurmans, and M. Norouzi. “An Optimistic Perspective on Offline Reinforcement Learning”. en. In: *ICML*. PMLR, Nov. 2020, pp. 104–114. URL: <https://proceedings.mlr.press/v119/agarwal20c.html>.
- [Ass+23] M. Assran, Q. Duval, I. Misra, P. Bojanowski, P. Vincent, M. Rabbat, Y. LeCun, and N. Ballas. “Self-Supervised Learning from Images with a Joint-Embedding Predictive Architecture”. In: *CVPR*. Jan. 2023. URL: <http://arxiv.org/abs/2301.08243>.
- [Ass+25] M. Assran et al. “V-JEPA 2: Self-supervised video models enable understanding, prediction and planning”. In: *arXiv [cs.AI]* (June 2025). URL: <http://arxiv.org/abs/2506.09985>.
- [AT17] R. Ackerman and V. A. Thompson. “Meta-reasoning: Monitoring and control of thinking and reasoning”. en. In: *Trends Cogn. Sci.* 21.8 (Aug. 2017), pp. 607–617. URL: <https://www.sciencedirect.com/science/article/pii/S1364661317301055>.
- [Att03] H. Attias. “Planning by Probabilistic Inference”. In: *AI-Stats*. 2003. URL: <http://research.goldenmetallic.com/aistats03.pdf>.

- [Aue+19] P. Auer, Y. Chen, P. Gajane, C.-W. Lee, H. Luo, R. Ortner, and C.-Y. Wei. “Achieving Optimal Dynamic Regret for Non-stationary Bandits without Prior Information”. en. In: *Conference on Learning Theory*. PMLR, June 2019, pp. 159–163. URL: <https://proceedings.mlr.press/v99/auer19b.html>.
- [Aum87] R. J. Aumann. “Correlated equilibrium as an expression of Bayesian rationality”. en. In: *Econometrica* 55.1 (Jan. 1987), p. 1. URL: <https://www.jstor.org/stable/1911154>.
- [Axe84] R. Axelrod. *The evolution of cooperation*. Basic Books, 1984.
- [Bad+20] A. P. Badia, B. Piot, S. Kapturowski, P. Sprechmann, A. Vitvitskyi, D. Guo, and C. Blundell. “Agent57: Outperforming the Atari Human Benchmark”. In: *ICML* 119 (Mar. 2020), pp. 507–517. URL: <https://proceedings.mlr.press/v119/badia20a/badia20a.pdf>.
- [Bai+21] Y. Bai, C. Jin, H. Wang, and C. Xiong. “Sample-Efficient Learning of Stackelberg Equilibria in General-Sum Games”. In: *NIPS*. Nov. 2021. URL: <https://openreview.net/forum?id=LZOG2YgDiRn>.
- [Bai95] L. C. Baird. “Residual Algorithms: Reinforcement Learning with Function Approximation”. In: *ICML*. 1995, pp. 30–37.
- [Bak+21] A. Bakhtin, D. J. Wu, A. Lerer, and N. Brown. “No-Press Diplomacy from Scratch”. In: *NIPS*. Nov. 2021. URL: <https://openreview.net/forum?id=Pq7wIzt30UE>.
- [Bak+22] B. Baker, I. Akkaya, P. Zhokhov, J. Huizinga, J. Tang, A. Ecoffet, B. Houghton, R. Sampedro, and J. Clune. “Video PreTraining (VPT): Learning to Act by Watching Unlabeled Online Videos”. In: *NIPS*. June 2022. URL: <https://arxiv.org/abs/2206.11795>.
- [Bal+23] P. J. Ball, L. Smith, I. Kostrikov, and S. Levine. “Efficient Online Reinforcement Learning with Offline Data”. en. In: *ICML*. PMLR, July 2023, pp. 1577–1594. URL: <https://proceedings.mlr.press/v202/ball123a.html>.
- [Bal+25a] A. Balashankar et al. “InfAlign: Inference-aware language model alignment”. In: *ICML*. Aug. 2025. URL: <https://arxiv.org/abs/2412.19792>.
- [Bal+25b] F. Baldassarre, M. Szafraniec, B. Terver, V. Khalidov, F. Massa, Y. LeCun, P. Labatut, M. Seitzer, and P. Bojanowski. “Back to the features: DINO as a foundation for video world models”. In: *arXiv [cs.CV]* (July 2025). URL: <http://arxiv.org/abs/2507.19468>.
- [Ban+23] D. Bansal, R. T. Q. Chen, M. Mukadam, and B. Amos. “TaskMet: Task-driven metric learning for model learning”. In: *NIPS*. Vol. abs/2312.05250. Dec. 2023, pp. 46505–46519. URL: https://proceedings.neurips.cc/paper_files/paper/2023/hash/91a5742235f70ae846436d9780e9f1d4-Abstract-Conference.html.
- [Bar+17] A. Barreto, W. Dabney, R. Munos, J. J. Hunt, T. Schaul, H. P. van Hasselt, and D. Silver. “Successor Features for Transfer in Reinforcement Learning”. In: *NIPS*. Vol. 30. 2017. URL: https://proceedings.neurips.cc/paper_files/paper/2017/file/350db081a661525235354dd3e19b8c05-Paper.pdf.
- [Bar+19] A. Barreto et al. “The Option Keyboard: Combining Skills in Reinforcement Learning”. In: *NIPS*. Vol. 32. 2019. URL: https://proceedings.neurips.cc/paper_files/paper/2019/file/251c5ffd6b62cc21c446c963c76cf214-Paper.pdf.
- [Bar+20] A. Barreto, S. Hou, D. Borsa, D. Silver, and D. Precup. “Fast reinforcement learning with generalized policy updates”. en. In: *PNAS* 117.48 (Dec. 2020), pp. 30079–30087. URL: <https://www.pnas.org/doi/abs/10.1073/pnas.1907370117>.
- [Bar+24] A. Bardes, Q. Garrido, J. Ponce, X. Chen, M. Rabbat, Y. Le Cun, M. Assran, and N. Ballas. “Revisiting Feature Prediction for Learning Visual Representations from Video”. In: *TMLR* (2024). URL: <https://openreview.net/forum?id=QaCCuDfBk2>.
- [Bau+23] K. Baumli et al. “Vision-language models as a source of rewards”. In: *arXiv [cs.LG]* (Dec. 2023). URL: <http://arxiv.org/abs/2312.09187>.

- [Bau+24] M. Bauza et al. “DemoStart: Demonstration-led auto-curriculum applied to sim-to-real with multi-fingered robots”. In: (Sept. 2024). URL: <http://arxiv.org/abs/2409.06613>.
- [BAVE25] N. M. Boffi, M. S. Albergo, and E. Vanden-Eijnden. “How to build a consistency model: Learning flow maps via self-distillation”. In: *NIPS*. Oct. 2025.
- [BBS95] A. G. Barto, S. J. Bradtke, and S. P. Singh. “Learning to act using real-time dynamic programming”. In: *AIJ* 72.1 (1995), pp. 81–138. URL: <http://www.sciencedirect.com/science/article/pii/0004370294000110>.
- [BDG00] C. Boutilier, R. Dearden, and M. Goldszmidt. “Stochastic dynamic programming with factored representations”. en. In: *Artif. Intell.* 121.1-2 (Aug. 2000), pp. 49–107. URL: [http://dx.doi.org/10.1016/S0004-3702\(00\)00033-3](http://dx.doi.org/10.1016/S0004-3702(00)00033-3).
- [BDM17] M. G. Bellemare, W. Dabney, and R. Munos. “A Distributional Perspective on Reinforcement Learning”. In: *ICML*. 2017. URL: <http://arxiv.org/abs/1707.06887>.
- [BDR23] M. G. Bellemare, W. Dabney, and M. Rowland. *Distributional Reinforcement Learning*. <http://www.distributional-rl.org>. MIT Press, 2023.
- [Bec+24] M. Beck, K. Pöppel, M. Spanring, A. Auer, O. Prudnikova, M. K. Kopp, G. Klambauer, J. Brandstetter, and S. Hochreiter. “xLSTM: Extended Long Short-Term Memory”. In: *NIPS*. Nov. 2024. URL: <https://openreview.net/forum?id=ARAxPPIAhq>.
- [Bec+25] M. Beck, K. Schweighofer, S. Böck, S. Lehner, and S. Hochreiter. “XLSTM scaling laws: Competitive performance with linear time-complexity”. In: *arXiv [cs.LG]* (Oct. 2025). URL: <http://arxiv.org/abs/2510.02228>.
- [Beh+25a] A. Behrouz, Z. Li, P. Kacham, M. Daliri, Y. Deng, P. Zhong, M. Razaviyayn, and V. Mirrokni. “ATLAS: Learning to optimally memorize the context at test time”. In: *arXiv [cs.CL]* (May 2025). URL: <http://arxiv.org/abs/2505.23735>.
- [Beh+25b] A. Behrouz, M. Razaviyayn, P. Zhong, and V. Mirrokni. “It’s all connected: A journey through test-time memorization, attentional bias, retention, and online optimization”. In: *arXiv [cs.LG]* (Apr. 2025). URL: <http://arxiv.org/abs/2504.13173>.
- [Bei+25] A. Beirami, A. Agarwal, J. Berant, A. D’Amour, J. Eisenstein, C. Nagpal, and A. T. Suresh. “Theoretical guarantees on the best-of-n alignment policy”. In: *ICML*. 2025. URL: <https://arxiv.org/abs/2401.01879>.
- [Bel+13] M. G. Bellemare, Y. Naddaf, J. Veness, and M. Bowling. “The Arcade Learning Environment: An Evaluation Platform for General Agents”. In: *JAIR* 47 (2013), pp. 253–279.
- [Bel+16] M. G. Bellemare, S. Srinivasan, G. Ostrovski, T. Schaul, D. Saxton, and R. Munos. “Unifying Count-Based Exploration and Intrinsic Motivation”. In: *NIPS*. 2016. URL: <http://arxiv.org/abs/1606.01868>.
- [Ber19] D. Bertsekas. *Reinforcement learning and optimal control*. Athena Scientific, 2019. URL: <http://www.mit.edu/~dimitrib/RLbook.html>.
- [Ber+21] G. Berseth, D. Geng, C. Devin, N. Rhinehart, C. Finn, D. Jayaraman, and S. Levine. “SMiRL: Surprise minimizing reinforcement learning in unstable environments”. In: *ICLR*. 2021. URL: <https://arxiv.org/pdf/1912.05510>.
- [Ber24] D. P. Bertsekas. “Model Predictive Control and Reinforcement Learning: A unified framework based on Dynamic Programming”. In: *arXiv [eess.SY]* (June 2024). URL: <http://arxiv.org/abs/2406.00592>.
- [BHP17] P.-L. Bacon, J. Harb, and D. Precup. “The Option-Critic Architecture”. In: *AAAI*. 2017.
- [BK20] A. Bagaria and G. Konidaris. “Option Discovery using Deep Skill Chaining”. In: *ICLR*. 2020. URL: <https://openreview.net/forum?id=BlgqipNYwH>.
- [BKH16] J. L. Ba, J. R. Kiros, and G. E. Hinton. “Layer Normalization”. In: (2016). arXiv: [1607.06450](https://arxiv.org/abs/1607.06450) [stat.ML]. URL: <http://arxiv.org/abs/1607.06450>.

- [BKM24] W Bradley Knox and J. MacGlashan. “How to Specify Reinforcement Learning Objectives”. In: *Finding the Frame: An RLC Workshop for Examining Conceptual Frameworks*. July 2024. URL: <https://openreview.net/pdf?id=2MGEQNrmDN>.
- [BLM16] S. Boucheron, G. Lugosi, and P. Massart. *Concentration Inequalities: A Nonasymptotic Theory of Independence*. Oxford University Press, 2016.
- [Blo+15] D. Bloembergen, K. Tuyls, D. Hennes, and M. Kaisers. “Evolutionary dynamics of multi-agent learning: A survey”. en. In: *JAIR* 53 (Aug. 2015), pp. 659–697. URL: <https://jair.org/index.php/jair/article/view/10952>.
- [BM+18] G. Barth-Maron, M. W. Hoffman, D. Budden, W. Dabney, D. Horgan, T. B. Dhruva, A. Muldal, N. Heess, and T. Lillicrap. “Distributed Distributional Deterministic Policy Gradients”. In: *ICLR*. 2018. URL: <https://openreview.net/forum?id=SyZipzbCb¬eId=SyZipzbCb>.
- [BMS11] S. Bubeck, R. Munos, and G. Stoltz. “Pure Exploration in Finitely-armed and Continuous-armed Bandits”. In: *Theoretical Computer Science* 412.19 (2011), pp. 1832–1852.
- [Bo+25] L. Bo et al. “SPIRAL: Self-play on zero-sum games incentivizes reasoning via multi-agent multi-turn reinforcement learning”. In: *arXiv [cs.AI]* (June 2025). URL: <http://arxiv.org/abs/2506.24119>.
- [Boe+05] P.-T. de Boer, D. P. Kroese, S. Mannor, and R. Y. Rubinstein. “A Tutorial on the Cross-Entropy Method”. en. In: *Ann. Oper. Res.* 134.1 (2005), pp. 19–67. URL: <https://link.springer.com/article/10.1007/s10479-005-5724-z>.
- [Boo+23] S. Booth, W. B. Knox, J. Shah, S. Niekum, P. Stone, and A. Allievi. “The Perils of Trial-and-Error Reward Design: Misdesign through Overfitting and Invalid Task Specifications”. In: *AAAI*. 2023. URL: https://slbooth.com/assets/projects/Reward_Design_Perils/.
- [Bor+19] D. Borsa, A. Barreto, J. Quan, D. J. Mankowitz, H. van Hasselt, R. Munos, D. Silver, and T. Schaul. “Universal Successor Features Approximators”. In: *ICLR*. 2019. URL: <https://openreview.net/pdf?id=S1VWjiRcKX>.
- [Bos16] N. Bostrom. *Superintelligence: Paths, Dangers, Strategies*. en. London, England: Oxford University Press, Mar. 2016. URL: <https://www.amazon.com/Superintelligence-Dangers-Strategies-Nick-Bostrom/dp/0198739834>.
- [Bou+23] K. Bousmalis et al. “RoboCat: A Self-Improving Generalist Agent for Robotic Manipulation”. In: *TMLR* (June 2023). URL: <http://arxiv.org/abs/2306.11706>.
- [Bow+23] M. Bowling, J. D. Martin, D. Abel, and W. Dabney. “Settling the reward hypothesis”. en. In: *ICML*. 2023. URL: <https://arxiv.org/abs/2212.10420>.
- [BPL22a] A. Bardes, J. Ponce, and Y. LeCun. “VICReg: Variance-Invariance-Covariance Regularization for Self-Supervised Learning”. In: *ICLR*. 2022. URL: <https://openreview.net/pdf?id=xm6YD62D1Ub>.
- [BPL22b] A. Bardes, J. Ponce, and Y. LeCun. “VICRegL: Self-supervised learning of local visual features”. In: *NIPS*. Oct. 2022. URL: <https://arxiv.org/abs/2210.01571>.
- [Bra+22] D. Brandfonbrener, A. Bietti, J. Buckman, R. Laroché, and J. Bruna. “When does return-conditioned supervised learning work for offline reinforcement learning?” In: *NIPS*. June 2022. URL: <http://arxiv.org/abs/2206.01079>.
- [Bro+12] C. B. Browne, E. Powley, D. Whitehouse, S. M. Lucas, P. I. Cowling, P. Rohlfshagen, S. Tavener, D. Perez, S. Samothrakis, and S. Colton. “A Survey of Monte Carlo Tree Search Methods”. In: *IEEE Transactions on Computational Intelligence and AI in Games* 4.1 (2012).
- [Bro+19] N. Brown, A. Lerer, S. Gross, and T. Sandholm. “Deep Counterfactual Regret Minimization”. en. In: *ICML*. PMLR, May 2019, pp. 793–802. URL: <https://proceedings.mlr.press/v97/brown19b.html>.

- [Bro24] W. Brown. *Generative AI Handbook: A Roadmap for Learning Resources*. 2024. URL: <https://genai-handbook.github.io>.
- [BS17] N. Brown and T. Sandholm. “Superhuman AI for heads-up no-limit poker: Libratus beats top professionals”. en. In: *Science* 359.6374 (2017), pp. 418–424. URL: <https://www.science.org/doi/10.1126/science.aao1733>.
- [BS19] N. Brown and T. Sandholm. “Superhuman AI for multiplayer poker”. en. In: *Science* 365.6456 (Aug. 2019), pp. 885–890. URL: <https://www.science.org/doi/10.1126/science.aay2400>.
- [BSA83] A. G. Barto, R. S. Sutton, and C. W. Anderson. “Neuronlike adaptive elements that can solve difficult learning control problems”. In: *SMC* 13.5 (1983), pp. 834–846. URL: <http://dx.doi.org/10.1109/TSMC.1983.6313077>.
- [BT12] M. Botvinick and M. Toussaint. “Planning as inference”. en. In: *Trends Cogn. Sci.* 16.10 (2012), pp. 485–488. URL: <https://pdfs.semanticscholar.org/2ba7/88647916f6206f7fcc137fe7866c58e6211e.pdf>.
- [Buc+17] C. L. Buckley, C. S. Kim, S. McGregor, and A. K. Seth. “The free energy principle for action and perception: A mathematical review”. In: *J. Math. Psychol.* 81 (2017), pp. 55–79. URL: <https://www.sciencedirect.com/science/article/pii/S0022249617300962>.
- [Bur+18] Y. Burda, H. Edwards, A. Storkey, and O. Klimov. “Exploration by random network distillation”. In: *ICLR*. Vol. abs/1810.12894. Sept. 2018.
- [Bur25] A. Burkov. *The Hundred-Page Language Models Book*. 2025. URL: <https://thelmlbook.com/>.
- [BV02] M. Bowling and M. Veloso. “Multiagent learning using a variable learning rate”. en. In: *Artif. Intell.* 136.2 (Apr. 2002), pp. 215–250. URL: [http://dx.doi.org/10.1016/S0004-3702\(02\)00121-2](http://dx.doi.org/10.1016/S0004-3702(02)00121-2).
- [BXS20] H. Bharadhwaj, K. Xie, and F. Shkurti. “Model-Predictive Control via Cross-Entropy and Gradient-Based Optimization”. en. In: *Learning for Dynamics and Control*. PMLR, July 2020, pp. 277–286. URL: <https://proceedings.mlr.press/v120/bharadhwaj20a.html>.
- [CA13] E. F. Camacho and C. B. Alba. *Model predictive control*. Springer, 2013.
- [Cao+24] Y. Cao, H. Zhao, Y. Cheng, T. Shu, G. Liu, G. Liang, J. Zhao, and Y. Li. “Survey on large language model-enhanced reinforcement learning: Concept, taxonomy, and methods”. In: *IEEE Transactions on Neural Networks and Learning Systems* (2024). URL: <http://arxiv.org/abs/2404.00282>.
- [Cao+25a] H. Cao, F. Feng, M. Fang, S. Dong, T. Yang, J. Huo, and Y. Gao. “Towards empowerment gain through causal structure learning in model-based RL”. In: Feb. 2025. URL: <http://arxiv.org/abs/2502.10077>.
- [Cao+25b] H. Cao, F. Feng, J. Huo, and Y. Gao. “Causal action empowerment for efficient reinforcement learning in embodied agents”. en. In: *Sci. China Inf. Sci.* 68.5 (May 2025), pp. 1–15. URL: <https://link.springer.com/article/10.1007/s11432-024-4396-3>.
- [Cao+25c] H. Cao, F. Feng, T. Yang, J. Huo, and Y. Gao. “Causal Information Prioritization for efficient Reinforcement Learning”. In: *arXiv [cs.AI]* (Feb. 2025). URL: <http://arxiv.org/abs/2502.10097>.
- [Car+23a] T. Carta, C. Romac, T. Wolf, S. Lamprier, O. Sigaud, and P.-Y. Oudeyer. “Grounding Large Language Models in interactive environments with online Reinforcement Learning”. In: *ICML*. Feb. 2023. URL: <https://arxiv.org/abs/2302.02662>.
- [Car+23b] W. C. Carvalho, A. Saraiva, A. Filos, A. Lampinen, L. Matthey, R. L. Lewis, H. Lee, S. Singh, D. Jimenez Rezende, and D. Zoran. “Combining Behaviors with the Successor Features Keyboard”. In: *NIPS*. Vol. 36. 2023, pp. 9956–9983. URL: https://proceedings.neurips.cc/paper_files/paper/2023/hash/1f69928210578f4cf5b538a8c8806798-Abstract-Conference.html.

- [Car+24] W. Carvalho, M. S. Tomov, W. de Cothi, C. Barry, and S. J. Gershman. “Predictive representations: building blocks of intelligence”. In: *Neural Comput.* (Feb. 2024). URL: <https://gershmanlab.com/pubs/Carvalho24.pdf>.
- [Cas11] P. S. Castro. “On planning, prediction and knowledge transfer in Fully and Partially Observable Markov Decision Processes”. en. PhD thesis. McGill, 2011. URL: <https://www.proquest.com/openview/d35984acba38c072359f8a8d5102c777/1?pq-origsite=gscholar&cbl=18750>.
- [Cas20] P. S. Castro. “Scalable methods for computing state similarity in deterministic Markov Decision Processes”. In: *AAAI*. 2020.
- [Cas+21] P. S. Castro, T. Kastner, P. Panangaden, and M. Rowland. “MICo: Improved representations via sampling-based state similarity for Markov decision processes”. In: *NIPS*. Nov. 2021. URL: <https://openreview.net/pdf?id=wFp6kmQELgu>.
- [Cas+23] P. S. Castro, T. Kastner, P. Panangaden, and M. Rowland. “A kernel perspective on behavioural metrics for Markov decision processes”. In: *TMLR abs/2310.19804* (Oct. 2023). URL: <https://openreview.net/pdf?id=nHfPX11ly7>.
- [Cat+23] E. Catt, J. Grau-Moya, M. Hutter, M. Aitchison, T. Genewein, G. Delétang, K. Li, and J. Veness. “Self-Predictive Universal AI”. In: *NIPS*. Vol. 36. 2023, pp. 27181–27198. URL: https://proceedings.neurips.cc/paper_files/paper/2023/hash/56a225639da77e8f7c0409f6d5ba996b-Abstract-Conference.html.
- [CC25] I. Chrysomallis and G. Chalkiadakis. “Imitation learning in the deep learning era: A novel taxonomy and recent advances”. In: *arXiv [cs.LG]* (Nov. 2025). URL: <http://arxiv.org/abs/2511.03565>.
- [Cet+24] E. Cetin, A. Tirinzoni, M. Pirotta, A. Lazaric, Y. Ollivier, and A. Touati. “Simple ingredients for offline reinforcement learning”. In: *arXiv [cs.LG]* (Mar. 2024). URL: <http://arxiv.org/abs/2403.13097>.
- [CH20] X. Chen and K. He. “Exploring simple Siamese representation learning”. In: *arXiv [cs.CV]* (Nov. 2020). URL: <http://arxiv.org/abs/2011.10566>.
- [Che+20] X. Chen, C. Wang, Z. Zhou, and K. W. Ross. “Randomized Ensembled Double Q-Learning: Learning Fast Without a Model”. In: *ICLR*. Oct. 2020. URL: <https://openreview.net/pdf?id=AY8zfZmOtDd>.
- [Che+21a] C. Chen, Y.-F. Wu, J. Yoon, and S. Ahn. “TransDreamer: Reinforcement Learning with Transformer World Models”. In: *Deep RL Workshop NeurIPS*. 2021. URL: <http://arxiv.org/abs/2202.09481>.
- [Che+21b] L. Chen, K. Lu, A. Rajeswaran, K. Lee, A. Grover, M. Laskin, P. Abbeel, A. Srinivas, and I. Mordatch. “Decision Transformer: Reinforcement Learning via Sequence Modeling”. In: *arXiv [cs.LG]* (June 2021). URL: <http://arxiv.org/abs/2106.01345>.
- [Che+24a] F. Che, C. Xiao, J. Mei, B. Dai, R. Gummadi, O. A. Ramirez, C. K. Harris, A. R. Mahmood, and D. Schuurmans. “Target networks and over-parameterization stabilize off-policy bootstrapping with function approximation”. In: *ICML*. May 2024. URL: <http://arxiv.org/abs/2405.21043>.
- [Che+24b] J. Chen, B. Ganguly, Y. Xu, Y. Mei, T. Lan, and V. Aggarwal. “Deep Generative Models for Offline Policy Learning: Tutorial, Survey, and Perspectives on Future Directions”. In: *TMLR* (Feb. 2024). URL: <https://openreview.net/forum?id=Mm2cMD19r5>.
- [Che+24c] J. Chen, B. Ganguly, Y. Xu, Y. Mei, T. Lan, and V. Aggarwal. “Deep Generative Models for Offline Policy Learning: Tutorial, Survey, and Perspectives on Future Directions”. In: *TMLR* (Feb. 2024). URL: <https://openreview.net/forum?id=Mm2cMD19r5>.
- [Che+24d] W. Chen, O. Mees, A. Kumar, and S. Levine. “Vision-language models provide promptable representations for reinforcement learning”. In: *arXiv [cs.LG]* (Feb. 2024). URL: <http://arxiv.org/abs/2402.02651>.

- [Che+25] W. Chen, J. Chen, H. Zhu, and J. Schneider. “Verlog: A Multi-turn RL framework for LLM agents”. In: (2025). URL: <https://blog.ml.cmu.edu/2025/09/15/verlog-a-multi-turn-rl-framework-for-llm-agents/>.
- [Chi+23] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song. “Diffusion Policy: Visuomotor Policy Learning via Action Diffusion”. In: *RSS*. Mar. 2023. URL: <http://arxiv.org/abs/2303.04137>.
- [Cho+21] J. Choi, A. Sharma, H. Lee, S. Levine, and S. S. Gu. “Variational empowerment as representation learning for goal-based reinforcement learning”. In: *ICML*. June 2021.
- [Chr19] P. Christodoulou. “Soft Actor-Critic for discrete action settings”. In: *arXiv [cs.LG]* (Oct. 2019). URL: <http://arxiv.org/abs/1910.07207>.
- [Chu+18] K. Chua, R. Calandra, R. McAllister, and S. Levine. “Deep Reinforcement Learning in a Handful of Trials using Probabilistic Dynamics Models”. In: *NIPS*. 2018. URL: <http://arxiv.org/abs/1805.12114>.
- [Chu+24] R. Chua, A. Ghosh, C. Kaplanis, B. A. Richards, and D. Precup. “Learning successor Features the simple way”. In: *NIPS*. Oct. 2024. URL: <https://arxiv.org/abs/2410.22133>.
- [Chu+25] T. Chu, Y. Zhai, J. Yang, S. Tong, S. Xie, S. Levine, and Y. Ma. “SFT Memorizes, RL Generalizes: A Comparative Study of Foundation Model Post-training”. In: *The Second Conference on Parsimony and Learning*. Mar. 2025. URL: <https://openreview.net/forum?id=d3E3LWmTar>.
- [CL11] O. Chapelle and L. Li. “An empirical evaluation of Thompson sampling”. In: *NIPS*. 2011.
- [CMS07] B. Colson, P. Marcotte, and G. Savard. “An overview of bilevel optimization”. en. In: *Ann. Oper. Res.* 153.1 (Sept. 2007), pp. 235–256. URL: <https://link.springer.com/article/10.1007/s10479-007-0176-2>.
- [Cob+19] K. Cobbe, O. Klimov, C. Hesse, T. Kim, and J. Schulman. “Quantifying Generalization in Reinforcement Learning”. en. In: *ICML*. May 2019, pp. 1282–1289. URL: <https://proceedings.mlr.press/v97/cobbe19a.html>.
- [Cob+20] K. Cobbe, J. Hilton, O. Klimov, and J. Schulman. “Phasic Policy Gradient”. In: *arXiv [cs.LG]* (Sept. 2020). URL: <http://arxiv.org/abs/2009.04416>.
- [Col+22] C. Colas, T. Karch, O. Sigaud, and P.-Y. Oudeyer. “Autotelic agents with intrinsically motivated goal-conditioned reinforcement learning: A short survey”. en. In: *JAIR* 74 (July 2022), pp. 1159–1199. URL: <https://www.jair.org/index.php/jair/article/view/13554>.
- [Col+23] C. Colas, L. Teodorescu, P.-Y. Oudeyer, X. Yuan, and M.-A. Côté. “Augmenting autotelic agents with large language models”. In: May 2023. URL: <https://proceedings.mlr.press/v232/colas23a/colas23a.pdf>.
- [Com+18] G. Comanici, D. Precup, A. Barreto, D. K. Toyama, E. Aygün, P. Hamel, S. Vezhnevets, S. Hou, and S. Mourad. “Knowledge Representation for Reinforcement Learning using General Value Functions”. In: (Sept. 2018). URL: <https://openreview.net/forum?id=rygvZ2RcYm>.
- [Con+24] V. Conitzer et al. “Position: Social Choice should guide AI alignment in dealing with diverse human feedback”. en. In: *ICML* 235 (Apr. 2024), pp. 9346–9360. URL: <https://dl.acm.org/doi/10.5555/3692070.3692441>.
- [CP19] Y. K. Cheung and G. Piliouras. “Vortices instead of equilibria in MinMax optimization: Chaos and butterfly effects of online learning in zero-sum games”. In: *COLT*. May 2019. URL: <https://arxiv.org/abs/1905.08396>.
- [CP20] N. Chopin and O. Papaspiliopoulos. *An Introduction to Sequential Monte Carlo*. en. 1st ed. Springer, 2020.
- [CR+18] J. D. Co-Reyes, Y. Liu, A. Gupta, B. Eysenbach, P. Abbeel, and S. Levine. “Self-consistent trajectory autoencoder: Hierarchical reinforcement learning with trajectory embeddings”. In: *ICML*. June 2018. URL: <https://arxiv.org/abs/1806.02813>.

- [CS04] I. Csiszár and P. C. Shields. “Information theory and statistics: A tutorial”. In: (2004).
- [Csi67] I. Csiszar. “Information-Type Measures of Difference of Probability Distributions and Indirect Observations”. In: *Studia Scientiarum Mathematicarum Hungarica* 2 (1967), pp. 299–318.
- [CSLZ23] W. C. Cheung, D. Simchi-Levi, and R. Zhu. “Nonstationary reinforcement learning: The blessing of (more) optimism”. en. In: *Manage. Sci.* 69.10 (Oct. 2023), pp. 5722–5739. URL: <https://pubsonline.informs.org/doi/abs/10.1287/mnsc.2023.4704>.
- [Dab+17] W. Dabney, M. Rowland, M. G. Bellemare, and R. Munos. “Distributional reinforcement learning with quantile regression”. In: *arXiv [cs.AI]* (Oct. 2017). URL: <http://arxiv.org/abs/1710.10044>.
- [Dab+18] W. Dabney, G. Ostrovski, D. Silver, and R. Munos. “Implicit quantile networks for distributional reinforcement learning”. In: *arXiv [cs.LG]* (June 2018). URL: <http://arxiv.org/abs/1806.06923>.
- [Dai+22] X. Dai et al. “CTRL: Closed-Loop Transcription to an LDR via Minimizing Rate Reduction”. en. In: *Entropy* 24.4 (Mar. 2022). URL: <http://dx.doi.org/10.3390/e24040456>.
- [Dai+24] N. Dainese, M. Merler, M. Alakuijala, and P. Marttinen. “Generating Code World Models with large language models guided by Monte Carlo Tree Search”. In: *NIPS*. May 2024. URL: <https://arxiv.org/abs/2405.15383>.
- [Dan+16] C. Daniel, H. van Hoof, J. Peters, and G. Neumann. “Probabilistic inference for determining options in reinforcement learning”. en. In: *Mach. Learn.* 104.2-3 (Sept. 2016), pp. 337–357. URL: <https://link.springer.com/article/10.1007/s10994-016-5580-x>.
- [Dan+22] I. Danihelka, A. Guez, J. Schrittwieser, and D. Silver. “Policy improvement by planning with Gumbel”. In: *ICLR*. 2022. URL: <https://openreview.net/forum?id=bERaNoegn0>.
- [Day93] P. Dayan. “Improving generalization for temporal difference learning: The successor representation”. en. In: *Neural Comput.* 5.4 (July 1993), pp. 613–624. URL: <https://ieeexplore.ieee.org/abstract/document/6795455>.
- [De+24] S. De et al. “Griffin: Mixing gated linear recurrences with local attention for efficient language models”. In: *arXiv [cs.LG]* (Feb. 2024). URL: <http://arxiv.org/abs/2402.19427>.
- [Dee24] DeepSeek-AI. “DeepSeek-V3 Technical Report”. In: *arXiv [cs.CL]* (Dec. 2024). URL: <http://arxiv.org/abs/2412.19437>.
- [Dee25] G. DeepMind. *AlphaEvolve: A Gemini-powered coding agent for designing advanced algorithms*. en. 2025. URL: <https://deepmind.google/discover/blog/alphaevolve-a-gemini-powered-coding-agent-for-designing-advanced-algorithms/>.
- [Dee25] DeepSeek-AI. “DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning”. In: *arXiv [cs.CL]* (Jan. 2025). URL: <http://arxiv.org/abs/2501.12948>.
- [Den+20] M. Dennis, N. Jaques, E. Vinitzky, A. Bayen, S. Russell, A. Critch, and S. Levine. “Emergent complexity and zero-shot transfer via Unsupervised Environment Design”. In: *NIPS*. Dec. 2020. URL: <https://arxiv.org/abs/2012.02096>.
- [DFR15] M. P. Deisenroth, D. Fox, and C. E. Rasmussen. “Gaussian Processes for Data-Efficient Learning in Robotics and Control”. en. In: *IEEE PAMI* 37.2 (2015), pp. 408–423. URL: <http://dx.doi.org/10.1109/TPAMI.2013.218>.
- [DH92] P. Dayan and G. E. Hinton. “Feudal Reinforcement Learning”. In: *NIPS* 5 (1992). URL: https://proceedings.neurips.cc/paper_files/paper/1992/file/d14220ee66aeec73c49038385428ec4c-Paper.pdf.
- [Die00] T. G. Dietterich. “Hierarchical reinforcement learning with the MAXQ value function decomposition”. en. In: *JAIR* 13 (Nov. 2000), pp. 227–303. URL: <https://www.jair.org/index.php/jair/article/view/10266>.

- [Die+07] M. Diehl, H. G. Bock, H. Diedam, and P.-B. Wieber. “Fast Direct Multiple Shooting Algorithms for Optimal Robot Control”. In: *Lecture Notes in Control and Inform. Sci.* 340 (2007). URL: https://www.researchgate.net/publication/29603798_Fast_Direct_Multiple_Shooting_Algorithms_for_Optimal_Robot_Control.
- [Din+25] J. Ding et al. “Understanding world or predicting future? A comprehensive survey of world models”. en. In: *ACM Comput. Surv.* (June 2025). URL: <https://dl.acm.org/doi/10.1145/3746449>.
- [DMKM22] G. Duran-Martin, A. Kara, and K. Murphy. “Efficient Online Bayesian Inference for Neural Bandits”. In: *AISTATS*. 2022. URL: <http://arxiv.org/abs/2112.00195>.
- [D’O+22] P. D’Oro, M. Schwarzer, E. Nikishin, P.-L. Bacon, M. G. Bellemare, and A. Courville. “Sample-Efficient Reinforcement Learning by Breaking the Replay Ratio Barrier”. In: *Deep Reinforcement Learning Workshop NeurIPS 2022*. Dec. 2022. URL: <https://openreview.net/pdf?id=4GBGwVIEYJ>.
- [DOB21] W. Dabney, G. Ostrovski, and A. Barreto. “Temporally-Extended epsilon-Greedy Exploration”. In: *ICLR*. 2021. URL: <https://openreview.net/pdf?id=ONBPHFZ7zG4>.
- [Don+24] Q. Dong et al. “A survey on in-context learning”. In: *Proceedings of the 2024 Conference on Empirical Methods in Natural Language Processing*. Stroudsburg, PA, USA: Association for Computational Linguistics, Nov. 2024, pp. 1107–1128. URL: <https://aclanthology.org/2024.emnlp-main.64/>.
- [DPA23] F. Deng, J. Park, and S. Ahn. “Facing off world model backbones: RNNs, Transformers, and S4”. In: *NIPS* abs/2307.02064 (July 2023), pp. 72904–72930. URL: https://proceedings.neurips.cc/paper_files/paper/2023/hash/e6c65eb9b56719c1aa45ff73874de317-Abstract-Conference.html.
- [DR11] M. P. Deisenroth and C. E. Rasmussen. “PILCO: A Model-Based and Data-Efficient Approach to Policy Search”. In: *ICML*. 2011. URL: http://www.icml-2011.org/papers/323_icmlpaper.pdf.
- [DR24] V. Dave and E. Rueckert. “Denoised Predictive Imagination: An Information-theoretic approach for learning World Models”. In: *Seventeenth European Workshop on Reinforcement Learning*. Oct. 2024. URL: <https://openreview.net/forum?id=q75NXQdzdJ>.
- [Du+21] C. Du, Z. Gao, S. Yuan, L. Gao, Z. Li, Y. Zeng, X. Zhu, J. Xu, K. Gai, and K.-C. Lee. “Exploration in Online Advertising Systems with Deep Uncertainty-Aware Learning”. In: *KDD. KDD ’21*. Association for Computing Machinery, 2021, pp. 2792–2801. URL: <https://doi.org/10.1145/3447548.3467089>.
- [Du+23a] Y. Du, K. Konyushkova, M. Denil, A. Raju, J. Landon, F. Hill, N. de Freitas, and S. Cabi. “Vision-Language Models as Success Detectors”. In: *The 2nd Conference on Lifelong Learning Agents*. Vol. 232. Proceedings of Machine Learning Research. PMLR, 2023, pp. 120–136. URL: <https://proceedings.mlr.press/v232/du23b.html>.
- [Du+23b] Y. Du, E. Kosoy, A. Dayan, M. Rufova, P. Abbeel, and A. Gopnik. “What can AI Learn from Human Exploration? Intrinsically-Motivated Humans and Agents in Open-World Exploration”. In: *NeurIPS 2023 workshop: Information-Theoretic Principles in Cognitive Systems*. Nov. 2023. URL: <https://openreview.net/pdf?id=5uU9TY2oME>.
- [Du+23c] Y. Du, O. Watkins, Z. Wang, C. Colas, T. Darrell, P. Abbeel, A. Gupta, and J. Andreas. “Guiding pretraining in reinforcement learning with large language models”. In: *ICML*. Feb. 2023. URL: <https://arxiv.org/abs/2302.06692>.
- [Duf02] M. Duff. “Optimal Learning: Computational procedures for Bayes-adaptive Markov decision processes”. PhD thesis. U. Mass. Dept. Comp. Sci., 2002. URL: <http://envy.cs.umass.edu/People/duff/diss.html>.

- [DVRZ22] S. Dong, B. Van Roy, and Z. Zhou. “Simple Agent, Complex Environment: Efficient Reinforcement Learning with Agent States”. In: *J. Mach. Learn. Res.* (2022). URL: <https://www.jmlr.org/papers/v23/21-0773.html>.
- [DWS12] T. Degris, M. White, and R. S. Sutton. “Off-Policy Actor-Critic”. In: *ICML*. 2012. URL: <http://arxiv.org/abs/1205.4839>.
- [Ebe+24] M. Eberhardinger, J. Goodman, A. Dockhorn, D. Perez-Liebana, R. D. Gaina, D. Cakmak, S. Maghsudi, and S. Lucas. “From code to play: Benchmarking program search for games using large language models”. In: *arXiv [cs.AI]* (Dec. 2024). URL: <http://arxiv.org/abs/2412.04057>.
- [Eco+19] A. Ecoffet, J. Huizinga, J. Lehman, K. O. Stanley, and J. Clune. “Go-Explore: a New Approach for Hard-Exploration Problems”. In: (2019). arXiv: [1901.10995 \[cs.LG\]](https://arxiv.org/abs/1901.10995). URL: <http://arxiv.org/abs/1901.10995>.
- [Eco+21] A. Ecoffet, J. Huizinga, J. Lehman, K. O. Stanley, and J. Clune. “First return, then explore”. en. In: *Nature* 590.7847 (Feb. 2021), pp. 580–586. URL: <https://www.nature.com/articles/s41586-020-03157-9>.
- [EGW05] D. Ernst, P. Geurts, and L. Wehenkel. “Tree-based batch mode reinforcement learning”. In: *J. Mach. Learn. Res.* 6.18 (Dec. 2005), pp. 503–556. URL: <https://jmlr.org/papers/v6/ernst05a.html>.
- [Emm+21] S. Emmons, B. Eysenbach, I. Kostrikov, and S. Levine. “RvS: What is essential for offline RL via Supervised Learning?” In: *arXiv [cs.LG]* (Dec. 2021). URL: <http://arxiv.org/abs/2112.10751>.
- [Eno+24] J. Enouen, H. Nakhost, S. Ebrahimi, S. Arik, Y. Liu, and T. Pfister. “TextGenSHAP: Scalable post-hoc explanations in text generation with long documents”. In: *Proc. ACL*. 2024, pp. 13984–14011. URL: <https://aclanthology.org/2024.findings-acl.832/>.
- [Esc+25] A. L. Escoriza, N. Hansen, S. Tao, T. Mu, and H. Su. “Multi-Stage Manipulation with Demonstration-Augmented Reward, Policy, and World Model Learning”. In: *ICML*. June 2025. URL: <https://openreview.net/forum?id=Bv7LUUY0iq&Id=exu8dliVKA>.
- [ESL21] B. Eysenbach, R. Salakhutdinov, and S. Levine. “C-Learning: Learning to Achieve Goals via Recursive Classification”. In: *ICLR*. 2021. URL: <https://openreview.net/pdf?id=tc5qisoB-C>.
- [Esp+18] L. Espeholt et al. “IMPALA: Scalable Distributed Deep-RL with Importance Weighted Actor-Learner Architectures”. en. In: *ICML*. PMLR, July 2018, pp. 1407–1416. URL: <https://proceedings.mlr.press/v80/espeholt18a.html>.
- [Eys+19] B. Eysenbach, A. Gupta, J. Ibarz, and S. Levine. “Diversity is All You Need: Learning Skills without a Reward Function”. In: *ICLR*. 2019. URL: <https://openreview.net/forum?id=SJx63jRqFm>.
- [Eys+20] B. Eysenbach, X. Geng, S. Levine, and R. Salakhutdinov. “Rewriting History with Inverse RL: Hindsight Inference for Policy Improvement”. In: *NIPS*. Feb. 2020.
- [Eys+22] B. Eysenbach, A. Khazatsky, S. Levine, and R. Salakhutdinov. “Mismatched No More: Joint Model-Policy Optimization for Model-Based RL”. In: *NIPS*. 2022. URL: <http://arxiv.org/abs/2110.02758>.
- [Fan+22] L. Fan, G. Wang, Y. Jiang, A. Mandlekar, Y. Yang, H. Zhu, A. Tang, D.-A. Huang, Y. Zhu, and A. Anandkumar. “MineDojo: Building Open-Ended Embodied Agents with Internet-Scale Knowledge”. In: *NIPS*. June 2022. URL: <http://arxiv.org/abs/2206.08853>.
- [Far+18] G. Farquhar, T. Rocktäschel, M. Igl, and S. Whiteson. “TreeQN and ATreeC: Differentiable Tree-Structured Models for Deep Reinforcement Learning”. In: *ICLR*. Feb. 2018. URL: <https://openreview.net/pdf?id=H1dh6Ax0Z>.

- [Far+23a] J. Farebrother, J. Greaves, R. Agarwal, C. Le Lan, R. Goroshin, P. S. Castro, and M. G. Bellemare. “Proto-Value Networks: Scaling Representation Learning with Auxiliary Tasks”. In: *ICLR*. 2023. URL: <https://openreview.net/pdf?id=oGDKSt9JrZi>.
- [Far+23b] G. Farina, J. Grand-Clément, C. Kroer, C.-W. Lee, and H. Luo. “Regret Matching+: (In)Stability and Fast Convergence in Games”. In: *NIPS*. Nov. 2023. URL: <https://openreview.net/forum?id=nYgs0qZJ97ÅñeId=26WFv0oE35>.
- [Far+24] J. Farebrother et al. “Stop regressing: Training value functions via classification for scalable deep RL”. In: *arXiv [cs.LG]* (Mar. 2024). URL: <http://arxiv.org/abs/2403.03950>.
- [Far+25] J. Farebrother, M. Pirotta, A. Tirinzoni, R. Munos, A. Lazaric, and A. Touati. “Temporal Difference Flows”. In: *ICML*. Mar. 2025. URL: <https://arxiv.org/abs/2503.09817>.
- [FC24] J. Farebrother and P. S. Castro. “CALE: Continuous Arcade Learning Environment”. In: *NIPS*. Oct. 2024. URL: <https://arxiv.org/abs/2410.23810>.
- [Fen+25] S. Feng, X. Kong, S. Ma, A. Zhang, D. Yin, C. Wang, R. Pang, and Y. Yang. “Step-by-Step Reasoning for Math Problems via Twisted Sequential Monte Carlo”. In: *ICLR*. 2025. URL: <https://openreview.net/forum?id=Ze4aPP0tIn>.
- [FG21] S. Fujimoto and S. s. Gu. “A Minimalist Approach to Offline Reinforcement Learning”. In: *NIPS*. Vol. 34. Dec. 2021, pp. 20132–20145. URL: https://proceedings.neurips.cc/paper_files/paper/2021/file/a8166da05c5a094f7dc03724b41886e5-Paper.pdf.
- [FHM18] S. Fujimoto, H. van Hoof, and D. Meger. “Addressing Function Approximation Error in Actor-Critic Methods”. In: *ICLR*. 2018. URL: <http://arxiv.org/abs/1802.09477>.
- [FL+18] V. François-Lavet, P. Henderson, R. Islam, M. G. Bellemare, and J. Pineau. “An Introduction to Deep Reinforcement Learning”. In: *Foundations and Trends in Machine Learning* 11.3 (2018). URL: <http://arxiv.org/abs/1811.12560>.
- [FLA16] C. Finn, S. Levine, and P. Abbeel. “Guided Cost Learning: Deep Inverse Optimal Control via Policy Optimization”. In: *ICML*. 2016, pp. 49–58.
- [FLL18] J. Fu, K. Luo, and S. Levine. “Learning Robust Rewards with Adversarial Inverse Reinforcement Learning”. In: *ICLR*. 2018.
- [FMR25] D. J. Foster, Z. Mhammedi, and D. Rohatgi. “Is a Good Foundation Necessary for Efficient Reinforcement Learning? The Computational Role of the Base Model in Exploration”. In: (Mar. 2025). URL: <http://arxiv.org/abs/2503.07453>.
- [For+18] M. Fortunato et al. “Noisy Networks for Exploration”. In: *ICLR*. 2018. URL: <http://arxiv.org/abs/1706.10295>.
- [For+22] S. Forestier, R. Portelas, Y. Mollard, and P.-Y. Oudeyer. “Intrinsically motivated goal exploration processes with automatic curriculum learning”. In: *JMLR* 23.1 (Jan. 2022), pp. 6818–6858. URL: <https://dl.acm.org/doi/10.5555/3586589.3586741>.
- [FPP04] N. Ferns, P. Panangaden, and D. Precup. “Metrics for finite Markov decision processes”. en. In: *UAI*. 2004. URL: <https://dl.acm.org/doi/10.5555/1036843.1036863>.
- [FR23] D. J. Foster and A. Rakhlin. “Foundations of reinforcement learning and interactive decision making”. In: *arXiv [cs.LG]* (Dec. 2023). URL: <http://arxiv.org/abs/2312.16730>.
- [Fra+24a] K. Frans, D. Hafner, S. Levine, and P. Abbeel. “One step diffusion via shortcut models”. In: *arXiv [cs.LG]* (Oct. 2024). URL: <http://arxiv.org/abs/2410.12557>.
- [Fra+24b] B. Frauenknecht, A. Eisele, D. Subhasish, F. Solowjow, and S. Trimpe. “Trust the Model Where It Trusts Itself - Model-Based Actor-Critic with Uncertainty-Aware Rollout Adaption”. In: *ICML*. June 2024. URL: <https://openreview.net/pdf?id=N0ntTjTfHb>.
- [Fre+24] B. Freed, T. Wei, R. Calandra, J. Schneider, and H. Choset. “Unifying Model-Based and Model-Free Reinforcement Learning with Equivalent Policy Sets”. In: *RL Conference*. 2024. URL: https://rlj.cs.umass.edu/2024/papers/RLJ_RLC_2024_37.pdf.

- [Fri03] K. Friston. “Learning and inference in the brain”. en. In: *Neural Netw.* 16.9 (2003), pp. 1325–1352. URL: <http://dx.doi.org/10.1016/j.neunet.2003.06.005>.
- [Fri09] K. Friston. “The free-energy principle: a rough guide to the brain?”. en. In: *Trends Cogn. Sci.* 13.7 (2009), pp. 293–301. URL: <http://dx.doi.org/10.1016/j.tics.2009.04.005>.
- [FS+19] H Francis Song et al. “V-MPO: On-Policy Maximum a Posteriori Policy Optimization for Discrete and Continuous Control”. In: *arXiv [cs.AI]* (Sept. 2019). URL: <http://arxiv.org/abs/1909.12238>.
- [FS25] G. Faria and N. A. Smith. “Sample, don’t search: Rethinking test-time alignment for language models”. In: *arXiv [cs.CL]* (Apr. 2025). URL: <http://arxiv.org/abs/2504.03790>.
- [FSW23] M. Fellows, M. J. A. Smith, and S. Whiteson. “Why Target Networks Stabilise Temporal Difference Methods”. en. In: *ICML*. PMLR, July 2023, pp. 9886–9909. URL: <https://proceedings.mlr.press/v202/fellows23a.html>.
- [Fu15] M. Fu, ed. *Handbook of Simulation Optimization*. 1st ed. Springer-Verlag New York, 2015. URL: <http://www.springer.com/us/book/9781493913831>.
- [Fu+20] J. Fu, A. Kumar, O. Nachum, G. Tucker, and S. Levine. *D4RL: Datasets for Deep Data-Driven Reinforcement Learning*. arXiv:2004.07219. 2020.
- [Fuj+19] S. Fujimoto, E. Conti, M. Ghavamzadeh, and J. Pineau. “Benchmarking batch deep reinforcement learning algorithms”. In: *Deep RL Workshop NeurIPS*. Oct. 2019. URL: <https://arxiv.org/abs/1910.01708>.
- [Fur+21] H. Furuta, T. Kozuno, T. Matsushima, Y. Matsuo, and S. S. Gu. “Co-Adaptation of Algorithmic and Implementational Innovations in Inference-based Deep Reinforcement Learning”. In: *NIPS*. Mar. 2021. URL: <http://arxiv.org/abs/2103.17258>.
- [Gal+24] M. Gallici, M. Fellows, B. Ellis, B. Pou, I. Masmitja, J. N. Foerster, and M. Martin. “Simplifying deep temporal difference learning”. In: *ICML*. July 2024.
- [Gan+25a] K. Gandhi, A. Chakravarthy, A. Singh, N. Lile, and N. D. Goodman. “Cognitive behaviors that enable self-improving reasoners, or, four habits of highly effective STaRs”. In: *arXiv [cs.CL]* (Mar. 2025). URL: <http://arxiv.org/abs/2503.01307>.
- [Gan+25b] K. Gandhi, M. Y. Li, L. Goodyear, L. Li, A. Bhaskar, M. Zaman, and N. D. Goodman. “BoxingGym: Benchmarking progress in automated experimental design and model discovery”. In: *arXiv [cs.LG]* (Jan. 2025). URL: <http://arxiv.org/abs/2501.01540>.
- [Gao+25] Y. Gao et al. “Seedance 1.0: Exploring the boundaries of video generation models”. In: *arXiv [cs.CV]* (June 2025). URL: <http://arxiv.org/abs/2506.09113>.
- [Gar23] R. Garnett. *Bayesian Optimization*. Cambridge University Press, 2023. URL: <https://bayesoptbook.com/>.
- [Gar+23] Q. Garrido, Y. Chen, A. Bardes, L. Najman, and Y. LeCun. “On the duality between contrastive and non-contrastive self-supervised learning”. In: *ICLR*. 2023. URL: <https://openreview.net/forum?id=kDEL91Dufpa>.
- [Gar+24] Q. Garrido, M. Assran, N. Ballas, A. Bardes, L. Najman, and Y. LeCun. “Learning and leveraging world models in visual representation learning”. In: *arXiv [cs.CV]* (Mar. 2024). URL: <http://arxiv.org/abs/2403.00504>.
- [Gar+25] S. Garcin, T. McInroe, P. S. Castro, C. G. Lucas, D. Abel, P. Panangaden, and S. V. Albrecht. “Studying the Interplay Between the Actor and Critic Representations in Reinforcement Learning”. In: *ICLR*. 2025. URL: <https://openreview.net/forum?id=tErHYBG1Wc>.
- [GBS22] C. Grimm, A. Barreto, and S. Singh. “Approximate Value Equivalence”. In: *NIPS*. Oct. 2022. URL: <https://openreview.net/pdf?id=S2Awu3Zn04v>.

- [GD22] S. Gronauer and K. Diepold. “Multi-agent deep reinforcement learning: a survey”. en. In: *Artif. Intell. Rev.* 55.2 (Feb. 2022), pp. 895–943. URL: <https://dl.acm.org/doi/10.1007/s10462-021-09996-w>.
- [GD23] A. Gu and T. Dao. “Mamba: Linear-Time Sequence Modeling with Selective State Spaces”. In: *arXiv [cs.LG]* (Dec. 2023). URL: <http://arxiv.org/abs/2312.00752>.
- [GDG03] R. Givan, T. Dean, and M. Greig. “Equivalence notions and model minimization in Markov decision processes”. en. In: *Artif. Intell.* 147.1-2 (July 2003), pp. 163–223. URL: <https://www.sciencedirect.com/science/article/pii/S0004370202003764>.
- [GDWF22] J. Grudzien, C. A. S. De Witt, and J. Foerster. “Mirror Learning: A Unifying Framework of Policy Optimisation”. In: *ICML*. Vol. 162. Proceedings of Machine Learning Research. PMLR, 2022, pp. 7825–7844. URL: <https://proceedings.mlr.press/v162/grudzien22a/grudzien22a.pdf>.
- [Ge+24] L. Ge, D. Halpern, E. Micha, A. D. Procaccia, I. Shapira, Y. Vorobeychik, and J. Wu. “Axioms for AI alignment from human feedback”. In: *arXiv [cs.GT]* (May 2024). URL: <http://arxiv.org/abs/2405.14758>.
- [Ger18] S. J. Gershman. “Deconstructing the human algorithms for exploration”. en. In: *Cognition* 173 (Apr. 2018), pp. 34–42. URL: <https://www.sciencedirect.com/science/article/abs/pii/S0010027717303359>.
- [Ger19] S. J. Gershman. “What does the free energy principle tell us about the brain?” In: *Neurons, Behavior, Data Analysis, and Theory* (2019). URL: <http://arxiv.org/abs/1901.07945>.
- [GF21] A. Grim and P. Felzenszwalb. “Convex Combination Belief Propagation Algorithms”. In: (May 2021). arXiv: 2105.12815 [cs.AI]. URL: <http://arxiv.org/abs/2105.12815>.
- [GFZ24] J. Grigsby, L. Fan, and Y. Zhu. “AMAGO: Scalable In-Context Reinforcement Learning for Adaptive Agents”. In: *ICLR*. 2024. URL: <https://openreview.net/forum?id=M6XWoEdmwf>.
- [GG24] M. K. Goddu and A. Gopnik. “The development of human causal learning and reasoning”. en. In: *Nat. Rev. Psychol.* 3.5 (Apr. 2024), pp. 319–339. URL: <https://www.nature.com/articles/s44159-024-00300-5>.
- [GGL22] D. Graves, J. Günther, and J. Luo. “Affordance as general value function: a computational model”. en. In: *Adapt. Behav.* 30.4 (Aug. 2022), pp. 307–327. URL: https://journals.sagepub.com/doi/full/10.1177/1059712321999421?casa_token=FVDN6RCZ6qMAAAAA%3ApmSzZw3cduXIQontJxYIH5uPQNNaS4Q1tmD8VWxDe1kUhZWwdJ1Hga4eFASQkjaMkkm5VQcIGGUi.
- [GGN22] S. K. S. Ghasemipour, S. S. Gu, and O. Nachum. “Why So Pessimistic? Estimating Uncertainties for Offline RL through Ensembles, and Why Their Independence Matters”. In: *NIPS*. Oct. 2022. URL: <https://openreview.net/pdf?id=z64kN1h1-rR>.
- [GGV24] L. Gui, C. Gârbacea, and V. Veitch. “BoNBon Alignment for large language models and the sweetness of best-of-n sampling”. In: *arXiv [cs.CL]* (June 2024). URL: <http://arxiv.org/abs/2406.00832>.
- [Gha+15] M. Ghavamzadeh, S. Mannor, J. Pineau, and A. Tamar. “Bayesian Reinforcement Learning: A Survey”. en. In: *Found. Trends® Mach. Learn.* 8.5-6 (Nov. 2015), pp. 359–483. URL: <https://arxiv.org/abs/1609.04436>.
- [Ghi+20] S. Ghiassian, A. Patterson, S. Garg, D. Gupta, A. White, and M. White. “Gradient temporal-difference learning with Regularized Corrections”. In: *ICML*. July 2020.
- [Gho+21] D. Ghosh, J. Rahme, A. Kumar, A. Zhang, R. P. Adams, and S. Levine. “Why Generalization in RL is Difficult: Epistemic POMDPs and Implicit Partial Observability”. In: *NIPS*. Vol. 34. Dec. 2021, pp. 25502–25515. URL: https://proceedings.neurips.cc/paper_files/paper/2021/file/d5ff135377d39f1de7372c95c74dd962-Paper.pdf.

- [Ghu+22] R. Ghugare, H. Bharadhwaj, B. Eysenbach, S. Levine, and R. Salakhutdinov. “Simplifying Model-based RL: Learning Representations, Latent-space Models, and Policies with One Objective”. In: *ICLR*. Sept. 2022. URL: <https://openreview.net/forum?id=MQcmfgRxf7a>.
- [Ghu+24] R. Ghugare, M. Geist, G. Berseth, and B. Eysenbach. “Closing the gap between TD learning and supervised learning – A generalisation Point of View”. In: *ICML*. Jan. 2024. URL: <https://arxiv.org/abs/2401.11237>.
- [Gil+25] N. Gillman, D. Aggarwal, M. Freeman, S. Singh, and C. Sun. “Fourier head: Helping large language models learn complex probability distributions”. In: *ICLR*. 2025. URL: <https://arxiv.org/abs/2410.22269>.
- [Git89] J. Gittins. *Multi-armed Bandit Allocation Indices*. Wiley, 1989.
- [GK19] L. Graesser and W. L. Keng. *Foundations of Deep Reinforcement Learning: Theory and Practice in Python*. en. 1 edition. Addison-Wesley Professional, 2019. URL: <https://www.amazon.com/Deep-Reinforcement-Learning-Python-Hands/dp/0135172381>.
- [GM+24] J. Grau-Moya et al. “Learning Universal Predictors”. In: *arXiv [cs.LG]* (Jan. 2024). URL: <https://arxiv.org/abs/2401.14953>.
- [Gon+24] X. Gong, D. Feng, K. Xu, B. Ding, and H. Wang. “Goal-conditioned on-policy reinforcement learning”. In: *NIPS* 37 (2024), pp. 45975–46001. URL: https://proceedings.neurips.cc/paper_files/paper/2024/hash/51c6e143b5da2bd6e4a618d8a5d7f38b-Abstract-Conference.html.
- [Gor95] G. J. Gordon. “Stable Function Approximation in Dynamic Programming”. In: *ICML*. 1995, pp. 261–268.
- [GP23] M. Gerstgrasser and D. C. Parkes. “Oracles & Followers: Stackelberg Equilibria in Deep Multi-Agent Reinforcement Learning”. In: *ICML*. Vol. 202. Proceedings of Machine Learning Research. PMLR, 2023, pp. 11213–11236. URL: <https://proceedings.mlr.press/v202/gerstgrasser23a.html>.
- [Gra+10] T. Graepel, J. Quinonero-Candela, T. Borchert, and R. Herbrich. “Web-Scale Bayesian Click-Through Rate Prediction for Sponsored Search Advertising in Microsoft’s Bing Search Engine”. In: *ICML*. 2010.
- [Gri+20a] J.-B. Grill et al. “Bootstrap your own latent: A new approach to self-supervised Learning”. In: *NIPS*. June 2020. URL: <http://arxiv.org/abs/2006.07733>.
- [Gri+20b] C. Grimm, A. Barreto, S. Singh, and D. Silver. “The Value Equivalence Principle for Model-Based Reinforcement Learning”. In: *NIPS* 33 (2020), pp. 5541–5552. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/3bb585ea00014b0e3ebe4c6dd165a358-Paper.pdf.
- [Gri+24] J. Grigsby, J. Sasek, S. Parajuli, D. Adeb, A. Zhang, and Y. Zhu. “AMAGO-2: Breaking the Multi-Task Barrier in Meta-Reinforcement Learning with Transformers”. In: *NIPS*. 2024. URL: <https://openreview.net/forum?id=OSHARf4TVU>.
- [GRW17] K. Gregor, D. J. Rezende, and D. Wierstra. “Variational intrinsic control”. In: *ICLR Workshop*. 2017. URL: <https://openreview.net/pdf?id=Skc-Fo4Yg>.
- [Gua+23] L. Guan, K. Valmeekam, S. Sreedharan, and S. Kambhampati. “Leveraging Pre-trained Large Language Models to Construct and Utilize World Models for Model-based Task Planning”. In: *NIPS*. May 2023. URL: <http://arxiv.org/abs/2305.14909>.
- [Gua+24] Z. Guan, X. Kong, F. Zhong, and Y. Wang. “Richelieu: Self-Evolving LLM-Based Agents for AI Diplomacy”. In: *NIPS*. Nov. 2024. URL: <https://openreview.net/forum?id=7Jb4NJS8Yk&Id=BBHQNOu0h>.
- [Gul+20] C. Gulcehre et al. *RL Unplugged: Benchmarks for Offline Reinforcement Learning*. arXiv:2006.13888. 2020.

- [Gul25] A. Gulli. *Agentic design patterns: A hands-on guide to building intelligent systems*. en. Springer, Nov. 2025. URL: https://www.google.com/books/edition/Agentic_Design_Patterns/QqR20QEACAAJ?hl=en.
- [Guo+22] Z. D. Guo et al. “BYOL-Explore: Exploration by Bootstrapped Prediction”. In: *NIPS*. Oct. 2022. URL: <https://openreview.net/pdf?id=qHGCH75usg>.
- [GZG19] S. K. S. Ghasemipour, R. S. Zemel, and S. Gu. “A Divergence Minimization Perspective on Imitation Learning Methods”. In: *CORL*. 2019, pp. 1259–1277.
- [Haa+18a] T. Haarnoja, A. Zhou, P. Abbeel, and S. Levine. “Soft Actor-Critic: Off-Policy Maximum Entropy Deep Reinforcement Learning with a Stochastic Actor”. In: *ICML*. 2018. URL: <http://arxiv.org/abs/1801.01290>.
- [Haa+18b] T. Haarnoja et al. “Soft Actor-Critic Algorithms and Applications”. In: (2018). arXiv: 1812.05905 [cs.LG]. URL: <http://arxiv.org/abs/1812.05905>.
- [Haf+19] D. Hafner, T. Lillicrap, I. Fischer, R. Villegas, D. Ha, H. Lee, and J. Davidson. “Learning Latent Dynamics for Planning from Pixels”. In: *ICML*. 2019. URL: <http://arxiv.org/abs/1811.04551>.
- [Haf+20] D. Hafner, T. Lillicrap, J. Ba, and M. Norouzi. “Dream to Control: Learning Behaviors by Latent Imagination”. In: *ICLR*. 2020. URL: <https://openreview.net/forum?id=S110TC4tDS>.
- [Haf+21] D. Hafner, T. Lillicrap, M. Norouzi, and J. Ba. “Mastering Atari with discrete world models”. In: *ICLR*. 2021.
- [Haf+25] D. Hafner, J. Pasukonis, J. Ba, and T. Lillicrap. “Mastering diverse control tasks through world models”. en. In: *Nature* 640.8059 (Apr. 2025), pp. 647–653. URL: <https://www.nature.com/articles/s41586-025-08744-2>.
- [Ham+21] J. B. Hamrick, A. L. Friesen, F. Behbahani, A. Guez, F. Viola, S. Witherspoon, T. Anthony, L. Buesing, P. Veličković, and T. Weber. “On the role of planning in model-based deep reinforcement learning”. In: *ICLR*. 2021. URL: <https://arxiv.org/abs/2011.04021>.
- [Ham+25] L. Hammond et al. “Multi-Agent Risks from Advanced AI”. In: *arXiv [cs.MA]* (Feb. 2025). URL: <http://arxiv.org/abs/2502.14143>.
- [Han+19] S. Hansen, W. Dabney, A. Barreto, D. Warde-Farley, T. Van de Wiele, and V. Mnih. “Fast Task Inference with Variational Intrinsic Successor Features”. In: *ICLR*. Sept. 2019. URL: <https://openreview.net/pdf?id=BJeAHkrYDS>.
- [Han+23] N. Hansen, Z. Yuan, Y. Ze, T. Mu, A. Rajeswaran, H. Su, H. Xu, and X. Wang. “On Pre-Training for Visuo-Motor Control: Revisiting a Learning-from-Scratch Baseline”. In: *ICML*. June 2023. URL: <https://openreview.net/pdf?id=dvp30Hrijj>.
- [Han+24] S. Han, Q. Zhang, Y. Yao, W. Jin, and Z. Xu. “LLM Multi-Agent Systems: Challenges and Open Problems”. In: *arXiv [cs.MA]* (Feb. 2024). URL: <http://arxiv.org/abs/2402.03578>.
- [Han+25] N. Hansen, S. V. Jyothir, V. Sobal, Y. LeCun, X. Wang, and H. Su. “Hierarchical World Models as Visual Whole-Body Humanoid Controllers”. In: *ICLR*. 2025. URL: <https://openreview.net/forum?id=7wuJMvK639>.
- [Hao+25] Q. Hao et al. “Reinforcement Learning in the Era of Large Language Models: Challenges and Opportunities”. In: (2025). URL: https://www.researchgate.net/profile/Qianyue-Hao/publication/395648618_Reinforcement_Learning_in_the_Era_of_Large_Language_Models_Challenges_and_Opportunities/links/68ccbc8911d348252ba644cc/Reinforcement-Learning-in-the-Era-of-Large-Language-Models-Challenges-and-Opportunities.pdf.
- [Har+16] A. Harutyunyan, M. G. Bellemare, T. Stepleton, and R. Munos. “Q(λ) with Off-Policy Corrections”. In: (2016).

- [Har+18] J. Harb, P.-L. Bacon, M. Klissarov, and D. Precup. “When waiting is not an option: Learning options with a deliberation cost”. en. In: *AAAI* 32.1 (Apr. 2018). URL: <https://ojs.aaai.org/index.php/AAAI/article/view/11831>.
- [Has10] H. van Hasselt. “Double Q-learning”. In: *NIPS*. Ed. by J. D. Lafferty, C. K. I. Williams, J. Shawe-Taylor, R. S. Zemel, and A. Culotta. Curran Associates, Inc., 2010, pp. 2613–2621. URL: <http://papers.nips.cc/paper/3964-double-q-learning.pdf>.
- [Has+16] H. van Hasselt, A. Guez, M. Hessel, V. Mnih, and D. Silver. “Learning values across many orders of magnitude”. In: *NIPS*. Feb. 2016.
- [Haz19] E. Hazan. *Introduction to online convex optimization*. MIT Press, Sept. 2019. URL: <https://arxiv.org/abs/1909.05207>.
- [HBZ04] E. A. Hansen, D. S. Bernstein, and S. Zilberstein. “Dynamic programming for partially observable stochastic games”. en. In: *Proceedings of the 19th national conference on Artificial intelligence*. AAAI’04. AAAI Press, July 2004, pp. 709–715. URL: <https://dl.acm.org/doi/10.5555/1597148.1597262>.
- [HDCM15] A. Hallak, D. Di Castro, and S. Mannor. “Contextual Markov decision processes”. In: *arXiv [stat.ML]* (Feb. 2015). URL: <http://arxiv.org/abs/1502.02259>.
- [HE16] J. Ho and S. Ermon. “Generative Adversarial Imitation Learning”. In: *NIPS*. 2016, pp. 4565–4573.
- [Hee+15] N. Heess, G. Wayne, D. Silver, T. Lillicrap, T. Erez, and Y. Tassa. “Learning Continuous Control Policies by Stochastic Value Gradients”. In: *NIPS* 28 (2015). URL: <https://proceedings.neurips.cc/paper/2015/hash/148510031349642de5ca0c544f31b2ef-Abstract.html>.
- [Hes+18] M. Hessel, J. Modayil, H. van Hasselt, T. Schaul, G. Ostrovski, W. Dabney, D. Horgan, B. Piot, M. Azar, and D. Silver. “Rainbow: Combining Improvements in Deep Reinforcement Learning”. In: *AAAI*. 2018. URL: <http://arxiv.org/abs/1710.02298>.
- [Hes+19] M. Hessel, H. Soyer, L. Espeholt, W. Czarnecki, S. Schmitt, and H. van Hasselt. “Multi-task deep reinforcement learning with PopArt”. In: *AAAI*. 2019.
- [HGS16] H. van Hasselt, A. Guez, and D. Silver. “Deep Reinforcement Learning with Double Q-Learning”. In: *AAAI*. AAAI’16. AAAI Press, 2016, pp. 2094–2100. URL: <http://dl.acm.org/citation.cfm?id=3016100.3016191>.
- [HHA19] H. van Hasselt, M. Hessel, and J. Aslanides. “When to use parametric models in reinforcement learning?” In: *NIPS*. 2019. URL: <http://arxiv.org/abs/1906.05243>.
- [HL04] D. R. Hunter and K. Lange. “A Tutorial on MM Algorithms”. In: *The American Statistician* 58 (2004), pp. 30–37.
- [HL20] O. van der Himst and P. Lanillos. “Deep active inference for partially observable MDPs”. In: *ECML workshop on active inference*. Sept. 2020. URL: <https://arxiv.org/abs/2009.03622>.
- [HLKT19] P. Hernandez-Leal, B. Kartal, and M. E. Taylor. “A survey and critique of multiagent deep reinforcement learning”. In: *Auton. Agent. Multi. Agent. Syst.* (2019). URL: <http://link.springer.com/10.1007/s10458-019-09421-1>.
- [HLS15] J. Heinrich, M. Lanctot, and D. Silver. “Fictitious Self-Play in Extensive-Form Games”. In: *ICML*. 2015.
- [HM20] M. Hosseini and A. Maida. “Hierarchical Predictive Coding Models in a Deep-Learning Framework”. In: (2020). arXiv: 2005.03230 [cs.CV]. URL: <http://arxiv.org/abs/2005.03230>.
- [HMC00] S. Hart and A. Mas-Colell. “A simple adaptive procedure leading to correlated equilibrium”. en. In: *Econometrica* 68.5 (Sept. 2000), pp. 1127–1150. URL: <https://www.jstor.org/stable/2999445>.
- [HMDH20] C. C.-Y. Hsu, C. Mendler-Dünner, and M. Hardt. “Revisiting design choices in proximal Policy Optimization”. In: *arXiv [cs.LG]* (Sept. 2020). URL: <http://arxiv.org/abs/2009.10897>.

- [Hof+23] M. D. Hoffman, D. Phan, D. Dohan, S. Douglas, T. A. Le, A. Parisi, P. Sountsov, C. Sutton, S. Vikram, and R. A. Saurous. “Training Chain-of-Thought via Latent-Variable Inference”. In: *NIPS*. Jan. 2023. URL: <https://openreview.net/forum?id=7p1t0Z13La>.
- [Hon+10] A. Honkela, T. Raiko, M. Kuusela, M. Tornio, and J. Karhunen. “Approximate Riemannian Conjugate Gradient Learning for Fixed-Form Variational Bayes”. In: *JMLR* 11. Nov (2010), pp. 3235–3268. URL: <http://www.jmlr.org/papers/volume11/honkela10a/honkela10a.pdf>.
- [Hon+23] M. Hong, H.-T. Wai, Z. Wang, and Z. Yang. “A two-timescale stochastic algorithm framework for bilevel optimization: Complexity analysis and application to actor-critic”. en. In: *SIAM J. Optim.* 33.1 (Mar. 2023), pp. 147–180. URL: <https://epubs.siam.org/doi/10.1137/20M1387341>.
- [Hou+11] N. Houlsby, F. Huszár, Z. Ghahramani, and M. Lengyel. “Bayesian active learning for classification and preference learning”. In: *arXiv [stat.ML]* (Dec. 2011). URL: <http://arxiv.org/abs/1112.5745>.
- [HQC24] M. Hutter, D. Quarel, and E. Catt. *An introduction to universal artificial intelligence*. Chapman and Hall, 2024. URL: <http://www.hutter1.net/ai/uaibook2.htm>.
- [HR11] R. Hafner and M. Riedmiller. “Reinforcement learning in feedback control: Challenges and benchmarks from technical process control”. en. In: *Mach. Learn.* 84.1-2 (July 2011), pp. 137–169. URL: <https://link.springer.com/article/10.1007/s10994-011-5235-x>.
- [HR17] C. Hoffmann and P. Rostalski. “Linear Optimal Control on Factor Graphs — A Message Passing Perspective”. In: *Intl. Federation of Automatic Control* 50.1 (2017), pp. 6314–6319. URL: <https://www.sciencedirect.com/science/article/pii/S2405896317313800>.
- [HS16] J. Heinrich and D. Silver. “Deep reinforcement learning from self-play in imperfect-information games”. In: *arXiv [cs.LG]* (Mar. 2016). URL: <http://arxiv.org/abs/1603.01121>.
- [HS18] D. Ha and J. Schmidhuber. “World Models”. In: *NIPS*. 2018. URL: <http://arxiv.org/abs/1803.10122>.
- [HS22] E. Hazan and K. Singh. “Introduction to online control”. In: *arXiv [cs.LG]* (Nov. 2022). URL: <http://arxiv.org/abs/2211.09619>.
- [HSW22] N. A. Hansen, H. Su, and X. Wang. “Temporal Difference Learning for Model Predictive Control”. en. In: *ICML*. June 2022, pp. 8387–8406. URL: <https://proceedings.mlr.press/v162/hansen22a.html>.
- [HSW24] N. Hansen, H. Su, and X. Wang. “TD-MPC2: Scalable, Robust World Models for Continuous Control”. In: *ICLR*. 2024. URL: <http://arxiv.org/abs/2310.16828>.
- [HT15] J. H. Huggins and J. B. Tenenbaum. “Risk and regret of hierarchical Bayesian learners”. In: *ICML*. 2015. URL: <http://proceedings.mlr.press/v37/hugginsb15.html>.
- [HTB18] G. Z. Holland, E. J. Talvitie, and M. Bowling. “The effect of planning shape on Dyna-style planning in high-dimensional state spaces”. In: *arXiv [cs.AI]* (June 2018). URL: <http://arxiv.org/abs/1806.01825>.
- [Hu+20] Y. Hu, W. Wang, H. Jia, Y. Wang, Y. Chen, J. Hao, F. Wu, and C. Fan. “Learning to Utilize Shaping Rewards: A New Approach of Reward Shaping”. In: *NIPS* 33 (2020), pp. 15931–15941. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/b710915795b9e9c02cf10d6d2bdb688c-Paper.pdf.
- [Hu+21] H. Hu, A. Lerer, N. Brown, and J. Foerster. “Learned Belief Search: Efficiently improving policies in partially observable settings”. In: *arXiv [cs.AI]* (June 2021). URL: <http://arxiv.org/abs/2106.09086>.
- [Hu+23a] A. Hu, L. Russell, H. Yeo, Z. Murez, G. Fedoseev, A. Kendall, J. Shotton, and G. Corrado. “GAIA-1: A Generative World Model for Autonomous Driving”. In: *arXiv [cs.CV]* (Sept. 2023). URL: <http://arxiv.org/abs/2309.17080>.

- [Hu+23b] E. S. Hu, R. Chang, O. Rybkin, and D. Jayaraman. “Planning Goals for Exploration”. In: *ICLR*. Mar. 2023. URL: <https://arxiv.org/abs/2303.13002>.
- [Hu+24] S. Hu, T. Huang, F. Ilhan, S. Tekin, G. Liu, R. Kompella, and L. Liu. *A Survey on Large Language Model-Based Game Agents*. 2024. arXiv: [2404.02039](https://arxiv.org/abs/2404.02039) [cs.AI].
- [Hu+25] J. Hu, J. K. Liu, H. Xu, and W. Shen. “REINFORCE++: An efficient RLHF algorithm with robustness to both prompt and reward models”. In: *arXiv [cs.CL]* (Aug. 2025). URL: <http://arxiv.org/abs/2501.03262>.
- [Hua+23] S. Huang et al. “Language is not all you need: Aligning perception with language models”. In: *arXiv [cs.CL]* (Feb. 2023). URL: <http://arxiv.org/abs/2302.14045>.
- [Hua+24a] S. Huang, M. Noukhovitch, A. Hosseini, K. Rasul, W. Wang, and L. Tunstall. “The N+ Implementation Details of RLHF with PPO: A Case Study on TL;DR Summarization”. In: *First Conference on Language Modeling*. Aug. 2024. URL: <https://openreview.net/pdf?id=kH02ZTa8e3>.
- [Hua+24b] S. Huang, J. Hu, Z. Yang, L. Yang, T. Luo, H. Chen, L. Sun, and B. Yang. “Decision Mamba: Reinforcement Learning via Hybrid Selective Sequence Modeling”. In: *NIPS*. Nov. 2024. URL: [https://openreview.net/forum?id=wFzIMbTsY7&referrer=%5Bthe%20profile%20of%20Tao%20Luo%5D\(%2Fprofile%3Fid%3D%5E~Tao_Luo2\)](https://openreview.net/forum?id=wFzIMbTsY7&referrer=%5Bthe%20profile%20of%20Tao%20Luo%5D(%2Fprofile%3Fid%3D%5E~Tao_Luo2)).
- [Hua+25a] A. Huang, A. Block, Q. Liu, N. Jiang, A. Krishnamurthy, and D. J. Foster. “Is best-of-N the best of them? Coverage, scaling, and optimality in inference-time alignment”. In: *arXiv [cs.AI]* (Mar. 2025). URL: <http://arxiv.org/abs/2503.21878>.
- [Hua+25b] A. Huang, W. Zhan, T. Xie, J. D. Lee, W. Sun, A. Krishnamurthy, and D. J. Foster. “Correcting the Mythos of KL-Regularization: Direct Alignment without Overoptimization via Chi-Squared Preference Optimization”. In: *ICLR*. 2025. URL: <https://openreview.net/forum?id=hXmOWu2U9K>.
- [Hub+21] T. Hubert, J. Schrittwieser, I. Antonoglou, M. Barekatin, S. Schmitt, and D. Silver. “Learning and planning in complex action spaces”. In: *arXiv [cs.LG]* (Apr. 2021). URL: <http://arxiv.org/abs/2104.06303>.
- [Hui+25] Z. Huichi et al. “Memento: Fine-tuning LLM agents without fine-tuning LLMs”. In: *arXiv [cs.LG]* (Aug. 2025). URL: <http://arxiv.org/abs/2508.16153>.
- [Hut05] M. Hutter. *Universal Artificial Intelligence: Sequential Decisions Based On Algorithmic Probability*. en. 2005th ed. Springer, 2005. URL: <http://www.hutter1.net/ai/uaibook.htm>.
- [HYL25] D. Hafner, W. Yan, and T. Lillicrap. “Training agents inside of scalable world models”. In: *arXiv [cs.AI]* (Sept. 2025). URL: <http://arxiv.org/abs/2509.24527>.
- [Ich+23] B. Ichter et al. “Do As I Can, Not As I Say: Grounding Language in Robotic Affordances”. en. In: *Conference on Robot Learning*. PMLR, Mar. 2023, pp. 287–318. URL: <https://proceedings.mlr.press/v205/ichter23a.html>.
- [ID19] S. Ivanov and A. D’yakonov. “Modern Deep Reinforcement Learning algorithms”. In: *arXiv [cs.LG]* (June 2019). URL: <http://arxiv.org/abs/1906.10025>.
- [IW18] E. Imani and M. White. “Improving Regression Performance with Distributional Losses”. en. In: *ICML*. PMLR, July 2018, pp. 2157–2166. URL: <https://proceedings.mlr.press/v80/imani18a.html>.
- [Jac+25] M. T. Jackson, U. Berdica, J. Liesen, S. Whiteson, and J. N. Foerster. “A Clean Slate for Offline Reinforcement Learning”. In: (Apr. 2025). URL: <http://arxiv.org/abs/2504.11453>.
- [Jad+17] M. Jaderberg, V. Mnih, W. M. Czarnecki, T. Schaul, J. Z. Leibo, D. Silver, and K. Kavukcuoglu. “Reinforcement Learning with Unsupervised Auxiliary Tasks”. In: *ICLR*. 2017. URL: <https://openreview.net/forum?id=SJ6yPD5xg>.

- [Jad+19] M. Jaderberg et al. “Human-level performance in 3D multiplayer games with population-based reinforcement learning”. en. In: *Science* 364.6443 (May 2019), pp. 859–865. URL: <https://www.science.org/doi/10.1126/science.aau6249>.
- [Jae00] H Jaeger. “Observable operator models for discrete stochastic time series”. en. In: *Neural Comput.* 12.6 (June 2000), pp. 1371–1398. URL: <https://direct.mit.edu/neco/article-pdf/12/6/1371/814514/089976600300015411.pdf>.
- [Jan+19a] M. Janner, J. Fu, M. Zhang, and S. Levine. “When to Trust Your Model: Model-Based Policy Optimization”. In: *NIPS*. 2019. URL: <http://arxiv.org/abs/1906.08253>.
- [Jan+19b] D. Janz, J. Hron, P. Mazur, K. Hofmann, J. M. Hernández-Lobato, and S. Tschitschek. “Successor Uncertainties: Exploration and Uncertainty in Temporal Difference Learning”. In: *NIPS*. Vol. 32. 2019. URL: https://proceedings.neurips.cc/paper_files/paper/2019/file/1b113258af3968aaf3969ca67e744ff8-Paper.pdf.
- [Jan+22] M. Janner, Y. Du, J. B. Tenenbaum, and S. Levine. “Planning with Diffusion for Flexible Behavior Synthesis”. In: *ICML*. May 2022. URL: <http://arxiv.org/abs/2205.09991>.
- [Jar+23] D. Jarrett, C. Tallec, F. Althé, T. Mesnard, R. Munos, and M. Valko. “Curiosity in Hindsight: Intrinsic Exploration in Stochastic Environments”. In: *ICML*. June 2023. URL: <https://openreview.net/pdf?id=fIH2G4fnSy>.
- [JCM24] M. Jones, P. Chang, and K. Murphy. “Bayesian online natural gradient (BONG)”. In: May 2024. URL: <http://arxiv.org/abs/2405.19681>.
- [JG24] B. Jaeger and A. Geiger. “An Invitation to Deep Reinforcement Learning”. In: *Found. trends optim.* (2024). URL: <http://arxiv.org/abs/2312.08365>.
- [JGP16] E. Jang, S. Gu, and B. Poole. “Categorical Reparameterization with Gumbel-Softmax”. In: (2016). arXiv: 1611.01144 [stat.ML]. URL: <http://arxiv.org/abs/1611.01144>.
- [Ji+25] Y. Ji, J. Li, H. Ye, K. Wu, K. Yao, J. Xu, L. Mo, and M. Zhang. “Test-time compute: From System-1 thinking to System-2 thinking”. In: *arXiv [cs.AI]* (Jan. 2025). URL: <http://arxiv.org/abs/2501.02497>.
- [Jia+15] N. Jiang, A. Kulesza, S. Singh, and R. Lewis. “The Dependence of Effective Planning Horizon on Model Accuracy”. en. In: *Proceedings of the 2015 International Conference on Autonomous Agents and Multiagent Systems*. AAMAS ’15. Richland, SC: International Foundation for Autonomous Agents and Multiagent Systems, May 2015, pp. 1181–1189. URL: <https://dl.acm.org/doi/10.5555/2772879.2773300>.
- [Jin+19] Y. Jinnai, J. W. Park, D. Abel, and G. Konidaris. “Discovering Options for Exploration by Minimizing Cover Time”. en. In: *ICML*. PMLR, May 2019, pp. 3130–3139. URL: <https://proceedings.mlr.press/v97/jinnai19b.html>.
- [Jin+20] Y. Jinnai, J. W. Park, M. C. Machado, and G. Konidaris. “Exploration in Reinforcement Learning with Deep Covering Options”. In: *ICLR*. 2020. URL: <https://openreview.net/forum?id=SkeIyaVtwB>.
- [Jin+22] L. Jing, P. Vincent, Y. LeCun, and Y. Tian. “Understanding Dimensional Collapse in Contrastive Self-supervised Learning”. In: *ICLR*. 2022. URL: <https://openreview.net/forum?id=YevsQ05DEN7>.
- [Jin+25] B. Jin, H. Zeng, Z. Yue, J. Yoon, S. O. Arik, D. Wang, H. Zamani, and J. Han. “Search-R1: Training LLMs to Reason and Leverage Search Engines with Reinforcement Learning”. In: *Second Conference on Language Modeling*. Aug. 2025. URL: <https://openreview.net/forum?id=Rwhi91ideu>.
- [JLL21] M. Janner, Q. Li, and S. Levine. “Offline Reinforcement Learning as One Big Sequence Modeling Problem”. In: *NIPS*. June 2021.
- [JM70] D. H. Jacobson and D. Q. Mayne. *Differential Dynamic Programming*. Elsevier Press, 1970.

- [JML20] M. Janner, I. Mordatch, and S. Levine. “Gamma-Models: Generative Temporal Difference Learning for Infinite-Horizon Prediction”. In: *NIPS*. Vol. 33. 2020, pp. 1724–1735. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/12ffb0968f2f56e51a59a6beb37b2859-Paper.pdf.
- [JOA10] T. Jaksch, R. Ortner, and P. Auer. “Near-optimal regret bounds for reinforcement learning”. In: *JMLR* 11.51 (2010), pp. 1563–1600. URL: <https://jmlr.org/papers/v11/jaksch10a.html>.
- [Jor+24] S. M. Jordan, A. White, B. C. da Silva, M. White, and P. S. Thomas. “Position: Benchmarking is Limited in Reinforcement Learning Research”. In: *ICML*. June 2024. URL: <https://arxiv.org/abs/2406.16241>.
- [JSJ94] T. Jaakkola, S. Singh, and M. Jordan. “Reinforcement Learning Algorithm for Partially Observable Markov Decision Problems”. In: *NIPS*. 1994.
- [KAG19] A. Kirsch, J. van Amersfoort, and Y. Gal. “BatchBALD: Efficient and Diverse Batch Acquisition for Deep Bayesian Active Learning”. In: *NIPS*. 2019. URL: <http://arxiv.org/abs/1906.08158>.
- [Kai+19] L. Kaiser et al. “Model-based reinforcement learning for Atari”. In: *arXiv [cs.LG]* (Mar. 2019). URL: <http://arxiv.org/abs/1903.00374>.
- [Kak01] S. M. Kakade. “A Natural Policy Gradient”. In: *NIPS*. Vol. 14. 2001. URL: https://proceedings.neurips.cc/paper_files/paper/2001/file/4b86abe48d358ecf194c56c69108433e-Paper.pdf.
- [Kal+18] D. Kalashnikov et al. “QT-Opt: Scalable Deep Reinforcement Learning for Vision-Based Robotic Manipulation”. In: *CORL*. 2018. URL: <http://arxiv.org/abs/1806.10293>.
- [Kam+25] S. Kambhampati, K. Stechly, K. Valmeekam, L. Saldyt, S. Bhambri, V. Palod, A. Gundawar, S. R. Samineni, D. Kalwar, and U. Biswas. “Stop anthropomorphizing intermediate tokens as reasoning/thinking traces!” In: *arXiv [cs.AI]* (Apr. 2025). URL: <http://arxiv.org/abs/2504.09762>.
- [Kap+18] S. Kapturowski, G. Ostrovski, J. Quan, R. Munos, and W. Dabney. “Recurrent Experience Replay in Distributed Reinforcement Learning”. In: *ICLR*. Sept. 2018. URL: <https://openreview.net/pdf?id=r1lyTjAqYX>.
- [Kap+22] S. Kapturowski, V. Campos, R. Jiang, N. Rakicevic, H. van Hasselt, C. Blundell, and A. P. Badia. “Human-level Atari 200x faster”. In: *ICLR*. Sept. 2022. URL: <https://openreview.net/pdf?id=JtC6yOHRoJJ>.
- [Kaw+25] K. Kawaharazuka, J. Oh, J. Yamada, I. Posner, and Y. Zhu. “Vision-language-action models for robotics: A review towards real-world applications”. en. In: *IEEE Access* 13 (2025), pp. 162467–162504. URL: <https://vla-survey.github.io/>.
- [Kaz+25] A. Kazemnejad, M. Aghajohari, E. Portelance, A. Sordoni, S. Reddy, A. Courville, and N. L. Roux. “VinePPO: Refining credit assignment in RL training of LLMs”. In: *ICML*. 2025. URL: <https://arxiv.org/abs/2410.01679>.
- [KB09] G. Konidaris and A. Barto. “Skill Discovery in Continuous Reinforcement Learning Domains using Skill Chaining”. In: *Advances in Neural Information Processing Systems* 22 (2009). URL: https://proceedings.neurips.cc/paper_files/paper/2009/file/e0cf1f47118daebc5b16269099ad7347-Paper.pdf.
- [KD18] S. Kamthe and M. P. Deisenroth. “Data-Efficient Reinforcement Learning with Probabilistic Model Predictive Control”. In: *AISTATS*. 2018. URL: <http://proceedings.mlr.press/v84/kamthe18a/kamthe18a.pdf>.
- [KD25] A. Karan and Y. Du. “Reasoning with sampling: Your base model is smarter than you think”. In: *arXiv [cs.LG]* (Oct. 2025). URL: <http://arxiv.org/abs/2510.14901>.
- [Ke+19] L. Ke, S. Choudhury, M. Barnes, W. Sun, G. Lee, and S. Srinivasa. *Imitation Learning as f-Divergence Minimization*. arXiv:1905.12888. 2019.

- [KF88] D. M. Kilgour and N. M. Fraser. “A taxonomy of all ordinal 2 x 2 games”. en. In: *Theory Decis.* 24.2 (Mar. 1988), pp. 99–117. URL: <https://link.springer.com/article/10.1007/BF00132457>.
- [KGO12] H. J. Kappen, V. Gómez, and M. Opper. “Optimal control as a graphical model inference problem”. In: *Mach. Learn.* 87.2 (2012), pp. 159–182. URL: <https://doi.org/10.1007/s10994-012-5278-7>.
- [Kha+25a] Z. Khan, A. Prasad, E. Stengel-Eskin, J. Cho, and M. Bansal. “One life to learn: Inferring symbolic world models for stochastic environments from unguided exploration”. In: *arXiv [cs.AI]* (Oct. 2025). URL: <http://arxiv.org/abs/2510.12088>.
- [Kha+25b] D. Khatri, L. Madaan, R. Tiwari, R. Bansal, S. S. Duvvuri, M. Zaheer, I. S. Dhillon, D. Brandfonbrener, and R. Agarwal. “The art of scaling reinforcement learning compute for LLMs”. In: *arXiv [cs.LG]* (Oct. 2025). URL: <http://arxiv.org/abs/2510.13786>.
- [Khe+20] K. Khetarpal, Z. Ahmed, G. Comanici, D. Abel, and D. Precup. “What can I do here? A Theory of Affordances in Reinforcement Learning”. In: *Proceedings of the 37th International Conference on Machine Learning*. Ed. by H. D. Iii and A. Singh. Vol. 119. Proceedings of Machine Learning Research. PMLR, 2020, pp. 5243–5253. URL: <https://proceedings.mlr.press/v119/khetarpal20a.html>.
- [Khe+25] K. Khetarpal, Z. D. Guo, B. A. Pires, Y. Tang, C. Lyle, M. Rowland, N. Heess, D. Borsa, A. Guez, and W. Dabney. “A unifying framework for action-conditional self-predictive Reinforcement Learning”. In: *AISTATS*. 2025. URL: <https://arxiv.org/abs/2406.02035>.
- [Kid+20] R. Kidambi, A. Rajeswaran, P. Netrapalli, and T. Joachims. “MOREL: Model-Based Offline Reinforcement Learning”. In: *NIPS*. Vol. 33. 2020, pp. 21810–21823. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/f7efa4f864ae9b88d43527f4b14f750f-Paper.pdf.
- [Kir+21] R. Kirk, A. Zhang, E. Grefenstette, and T. Rocktäschel. “A survey of zero-shot generalisation in deep Reinforcement Learning”. In: *JAIR* (Nov. 2021). URL: <http://jair.org/index.php/jair/article/view/14174>.
- [KL02] S. Kakade and J. Langford. “Approximately Optimal Approximate Reinforcement Learning”. In: *ICML*. ICML ’02. Morgan Kaufmann Publishers Inc., 2002, pp. 267–274. URL: <http://dl.acm.org/citation.cfm?id=645531.656005>.
- [KL93] E. Kalai and E. Lehrer. “Rational learning leads to Nash equilibrium”. en. In: *Econometrica* 61.5 (Sept. 1993), p. 1019. URL: <https://www.jstor.org/stable/2951492>.
- [KLC98] L. P. Kaelbling, M. Littman, and A. Cassandra. “Planning and acting in Partially Observable Stochastic Domains”. In: *AIJ* 101 (1998).
- [Kli+24] M. Klissarov, P. D’Oro, S. Sodhani, R. Raileanu, P.-L. Bacon, P. Vincent, A. Zhang, and M. Henaff. “Motif: Intrinsic motivation from artificial intelligence feedback”. In: *ICLR*. 2024.
- [Kli+25a] M. Klissarov, A. Bagaria, Z. Luo, G. Konidaris, D. Precup, and M. C. Machado. “Discovering temporal structure: An overview of hierarchical reinforcement learning”. In: *arXiv [cs.AI]* (June 2025). URL: <http://arxiv.org/abs/2506.14045>.
- [Kli+25b] M. Klissarov, R. Devon Hjelm, A. T. Toshev, and B. Mazoure. “On the Modeling Capabilities of Large Language Models for Sequential Decision Making”. In: *ICLR*. 2025. URL: <https://openreview.net/forum?id=vodsIF3o7N>.
- [KLP11] L. P. Kaelbling and T. Lozano-Pérez. “Hierarchical task and motion planning in the now”. In: *ICRA*. 2011, pp. 1470–1477. URL: <http://dx.doi.org/10.1109/ICRA.2011.5980391>.
- [KM23] M. Klissarov and M. C. Machado. “Deep Laplacian-based options for temporally-extended exploration”. In: *ICML*. Jan. 2023. URL: <https://arxiv.org/abs/2301.11181>.

- [KMN99] M. Kearns, Y. Mansour, and A. Ng. “A Sparse Sampling Algorithm for Near-Optimal Planning in Large Markov Decision Processes”. In: *IJCAI*. 1999.
- [Kon+24] D. Kong, D. Xu, M. Zhao, B. Pang, J. Xie, A. Lizarraga, Y. Huang, S. Xie, and Y. N. Wu. “Latent Plan Transformer for trajectory abstraction: Planning as latent space inference”. In: *NIPS*. Feb. 2024.
- [Kor+22] T. Korbak, H. Elsahar, G. Kruszewski, and M. Dymetman. “On Reinforcement Learning and Distribution Matching for fine-tuning language models with no catastrophic forgetting”. In: *NIPS*. June 2022. URL: <https://arxiv.org/abs/2206.00761>.
- [Kov+22] V. Kovařík, M. Schmid, N. Burch, M. Bowling, and V. Lisý. “Rethinking formal models of partially observable multiagent decision making”. In: *Artificial Intelligence* (2022). URL: <http://arxiv.org/abs/1906.11110>.
- [Koz+21] T. Kozuno, Y. Tang, M. Rowland, R. Munos, S. Kapturowski, W. Dabney, M. Valko, and D. Abel. “Revisiting Peng’s Q-lambda for modern reinforcement learning”. In: *ICML* 139 (Feb. 2021). Ed. by M. Meila and T. Zhang, pp. 5794–5804. URL: <https://proceedings.mlr.press/v139/kozuno21a/kozuno21a.pdf>.
- [KP19] K. Khetarpal and D. Precup. “Learning options with interest functions”. en. In: *AAAI* 33.01 (July 2019), pp. 9955–9956. URL: <https://ojs.aaai.org/index.php/AAAI/article/view/5114>.
- [KPB22] T. Korbak, E. Perez, and C. L. Buckley. “RL with KL penalties is better viewed as Bayesian inference”. In: *EMNLP*. May 2022. URL: <http://arxiv.org/abs/2205.11275>.
- [KPL19] A. Kumar, X. B. Peng, and S. Levine. “Reward-Conditioned Policies”. In: *arXiv [cs.LG]* (Dec. 2019). URL: <http://arxiv.org/abs/1912.13465>.
- [KPN05] A. S. Klyubin, D. Polani, and C. L. Nehaniv. “Empowerment: A universal agent-centric measure of control”. en. In: *2005 IEEE Congress on Evolutionary Computation*. Vol. 1. IEEE, 2005, 128–135 Vol.1. URL: <https://ieeexplore.ieee.org/abstract/document/1554676>.
- [KS02] M. Kearns and S. Singh. “Near-Optimal Reinforcement Learning in Polynomial Time”. en. In: *MLJ* 49.2/3 (Nov. 2002), pp. 209–232. URL: <https://link.springer.com/article/10.1023/A:1017984413808>.
- [KSS23] T. Kneib, A. Silbersdorff, and B. Säcken. “Rage Against the Mean – A Review of Distributional Regression Approaches”. In: *Econometrics and Statistics* 26 (Apr. 2023), pp. 99–123. URL: <https://www.sciencedirect.com/science/article/pii/S2452306221000824>.
- [Kub+25] J. G. Kuba, M. Gu, Q. Ma, Y. Tian, and V. Mohan. “Language Self-play for data-free training”. In: *arXiv [cs.AI]* (Sept. 2025). URL: <http://arxiv.org/abs/2509.07414>.
- [Kuh51] H. W. Kuhn. “A SIMPLIFIED TWO-PERSON POKER”. en. In: *Contributions to the Theory of Games (AM-24), Volume I*. Ed. by H. W. Kuhn and A. W. Tucker. Princeton: Princeton University Press, Dec. 1951, pp. 97–104. URL: <https://www.degruyter.com/document/doi/10.1515/9781400881727-010/html>.
- [Kul+16] T. D. Kulkarni, K. Narasimhan, A. Saeedi, and J. Tenenbaum. “Hierarchical Deep Reinforcement Learning: Integrating Temporal Abstraction and Intrinsic Motivation”. In: *NIPS* 29 (2016). URL: <https://proceedings.neurips.cc/paper/2016/hash/f442d33fa06832082290ad8544a8da27-Abstract.html>.
- [Kum+19] A. Kumar, J. Fu, M. Soh, G. Tucker, and S. Levine. “Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction”. In: *NIPS*. Vol. 32. 2019. URL: https://proceedings.neurips.cc/paper_files/paper/2019/file/c2073ffa77b5357a498057413bb09d3a-Paper.pdf.
- [Kum+20] A. Kumar, A. Zhou, G. Tucker, and S. Levine. “Conservative Q-Learning for Offline Reinforcement Learning”. In: *NIPS*. June 2020.
- [Kum+23] A. Kumar, R. Agarwal, X. Geng, G. Tucker, and S. Levine. “Offline Q-Learning on Diverse Multi-Task Data Both Scales And Generalizes”. In: *ICLR*. 2023. URL: <http://arxiv.org/abs/2211.15144>.

- [Kum+24] S. Kumar, H. J. Jeon, A. Lewandowski, and B. Van Roy. “The Need for a Big World Simulator: A Scientific Challenge for Continual Learning”. In: *Finding the Frame: An RLC Workshop for Examining Conceptual Frameworks*. July 2024. URL: <https://openreview.net/pdf?id=10XMwt1nMJ>.
- [Kur+19] T. Kurutach, I. Clavera, Y. Duan, A. Tamar, and P. Abbeel. “Model-Ensemble Trust-Region Policy Optimization”. In: *ICLR*. 2019. URL: <http://arxiv.org/abs/1802.10592>.
- [KWW22] M. J. Kochenderfer, T. A. Wheeler, and K. Wray. *Algorithms for Decision Making*. The MIT Press, 2022. URL: <https://github.com/sisl/algorithmsbook/>.
- [LA21] H. Liu and P. Abbeel. “APS: Active Pretraining with Successor Features”. en. In: *ICML*. PMLR, July 2021, pp. 6736–6747. URL: <https://proceedings.mlr.press/v139/liu21b.html>.
- [Lab+25] P. Laban, H. Hayashi, Y. Zhou, and J. Neville. “LLMs Get Lost In Multi-Turn Conversation”. In: *arXiv [cs.CL]* (May 2025). URL: <http://arxiv.org/abs/2505.06120>.
- [Lad+22] P. Ladosz, L. Weng, M. Kim, and H. Oh. “Exploration in deep reinforcement learning: A survey”. en. In: *Inf. Fusion* 85 (Sept. 2022), pp. 1–22. URL: <http://dx.doi.org/10.1016/j.inffus.2022.03.003>.
- [Lai+21] H. Lai, J. Shen, W. Zhang, Y. Huang, X. Zhang, R. Tang, Y. Yu, and Z. Li. “On effective scheduling of model-based reinforcement learning”. In: *NIPS* 34 (Nov. 2021). Ed. by M. Ranzato, A. Beygelzimer, Y. Dauphin, P. S. Liang, and J. W. Vaughan, pp. 3694–3705. URL: https://proceedings.neurips.cc/paper_files/paper/2021/hash/1e4d36177d71bbb3558e43af9577d70e-Abstract.html.
- [Lai+25] C. Laidlaw, E. Bronstein, T. Guo, D. Feng, L. Berglund, J. Svegliato, S. Russell, and A. Dragan. “AssistanceZero: Scalably Solving Assistance Games”. In: *ICML*. Apr. 2025. URL: <https://arxiv.org/abs/2504.07091>.
- [Lam+20] N. Lambert, B. Amos, O. Yadan, and R. Calandra. “Objective Mismatch in Model-based Reinforcement Learning”. In: *Conf. on Learning for Dynamics and Control (L4DC)*. Feb. 2020.
- [Lam+24] N. Lambert et al. “TÜLU 3: Pushing frontiers in open language model post-training”. In: *arXiv [cs.CL]* (Nov. 2024). URL: <http://arxiv.org/abs/2411.15124>.
- [Lam25] N. Lambert. *Reinforcement Learning from Human Feedback*. 2025. URL: <https://rlhfbook.com/book.pdf>.
- [Lan+09] M. Lanctot, K. Waugh, M. Zinkevich, and M. Bowling. “Monte Carlo Sampling for Regret Minimization in Extensive Games”. In: *NIPS* 22 (2009). URL: https://proceedings.neurips.cc/paper_files/paper/2009/file/00411460f7c92d2124a67ea0f4cb5f85-Paper.pdf.
- [Lan+17] M. Lanctot, V. Zambaldi, A. Gruslys, A. Lazaridou, K. Tuyls, J. Perolat, D. Silver, and T. Graepel. “A unified game-theoretic approach to multiagent reinforcement learning”. In: *NIPS*. Nov. 2017. URL: <https://arxiv.org/abs/1711.00832>.
- [Law+18] D. Lawson, G. Tucker, C. A. Naesseth, C. J. Maddison, R. P. Adams, and Y. W. Teh. “Twisted Variational Sequential Monte Carlo”. In: *BDL Workshop*. 2018. URL: <https://bayesiandeeplearning.org/2018/papers/111.pdf>.
- [Law+22] D. Lawson, A. Raventos, A. Warrington, and S. Linderman. “SIXO: Smoothing Inference with Twisted Objectives”. In: *Advances in Neural Information Processing Systems*. Oct. 2022. URL: <https://openreview.net/forum?id=bDyLgfvZ0qJ>.
- [Law+25] N. P. Lawrence, P. D. Loewen, M. G. Forbes, R. B. Gopaluni, and A. Mesbah. “A view on learning robust goal-conditioned value functions: Interplay between RL and MPC”. In: *arXiv [eess.SY]* (Feb. 2025). URL: <http://arxiv.org/abs/2502.06996>.
- [LBS08] K. Leyton-Brown and Y. Shoham. *Essentials of game theory: A concise, multidisciplinary introduction*. Synthesis lectures on artificial intelligence and machine learning. Cham: Springer International Publishing, 2008. URL: <https://www.gtessentials.org/>.

- [LCZ25] J. Liu, B. Chen, and C. Zhang. “Speculative Prefill: Turbocharging TTFT with Lightweight and Training-Free Token Importance Estimation”. In: *ICML*. June 2025. URL: <https://openreview.net/forum?id=bzbuZ0ItBq>.
- [Le+20] T. A. Le, A. R. Kosiorek, N. Siddharth, Y. W. Teh, and F. Wood. “Revisiting Reweighted Wake-Sleep for Models with Stochastic Control Flow”. en. In: *UAI*. PMLR, Aug. 2020, pp. 1039–1049. URL: <https://proceedings.mlr.press/v115/le20a.html>.
- [LeC22] Y. LeCun. *A Path Towards Autonomous Machine Intelligence*. 2022. URL: <https://openreview.net/pdf?id=BZ5a1r-kVsf>.
- [Lee+20] K.-H. Lee, I. Fischer, A. Liu, Y. Guo, H. Lee, J. Canny, and S. Guadarrama. “Predictive Information Accelerates Learning in RL”. In: *NIPS* 33 (2020), pp. 11890–11901. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/89b9e0a6f6d1505fe13dea0f18a2dcfa-Paper.pdf.
- [Lee+22] K.-H. Lee et al. “Multi-Game Decision Transformers”. In: *NIPS* abs/2205.15241 (May 2022), pp. 27921–27936. URL: <http://dx.doi.org/10.48550/arXiv.2205.15241>.
- [Leh24] M. Lehmann. “The definitive guide to policy gradients in deep reinforcement learning: Theory, algorithms and implementations”. In: *arXiv [cs.LG]* (Jan. 2024). URL: <http://arxiv.org/abs/2401.13662>.
- [Lev18] S. Levine. “Reinforcement Learning and Control as Probabilistic Inference: Tutorial and Review”. In: (2018). arXiv: [1805.00909 \[cs.LG\]](http://arxiv.org/abs/1805.00909). URL: <http://arxiv.org/abs/1805.00909>.
- [Lev+18] A. Levy, G. Konidaris, R. Platt, and K. Saenko. “Learning Multi-Level Hierarchies with Hindsight”. In: *ICLR*. Sept. 2018. URL: <https://openreview.net/pdf?id=ryzEC0AcY7>.
- [Lev+20a] S. Levine, A. Kumar, G. Tucker, and J. Fu. “Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems”. In: (2020). arXiv: [2005.01643 \[cs.LG\]](http://arxiv.org/abs/2005.01643). URL: <http://arxiv.org/abs/2005.01643>.
- [Lev+20b] S. Levine, A. Kumar, G. Tucker, and J. Fu. *Offline Reinforcement Learning: Tutorial, Review, and Perspectives on Open Problems*. arXiv:2005.01643. 2020.
- [Lew+23] A. K. Lew, T. Zhi-Xuan, G. Grand, and V. K. Mansinghka. “Sequential Monte Carlo Steering of Large Language Models using Probabilistic Programs”. In: *arXiv [cs.AI]* (June 2023). URL: <http://arxiv.org/abs/2306.03081>.
- [LGR12] S. Lange, T. Gabel, and M. Riedmiller. “Batch reinforcement learning”. en. In: *Adaptation, Learning, and Optimization*. Adaptation, learning, and optimization. Berlin, Heidelberg: Springer Berlin Heidelberg, 2012, pp. 45–73. URL: https://link.springer.com/chapter/10.1007/978-3-642-27645-3_2.
- [LHC25] C. Lu, S. Hu, and J. Clune. “Intelligent Go-Explore: Standing on the shoulders of giant foundation models”. In: *ICLR* (May 2025). URL: <http://arxiv.org/abs/2405.15143>.
- [LHP22] M. Lauri, D. Hsu, and J. Pajarinen. “Partially Observable Markov Decision Processes in Robotics: A Survey”. In: *IEEE Trans. Rob.* (Sept. 2022). URL: <http://arxiv.org/abs/2209.10342>.
- [LHS13] T. Lattimore, M. Hutter, and P. Sunehag. “The Sample-Complexity of General Reinforcement Learning”. en. In: *ICML*. PMLR, May 2013, pp. 28–36. URL: <https://proceedings.mlr.press/v28/lattimore13.html>.
- [Li+10] L. Li, W. Chu, J. Langford, and R. E. Schapire. “A contextual-bandit approach to personalized news article recommendation”. In: *WWW*. 2010.
- [Li18] Y. Li. “Deep Reinforcement Learning”. In: (2018). arXiv: [1810.06339 \[cs.LG\]](http://arxiv.org/abs/1810.06339). URL: <http://arxiv.org/abs/1810.06339>.
- [Li23] S. E. Li. *Reinforcement learning for sequential decision and optimal control*. en. Singapore: Springer Nature Singapore, 2023. URL: <https://link.springer.com/book/10.1007/978-981-19-7784-8>.

- [Li+23a] T. Li, A. Beirami, M. Sanjabi, and V. Smith. “On Tilted Losses in Machine Learning: Theory and Applications”. In: *TMLR* (2023). URL: <http://arxiv.org/abs/2109.06141>.
- [Li+23b] Z. Li, M. Lanctot, K. R. McKee, L. Marris, I. Gemp, D. Hennes, P. Muller, K. Larson, Y. Bachrach, and M. P. Wellman. “Combining tree-search, generative models, and Nash bargaining concepts in game-theoretic reinforcement learning”. In: *arXiv [cs.AI]* (Feb. 2023). URL: <http://arxiv.org/abs/2302.00797>.
- [Li+24a] H. Li, X. Yang, Z. Wang, X. Zhu, J. Zhou, Y. Qiao, X. Wang, H. Li, L. Lu, and J. Dai. “Auto MC-Reward: Automated Dense Reward Design with Large Language Models for Minecraft”. In: *CVPR*. 2024, pp. 16426–16435. URL: https://openaccess.thecvf.com/content/CVPR2024/papers/Li_Auto_MC-Reward_Automated_Dense_Reward_Design_with_Large_Language_Models_CVPR_2024_paper.pdf.
- [Li+24b] J. Li, Z. Zang, Q. Ji, C. Sun, W. Qiang, J. Zhang, C. Zheng, F. Sun, and H. Xiong. “Rethinking generalizability and discriminability of self-supervised learning from evolutionary game theory perspective”. In: *Int. J. Comput. Vis.* (2024). URL: <http://arxiv.org/abs/2412.00542>.
- [Li+24c] Z. Li, H. Liu, D. Zhou, and T. Ma. “Chain of thought empowers transformers to solve inherently serial problems”. In: *ICLR*. Feb. 2024. URL: <https://arxiv.org/abs/2402.12875>.
- [Li+25] Y. Li, X. Shen, X. Yao, X. Ding, Y. Miao, R. Krishnan, and R. Padman. “Beyond single-turn: A survey on multi-turn interactions with large language models”. In: *arXiv [cs.CL]* (Apr. 2025). URL: <http://arxiv.org/abs/2504.04717>.
- [Lia+25] J. Liao, M. Wen, J. Wang, and W. Zhang. “MARFT: Multi-Agent Reinforcement Fine-Tuning”. In: *arXiv [cs.MA]* (Apr. 2025). URL: <http://arxiv.org/abs/2504.16129>.
- [Lil+16] T. P. Lillicrap, J. J. Hunt, A. Pritzel, N. Heess, T. Erez, Y. Tassa, D. Silver, and D. Wierstra. “Continuous control with deep reinforcement learning”. In: *ICLR*. 2016. URL: <http://arxiv.org/abs/1509.02971>.
- [Lim+23] M. H. Lim, T. J. Becker, M. J. Kochenderfer, C. J. Tomlin, and Z. N. Sunberg. “Optimality guarantees for particle belief approximation of POMDPs”. In: *J. Artif. Intell. Res.* (2023). URL: <https://jair.org/index.php/jair/article/view/14525>.
- [Lin+19] C. Linke, N. M. Ady, M. White, T. Degris, and A. White. “Adapting behaviour via intrinsic reward: A survey and empirical study”. In: *J. Artif. Intell. Res.* (June 2019). URL: <http://arxiv.org/abs/1906.07865>.
- [Lin+23] F. Lin, S. Huang, T. Pearce, W. Chen, and W.-W. Tu. “TiZero: Mastering multi-agent football with curriculum learning and self-play”. In: *22nd International Conference on Autonomous Agents and Multiagent Systems*. Feb. 2023. URL: <https://arxiv.org/abs/2302.07515>.
- [Lin92] L.-J. Lin. “Self-Improving Reactive Agents Based on Reinforcement Learning, Planning and Teaching”. In: *Mach. Learn.* 8.3-4 (1992), pp. 293–321. URL: <https://doi.org/10.1007/BF00992699>.
- [Lio+22] V. Lioutas, J. W. Lavington, J. Sefas, M. Niedoba, Y. Liu, B. Zwartsenberg, S. Dabiri, F. Wood, and A. Scibior. “Critic Sequential Monte Carlo”. In: *ICLR*. Sept. 2022. URL: <https://openreview.net/pdf?id=0btGcyKmwNa>.
- [Lit94] M. Littman. “Markov games as a framework for multi-agent reinforcement learning”. en. In: *Machine Learning Proceedings 1994*. Elsevier, Jan. 1994, pp. 157–163. URL: <https://www.sciencedirect.com/science/article/abs/pii/B9781558603356500271>.
- [Liu+24] Z. Liu, C. Du, W. S. Lee, and M. Lin. “Locality Sensitive Sparse Encoding for Learning World Models Online”. In: *The Twelfth International Conference on Learning Representations*. 2024. URL: <https://openreview.net/pdf?id=i8PjQT3Uig>.
- [Liu+25a] M. Liu, L. Jiang, Y. Liang, S. S. Du, Y. Choi, T. Althoff, and N. Jaques. “Chasing moving targets with online self-play reinforcement learning for safer language models”. In: *arXiv [cs.LG]* (June 2025). URL: <http://arxiv.org/abs/2506.07468>.

- [Liu+25b] Z. Liu, C. Chen, W. Li, P. Qi, T. Pang, C. Du, W. S. Lee, and M. Lin. “Understanding R1-Zero-Like Training: A Critical Perspective”. In: *COLM*. 2025. URL: <https://github.com/sail-sg/understand-r1-zero/blob/main/understand-r1-zero.pdf>.
- [Liu+25c] Z. Liu, C. Chen, W. Li, P. Qi, T. Pang, C. Du, W. S. Lee, and M. Lin. “Understanding R1-zero-like training: A critical perspective”. In: *arXiv [cs.LG]* (Mar. 2025). URL: <http://arxiv.org/abs/2503.20783>.
- [Liu+25d] Z. Liu et al. “GEM: A Gym for Agentic LLMs”. In: *arXiv [cs.LG]* (Oct. 2025). URL: <http://arxiv.org/abs/2510.01051>.
- [Liu+25e] Z. Liu et al. “Part I: Tricks or traps? A deep dive into RL for LLM reasoning”. In: *arXiv [cs.LG]* (Aug. 2025). URL: <http://arxiv.org/abs/2508.08221>.
- [LL17] S. Lundberg and S.-I. Lee. “A unified approach to interpreting model predictions”. In: *NIPS*. May 2017. URL: <https://arxiv.org/abs/1705.07874>.
- [LMLFP11] G. J. Laurent, L. Matignon, and N Le Fort-Piat. “The world of independent learners is not markovian”. en. In: *International Journal of Knowledge-Based and Intelligent Engineering Systems* (2011). URL: <https://hal.science/hal-00601941/document>.
- [LMW24] B. Li, N. Ma, and Z. Wang. “Rewarded Region Replay (R3) for policy learning with discrete action space”. In: *arXiv [cs.LG]* (May 2024). URL: <http://arxiv.org/abs/2405.16383>.
- [Lor24] J. Lorraine. “Scalable nested optimization for deep learning”. In: *arXiv [cs.LG]* (July 2024). URL: <http://arxiv.org/abs/2407.01526>.
- [Lou+25] J. Loula et al. “Syntactic and semantic control of large language models via sequential Monte Carlo”. In: *ICLR*. Apr. 2025. URL: <https://arxiv.org/abs/2504.13139>.
- [LÖW21] T. van de Laar, A. Özçelikkale, and H. Wymeersch. “Application of the Free Energy Principle to Estimation and Control”. In: *IEEE Trans. Signal Process.* 69 (2021), pp. 4234–4244. URL: <http://dx.doi.org/10.1109/TSP.2021.3095711>.
- [LPC22] N. Lambert, K. Pister, and R. Calandra. “Investigating Compounding Prediction Errors in Learned Dynamics Models”. In: *arXiv [cs.LG]* (Mar. 2022). URL: <http://arxiv.org/abs/2203.09637>.
- [LR10] S. Lange and M. Riedmiller. “Deep auto-encoder neural networks in reinforcement learning”. en. In: *IJCNN*. IEEE, July 2010, pp. 1–8. URL: <https://ieeexplore.ieee.org/abstract/document/5596468>.
- [LS01] M. Littman and R. S. Sutton. “Predictive Representations of State”. In: *NIPS 14* (2001). URL: https://proceedings.neurips.cc/paper_files/paper/2001/file/1e4d36177d71bbb3558e43af9577d70e-Paper.pdf.
- [LS19] T. Lattimore and C. Szepesvari. *Bandit Algorithms*. Cambridge, 2019.
- [Lu+23] X. Lu, B. Van Roy, V. Dwaracherla, M. Ibrahimi, I. Osband, and Z. Wen. “Reinforcement Learning, Bit by Bit”. In: *Found. Trends® Mach. Learn.* (2023). URL: <https://www.nowpublishers.com/article/Details/MAL-097>.
- [Luo+22] F.-M. Luo, T. Xu, H. Lai, X.-H. Chen, W. Zhang, and Y. Yu. “A survey on model-based reinforcement learning”. In: *arXiv [cs.LG]* (June 2022). URL: <http://arxiv.org/abs/2206.09328>.
- [Luo+25] J. Luo et al. “Large Language Model agent: A survey on methodology, applications and challenges”. In: *arXiv [cs.CL]* (Mar. 2025). URL: <http://arxiv.org/abs/2503.21460>.
- [LV06] F. Liese and I. Vajda. “On divergences and informations in statistics and information theory”. In: *IEEE Transactions on Information Theory* 52.10 (2006), pp. 4394–4412.
- [LWL06] L. Li, T. J. Walsh, and M. L. Littman. “Towards a Unified Theory of State Abstraction for MDPs”. In: (2006). URL: <https://thomasjwalsh.net/pub/aima06Towards.pdf>.

- [LZZ22] M. Liu, M. Zhu, and W. Zhang. “Goal-conditioned reinforcement learning: Problems and solutions”. In: *IJCAI*. Jan. 2022. URL: <https://arxiv.org/abs/2201.08299>.
- [Ma+24] Y. J. Ma, W. Liang, G. Wang, D.-A. Huang, O. Bastani, D. Jayaraman, Y. Zhu, L. Fan, and A. Anandkumar. “Eureka: Human-Level Reward Design via Coding Large Language Models”. In: *ICLR*. 2024.
- [MA93] A. W. Moore and C. G. Atkeson. “Prioritized Sweeping: Reinforcement Learning with Less Data and Less Time”. In: *Machine Learning* 13.1 (1993), 103—130.
- [Mac+18a] M. C. Machado, M. G. Bellemare, E. Talvitie, J. Veness, M. Hausknecht, and M. Bowling. “Revisiting the Arcade Learning Environment: Evaluation Protocols and Open Problems for General Agents”. In: *JAIR* (2018). URL: <http://arxiv.org/abs/1709.06009>.
- [Mac+18b] M. C. Machado, C. Rosenbaum, X. Guo, M. Liu, G. Tesauro, and M. Campbell. “Eigenoption Discovery through the Deep Successor Representation”. In: *ICLR*. Feb. 2018. URL: <https://openreview.net/pdf?id=Bk8ZcAxR->.
- [Mac+23] M. C. Machado, A. Barreto, D. Precup, and M. Bowling. “Temporal Abstraction in Reinforcement Learning with the Successor Representation”. In: *JMLR* 24.80 (2023), pp. 1–69. URL: <http://jmlr.org/papers/v24/21-1213.html>.
- [Mac+24] M. Macfarlane, E. Toledo, D. J. Byrne, P. Duckworth, and A. Laterre. “SPO: Sequential Monte Carlo Policy Optimisation”. In: *NIPS*. Nov. 2024. URL: <https://openreview.net/pdf?id= XKvYcPPH5G>.
- [Mae+09] H. Maei, C. Szepesvári, S. Bhatnagar, D. Precup, D. Silver, and R. S. Sutton. “Convergent Temporal-Difference Learning with Arbitrary Smooth Function Approximation”. In: *NIPS*. Vol. 22. 2009. URL: https://proceedings.neurips.cc/paper_files/paper/2009/file/3a15c7d0bbe60300a39f76f8a5ba6896-Paper.pdf.
- [MAF22] V. Micheli, E. Alonso, and F. Fleuret. “Transformers are Sample-Efficient World Models”. In: *ICLR*. Sept. 2022.
- [MAF24] V. Micheli, E. Alonso, and F. Fleuret. “Efficient world models with context-aware tokenization”. In: *ICML*. June 2024.
- [Mai+25] X. Mai, H. Xu, Z.-Z. Li, W. Xing, W. Wang, J. Hu, Y. Zhang, and W. Zhang. “Agent RL scaling law: Agent RL with spontaneous code execution for mathematical problem solving”. In: *arXiv [cs.AI]* (Aug. 2025). URL: <http://arxiv.org/abs/2505.07773>.
- [Maj21] S. J. Majeed. “Abstractions of general reinforcement learning: An inquiry into the scalability of generally intelligent agents”. PhD thesis. ANU, Dec. 2021. URL: <https://arxiv.org/abs/2112.13404>.
- [Man+19] D. J. Mankowitz, N. Levine, R. Jeong, Y. Shi, J. Kay, A. Abdolmaleki, J. T. Springenberg, T. Mann, T. Hester, and M. Riedmiller. “Robust Reinforcement Learning for Continuous Control with Model Misspecification”. In: (2019). arXiv: [1906.07516](https://arxiv.org/abs/1906.07516) [cs.LG]. URL: <http://arxiv.org/abs/1906.07516>.
- [Mar10] J. Martens. “Deep learning via Hessian-free optimization”. In: *ICML*. 2010. URL: http://www.cs.toronto.edu/~asamir/cifar/HFO_James.pdf.
- [Mar16] J. Martens. “Second-order optimization for neural networks”. PhD thesis. Toronto, 2016. URL: http://www.cs.toronto.edu/~jmartens/docs/thesis_phd_martens.pdf.
- [Mar20] J. Martens. “New insights and perspectives on the natural gradient method”. In: *JMLR* (2020). URL: <http://arxiv.org/abs/1412.1193>.
- [Mar21] J. Marino. “Predictive Coding, Variational Autoencoders, and Biological Connections”. en. In: *Neural Comput.* 34.1 (2021), pp. 1–44. URL: http://dx.doi.org/10.1162/neco_a_01458.

- [Maz+22] P. Mazzaglia, T. Verbelen, O. Çatal, and B. Dhoedt. “The Free Energy Principle for Perception and Action: A Deep Learning Perspective”. en. In: *Entropy* 24.2 (2022). URL: <http://dx.doi.org/10.3390/e24020301>.
- [MB01] A. McGovern and A. G. Barto. “Automatic discovery of subgoals in reinforcement learning using di- verse density”. In: (2001).
- [MBB20] M. C. Machado, M. G. Bellemare, and M. Bowling. “Count-based exploration with the successor representation”. en. In: *AAAI* 34.04 (Apr. 2020), pp. 5125–5133. URL: <https://ojs.aaai.org/index.php/AAAI/article/view/5955>.
- [McM+13] H. B. McMahan, G. Holt, D. Sculley, M. Young, D. Ebner, J. Grady, L. Nie, T. Phillips, E. Davydov, D. Golovin, et al. “Ad click prediction: a view from the trenches”. In: *KDD*. 2013, pp. 1222–1230.
- [Men+21] R. Mendonca, O. Rybkin, K. Daniilidis, D. Hafner, and D. Pathak. “Discovering and Achieving Goals via World Models”. In: *NIPS* 34 (Dec. 2021), pp. 24379–24391. URL: <https://proceedings.neurips.cc/paper/2021/hash/cc4af25fa9d2d5c953496579b75f6f6c-Abstract.html>.
- [Men+23] W. Meng, Q. Zheng, G. Pan, and Y. Yin. “Off-Policy Proximal Policy Optimization”. en. In: *AAAI* 37.8 (June 2023), pp. 9162–9170. URL: <https://ojs.aaai.org/index.php/AAAI/article/view/26099>.
- [Met+17] L. Metz, J. Ibarz, N. Jaitly, and J. Davidson. “Discrete Sequential Prediction of Continuous Actions for Deep RL”. In: (2017). arXiv: [1705.05035 \[cs.LG\]](https://arxiv.org/abs/1705.05035). URL: <http://arxiv.org/abs/1705.05035>.
- [Met+22] Meta Fundamental AI Research Diplomacy Team (FAIR)[†] et al. “Human-level play in the game of Diplomacy by combining language models with strategic reasoning”. en. In: *Science* 378.6624 (Dec. 2022), pp. 1067–1074. URL: <http://dx.doi.org/10.1126/science.ade9097>.
- [Mey22] S. Meyn. *Control Systems and Reinforcement Learning*. Cambridge, 2022. URL: <https://meyn.ece.ufl.edu/2021/08/01/control-systems-and-reinforcement-learning/>.
- [MG14] D. B. Markant and T. M. Gureckis. “Is it better to select or to receive? Learning via active and passive hypothesis testing”. en. In: *J. Exp. Psychol. Gen.* 143.1 (Feb. 2014), pp. 94–122. URL: <https://psycnet.apa.org/doiLanding?doi=10.1037/a0032108>.
- [MG15] J. Martens and R. Grosse. “Optimizing Neural Networks with Kronecker-factored Approximate Curvature”. In: *ICML*. 2015. URL: <http://arxiv.org/abs/1503.05671>.
- [MGR18] H. Mania, A. Guy, and B. Recht. “Simple random search of static linear policies is competitive for reinforcement learning”. In: *NIPS*. Ed. by S Bengio, H Wallach, H Larochelle, K Grauman, N Cesa-Bianchi, and R Garnett. Curran Associates, Inc., 2018, pp. 1800–1809. URL: <http://papers.nips.cc/paper/7451-simple-random-search-of-static-linear-policies-is-competitive-for-reinforcement-learning.pdf>.
- [Mik+20] V. Mikulik, G. Delétang, T. McGrath, T. Genewein, M. Martic, S. Legg, and P. Ortega. “Meta-trained agents implement Bayes-optimal agents”. In: *NIPS* 33 (2020), pp. 18691–18703. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/d902c3ce47124c66ce615d5ad9ba304f-Paper.pdf.
- [Mil20] B. Millidge. “Deep Active Inference as Variational Policy Gradients”. In: *J. Mathematical Psychology* (2020). URL: <http://arxiv.org/abs/1907.03876>.
- [Mil+20] B. Millidge, A. Tschantz, A. K. Seth, and C. L. Buckley. “On the Relationship Between Active Inference and Control as Inference”. In: *International Workshop on Active Inference*. 2020. URL: <http://arxiv.org/abs/2006.12964>.
- [Min+24] M. J. Min, Y. Ding, L. Buratti, S. Pujar, G. Kaiser, S. Jana, and B. Ray. “Beyond Accuracy: Evaluating Self-Consistency of Code Large Language Models with IdentityChain”. In: *ICLR*. 2024. URL: <https://openreview.net/forum?id=caW7LdAALh>.

- [ML25] D. Mittal and W. S. Lee. “Differentiable Tree Search Network”. In: *ICLR*. 2025. URL: <https://arxiv.org/abs/2401.11660>.
- [MLR24] D. McNamara, J. Loper, and J. Regier. “Sequential Monte Carlo for Inclusive KL Minimization in Amortized Variational Inference”. en. In: *AISTATS*. PMLR, Apr. 2024, pp. 4312–4320. URL: <https://proceedings.mlr.press/v238/mcnamara24a.html>.
- [MM07] S. Mahadevan and M. Maggioni. “Proto-value functions: A Laplacian framework for learning representation and control in Markov decision processes”. In: *JMLR* 8.74 (Dec. 2007), pp. 2169–2231. URL: <https://jmlr.org/papers/v8/mahadevan07a.html>.
- [MM90] D. Q. Mayne and H. Michalska. “Receding horizon control of nonlinear systems”. In: *IEEE Trans. Automat. Contr.* 35.7 (1990), pp. 814–824.
- [MMT17] C. J. Maddison, A. Mnih, and Y. W. Teh. “The Concrete Distribution: A Continuous Relaxation of Discrete Random Variables”. In: *ICLR*. 2017. URL: <http://arxiv.org/abs/1611.00712>.
- [MMT24] S. Mannor, Y. Mansour, and A. Tamar. *Reinforcement Learning: Foundations*. 2024. URL: <https://sites.google.com/corp/view/rlfoundations/home>.
- [Mni+15] V. Mnih et al. “Human-level control through deep reinforcement learning”. In: *Nature* 518.7540 (2015), pp. 529–533.
- [Mni+16] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. P. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu. “Asynchronous Methods for Deep Reinforcement Learning”. In: *ICML*. 2016. URL: <http://arxiv.org/abs/1602.01783>.
- [Mo+25] Z. Mo, X. Li, Y. Chen, and L. Bing. “Multi-Agent Tool-Integrated Policy Optimization”. In: *arXiv [cs.CL]* (Oct. 2025). URL: <http://arxiv.org/abs/2510.04678>.
- [Moe+23] T. M. Moerland, J. Broekens, A. Plaat, and C. M. Jonker. “Model-based Reinforcement Learning: A Survey”. In: *Foundations and Trends in Machine Learning* 16.1 (2023), pp. 1–118. URL: <https://arxiv.org/abs/2006.16712>.
- [Moe+25] A. Moeini, J. Wang, J. Beck, E. Blaser, S. Whiteson, R. Chandra, and S. Zhang. “A survey of in-context reinforcement learning”. In: *arXiv [cs.LG]* (Feb. 2025). URL: <http://arxiv.org/abs/2502.07978>.
- [Moh+20] S. Mohamed, M. Rosca, M. Figurnov, and A. Mnih. “Monte Carlo Gradient Estimation in Machine Learning”. In: *JMLR* 21.132 (2020), pp. 1–62. URL: <http://jmlr.org/papers/v21/19-346.html>.
- [Mor+17] M. Moravčík, M. Schmid, N. Burch, V. Lisý, D. Morrill, N. Bard, T. Davis, K. Waugh, M. Johanson, and M. Bowling. “DeepStack: Expert-level artificial intelligence in heads-up no-limit poker”. en. In: *Science* 356.6337 (May 2017), pp. 508–513. URL: <http://dx.doi.org/10.1126/science.aam6960>.
- [Mor63] T. Morimoto. “Markov Processes and the H-Theorem”. In: *J. Phys. Soc. Jpn.* 18.3 (1963), pp. 328–331. URL: <https://doi.org/10.1143/JPSJ.18.328>.
- [Mos+24] R. J. Moss, A. Corso, J. Caers, and M. J. Kochenderfer. “BetaZero: Belief-state planning for long-horizon POMDPs using learned approximations”. In: *RL Conference*. 2024. URL: <https://arxiv.org/abs/2306.00249>.
- [MP+22] A. Mavor-Parker, K. Young, C. Barry, and L. Griffin. “How to Stay Curious while avoiding Noisy TVs using Aleatoric Uncertainty Estimation”. en. In: *ICML*. PMLR, June 2022, pp. 15220–15240. URL: <https://proceedings.mlr.press/v162/mavor-parker22a.html>.
- [MP95] R. D. McKelvey and T. R. Palfrey. “Quantal response equilibria for normal form games”. en. In: *Games Econ. Behav.* 10.1 (July 1995), pp. 6–38. URL: <https://www.sciencedirect.com/science/article/abs/pii/S0899825685710238>.

- [MP98] R. D. McKelvey and T. R. Palfrey. “Quantal response equilibria for extensive form games”. en. In: *Exp. Econ.* 1.1 (June 1998), pp. 9–41. URL: <https://link.springer.com/article/10.1023/A:1009905800005>.
- [Mro25] Y. Mroueh. “Reinforcement learning with verifiable rewards: GRPO’s effective loss, dynamics, and success amplification”. In: *arXiv [cs.LG]* (Mar. 2025). URL: <http://arxiv.org/abs/2503.06639>.
- [MS14] J. Modayil and R. Sutton. “Prediction driven behavior: Learning predictions that drive fixed responses”. In: *National Conference on Artificial Intelligence*. June 2014. URL: <https://www.semanticscholar.org/paper/Prediction-Driven-Behavior%3A-Learning-Predictions-Modayil-Sutton/22162abb8f5868938f8da391d3a1d603b3d8ac4c>.
- [MS24] W. Merrill and A. Sabharwal. “The Expressive Power of Transformers with Chain of Thought”. In: *ICLR*. 2024. URL: <https://arxiv.org/abs/2310.07923>.
- [MSB21] B. Millidge, A. Seth, and C. L. Buckley. “Predictive Coding: a Theoretical and Experimental Review”. In: (2021). arXiv: [2107.12979](https://arxiv.org/abs/2107.12979) [cs.AI]. URL: <http://arxiv.org/abs/2107.12979>.
- [MT24] S. Mo and S. Tong. “Connecting Joint-Embedding Predictive Architecture with Contrastive self-supervised learning”. In: *arXiv [cs.CV]* (Oct. 2024). URL: <http://arxiv.org/abs/2410.19560>.
- [Mun14] R. Munos. “From Bandits to Monte-Carlo Tree Search: The Optimistic Principle Applied to Optimization and Planning”. In: *Foundations and Trends in Machine Learning* 7.1 (2014), pp. 1–129. URL: <http://dx.doi.org/10.1561/22000000038>.
- [Mun+16] R. Munos, T. Stepleton, A. Harutyunyan, and M. G. Bellemare. “Safe and Efficient Off-Policy Reinforcement Learning”. In: *NIPS*. 2016, pp. 1046–1054.
- [Mur00] K. Murphy. *A Survey of POMDP Solution Techniques*. Tech. rep. Comp. Sci. Div., UC Berkeley, 2000. URL: <https://www.cs.ubc.ca/~murphyk/Papers/pomdp.pdf>.
- [Mur23] K. P. Murphy. *Probabilistic Machine Learning: Advanced Topics*. MIT Press, 2023. URL: <https://probml.github.io/pml-book/book2.html>.
- [MWS14] J. Modayil, A. White, and R. S. Sutton. “Multi-timescale nexting in a reinforcement learning robot”. en. In: *Adapt. Behav.* 22.2 (Apr. 2014), pp. 146–160. URL: <https://sites.ualberta.ca/~amw8/nexting.pdf>.
- [Nac+18] O. Nachum, S. Gu, H. Lee, and S. Levine. “Data-Efficient Hierarchical Reinforcement Learning”. In: *NIPS*. May 2018. URL: <https://proceedings.neurips.cc/paper/2018/hash/e6384711491713d29bc63fc5eeb5ba4f-Abstract.html>.
- [Nac+19] O. Nachum, S. Gu, H. Lee, and S. Levine. “Near-Optimal Representation Learning for Hierarchical Reinforcement Learning”. In: *ICLR*. 2019. URL: <https://openreview.net/pdf?id=H1emus0qF7>.
- [Nai+20] A. Nair, A. Gupta, M. Dalal, and S. Levine. “AWAC: Accelerating Online Reinforcement Learning with Offline Datasets”. In: *arXiv [cs.LG]* (June 2020). URL: <http://arxiv.org/abs/2006.09359>.
- [Nai+21] A. Naik, Z. Abbas, A. White, and R. S. Sutton. “Towards Reinforcement Learning in the Continuing Setting”. In: *Never-Ending Reinforcement Learning (NERL) Workshop at ICLR*. 2021. URL: https://drive.google.com/file/d/1xh7WjGP2VI_QdpjVWygRC1BuH6WB_gqi/view.
- [Nai+24] A. Naik, Y. Wan, M. Tomar, and R. S. Sutton. “Reward Centering”. In: *RL Conference*. May 2024. URL: <https://arxiv.org/abs/2405.09999>.
- [Nak+23] M. Nakamoto, Y. Zhai, A. Singh, M. S. Mark, Y. Ma, C. Finn, A. Kumar, and S. Levine. “Cal-QL: Calibrated offline RL pre-training for efficient online fine-tuning”. In: *arXiv [cs.LG]* (Mar. 2023). URL: <http://arxiv.org/abs/2303.05479>.

- [Ngu+25] X.-P. Nguyen, S. Pandit, R. G. Reddy, A. Xu, S. Savarese, C. Xiong, and S. Joty. “SFR-DeepResearch: Towards effective reinforcement learning for autonomously reasoning single agents”. In: *arXiv [cs.AI]* (Sept. 2025). URL: <http://arxiv.org/abs/2509.06283>.
- [NHR99] A. Ng, D. Harada, and S. Russell. “Policy invariance under reward transformations: Theory and application to reward shaping”. In: *ICML*. 1999.
- [Ni+24] T. Ni, B. Eysenbach, E. Seyedsalehi, M. Ma, C. Gehring, A. Mahajan, and P.-L. Bacon. “Bridging State and History Representations: Understanding Self-Predictive RL”. In: *ICLR*. Jan. 2024. URL: <http://arxiv.org/abs/2401.08898>.
- [Nik+22] E. Nikishin, R. Abachi, R. Agarwal, and P.-L. Bacon. “Control-oriented model-based reinforcement learning with implicit differentiation”. en. In: *AAAI*. Vol. 36. June 2022, pp. 7886–7894. URL: <https://arxiv.org/abs/2106.03273>.
- [NLS19] C. A. Naesseth, F. Lindsten, and T. B. Schön. “Elements of Sequential Monte Carlo”. In: *Foundations and Trends in Machine Learning* (2019). URL: <http://arxiv.org/abs/1903.04797>.
- [NR00] A. Ng and S. Russell. “Algorithms for inverse reinforcement learning”. In: *ICML*. 2000.
- [NWJ10] X. Nguyen, M. J. Wainwright, and M. I. Jordan. “Estimating Divergence Functionals and the Likelihood Ratio by Convex Risk Minimization”. In: *IEEE Trans. Inf. Theory* 56.11 (2010), pp. 5847–5861. URL: <http://dx.doi.org/10.1109/TIT.2010.2068870>.
- [OA14] F. A. Oliehoek and C. Amato. “Best Response Bayesian Reinforcement Learning for Multiagent Systems with State Uncertainty”. en. In: *Proceedings of the Ninth AAMAS Workshop on Multi-Agent Sequential Decision Making in Uncertain Domains (MSDM)*. University of Liverpool, 2014. URL: <https://livrepository.liverpool.ac.uk/3000453/1/Oliehoek14MSDM.pdf>.
- [OA16] F. A. Oliehoek and C. Amato. *A Concise Introduction to Decentralized POMDPs*. en. 1st ed. SpringerBriefs in Intelligent Systems. Cham, Switzerland: Springer International Publishing, June 2016. URL: <https://www.fransoliehoek.net/docs/OliehoekAmato16book.pdf>.
- [OCD21] G. Ostrovski, P. S. Castro, and W. Dabney. “The Difficulty of Passive Learning in Deep Reinforcement Learning”. In: *NIPS*. Vol. 34. Dec. 2021, pp. 23283–23295. URL: https://proceedings.neurips.cc/paper_files/paper/2021/file/c3e0c62ee91db8dc7382bde7419bb573-Paper.pdf.
- [O’D+16] B. O’Donoghue, R. Munos, K. Kavukcuoglu, and V. Mnih. “Combining policy gradient and Q-learning”. In: *arXiv [cs.LG]* (Nov. 2016). URL: <http://arxiv.org/abs/1611.01626>.
- [OK22] A. Ororbia and D. Kifer. “The neural coding framework for learning generative models”. en. In: *Nat. Commun.* 13.1 (Apr. 2022), p. 2064. URL: <https://www.nature.com/articles/s41467-022-29632-7>.
- [Ope+20] OpenAI et al. “Asymmetric self-play for automatic goal discovery in robotic manipulation”. In: *NIPS Deep RL Workshop*. 2020. URL: <https://arxiv.org/abs/2101.04882>.
- [Oqu+24] M. Oquab et al. “DINOv2: Learning Robust Visual Features without Supervision”. In: *Transactions on Machine Learning Research* (2024). URL: <https://openreview.net/forum?id=a68SUt6zFt>.
- [Ort+21] P. A. Ortega et al. “Shaking the foundations: delusions in sequence models for interaction and control”. In: *arXiv [cs.LG]* (Oct. 2021). URL: <http://arxiv.org/abs/2110.10819>.
- [ORVR13] I. Osband, D. Russo, and B. Van Roy. “(More) Efficient Reinforcement Learning via Posterior Sampling”. In: *NIPS*. 2013. URL: <http://arxiv.org/abs/1306.0940>.
- [Osb+16] I. Osband, C. Blundell, A. Pritzel, and B. Van Roy. “Deep Exploration via Bootstrapped DQN”. In: *NIPS*. 2016. URL: <http://arxiv.org/abs/1602.04621>.

- [Osb+19] I. Osband, B. Van Roy, D. J. Russo, and Z. Wen. “Deep exploration via randomized value functions”. In: *JMLR* 20.124 (2019), pp. 1–62. URL: <http://jmlr.org/papers/v20/18-339.html>.
- [Osb+23a] I. Osband, Z. Wen, S. M. Asghari, V. Dwaracherla, M. Ibrahimi, X. Lu, and B. Van Roy. “Approximate Thompson Sampling via Epistemic Neural Networks”. en. In: *UAI*. PMLR, July 2023, pp. 1586–1595. URL: <https://proceedings.mlr.press/v216/osband23a.html>.
- [Osb+23b] I. Osband, Z. Wen, S. M. Asghari, V. Dwaracherla, M. Ibrahimi, X. Lu, and B. Van Roy. “Epistemic Neural Networks”. In: *NIPS*. 2023. URL: https://proceedings.neurips.cc/paper_files/paper/2023/file/07fbde96bee50f4e09303fd4f877c2f3-Paper-Conference.pdf.
- [OSL17] J. Oh, S. Singh, and H. Lee. “Value Prediction Network”. In: *NIPS*. July 2017.
- [Osw+25] J. von Oswald et al. “MesaNet: Sequence modeling by locally optimal test-time training”. In: *arXiv [cs.LG]* (June 2025). URL: <http://arxiv.org/abs/2506.05233>.
- [OT22] M. Okada and T. Taniguchi. “DreamingV2: Reinforcement learning with discrete world models without reconstruction”. en. In: *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. IEEE, Oct. 2022, pp. 985–991. URL: <https://ieeexplore.ieee.org/abstract/document/9981405>.
- [OVR17] I. Osband and B. Van Roy. “Why is posterior sampling better than optimism for reinforcement learning?” In: *ICML*. 2017, pp. 2701–2710.
- [Par+23] J. S. Park, J. O’Brien, C. J. Cai, M. R. Morris, P. Liang, and M. S. Bernstein. “Generative agents: Interactive simulacra of human behavior”. en. In: *Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology*. New York, NY, USA: ACM, Oct. 2023. URL: <https://dl.acm.org/doi/10.1145/3586183.3606763>.
- [Par+24a] S. Park, K. Frans, B. Eysenbach, and S. Levine. “OGBench: Benchmarking Offline Goal-Conditioned RL”. In: *arXiv [cs.LG]* (Oct. 2024). URL: <http://arxiv.org/abs/2410.20092>.
- [Par+24b] S. Park, K. Frans, S. Levine, and A. Kumar. “Is value learning really the main bottleneck in offline RL?” In: *NIPS*. June 2024. URL: <https://arxiv.org/abs/2406.09329>.
- [Pat+17] D. Pathak, P. Agrawal, A. A. Efros, and T. Darrell. “Curiosity-driven Exploration by Self-supervised Prediction”. In: *ICML*. 2017. URL: <http://arxiv.org/abs/1705.05363>.
- [Pat+22] S. Pateria, B. Subagdja, A.-H. Tan, and C. Quek. “Hierarchical Reinforcement Learning: A comprehensive survey”. en. In: *ACM Comput. Surv.* 54.5 (June 2022), pp. 1–35. URL: <https://dl.acm.org/doi/10.1145/3453160>.
- [Pat+24] A. Patterson, S. Neumann, M. White, and A. White. “Empirical design in reinforcement learning”. In: *JMLR* (2024). URL: <http://arxiv.org/abs/2304.01315>.
- [PB+14] N. Parikh, S. Boyd, et al. “Proximal algorithms”. In: *Foundations and Trends in Optimization* 1.3 (2014), pp. 127–239.
- [PCA21] G. Papoudakis, F. Christianos, and S. V. Albrecht. “Agent modelling under partial observability for deep reinforcement learning”. In: *NIPS*. 2021. URL: <https://arxiv.org/abs/2006.09447>.
- [Pea94] B. A. Pearlmutter. “Fast Exact Multiplication by the Hessian”. In: *Neural Comput.* 6.1 (1994), pp. 147–160. URL: <https://doi.org/10.1162/neco.1994.6.1.147>.
- [Pen+19] X. B. Peng, A. Kumar, G. Zhang, and S. Levine. “Advantage-weighted regression: Simple and scalable off-policy reinforcement learning”. In: *arXiv [cs.LG]* (Sept. 2019). URL: <http://arxiv.org/abs/1910.00177>.
- [Pet08] A. S. I. R. Petersen. “Formulas for Discrete Time LQR, LQG, LEQG and Minimax LQG Optimal Control Problems”. In: *IFAC Proceedings Volumes* 41.2 (Jan. 2008), pp. 8773–8778. URL: <https://www.sciencedirect.com/science/article/pii/S1474667016403629>.

- [Pet+25] T. Pethick, W. Xie, K. Antonakopoulos, Z. Zhu, A. Silveti-Falls, and V. Cevher. “Training deep learning models with norm-constrained LMOs”. In: *arXiv [cs.LG]* (Feb. 2025). URL: <http://arxiv.org/abs/2502.07529>.
- [Pfa+25] D. Pfau, I. Davies, D. Borsa, J. Araujo, B. Tracey, and H. van Hasselt. “Wasserstein Policy Optimization”. In: *ICML*. May 2025. URL: <https://arxiv.org/abs/2505.00663>.
- [PH+23] J. Parker-Holder, M. Jiang, M. Dennis, M. Samvelyan, J. Foerster, E. Grefenstette, and T. Rocktäschel. “Evolving Curricula with Regret-Based Environment Design”. In: *ICLR*. 2023. URL: <https://arxiv.org/abs/2203.01302>.
- [Pic+19] A. Piche, V. Thomas, C. Ibrahim, Y. Bengio, and C. Pal. “Probabilistic Planning with Sequential Monte Carlo methods”. In: *ICLR*. 2019. URL: <https://openreview.net/pdf?id=ByetGn0cYX>.
- [Pir+24] W. T. Piriyaakulkij, C. Langenfeld, T. A. Le, and K. Ellis. “Doing experiments and revising rules with natural language and probabilistic reasoning”. In: *NIPS*. Feb. 2024. URL: <https://arxiv.org/abs/2402.06025>.
- [Pir+25] W. T. Piriyaakulkij, Y. Liang, H. Tang, A. Weller, M. Kryven, and K. Ellis. “PoE-world: Compositional world modeling with products of programmatic experts”. In: *NIPS*. Oct. 2025. URL: <https://arxiv.org/abs/2505.10819>.
- [PKP21] A. Plaat, W. Kusters, and M. Preuss. “High-Accuracy Model-Based Reinforcement Learning, a Survey”. In: (2021). arXiv: [2107.08241](https://arxiv.org/abs/2107.08241) [cs.LG]. URL: <http://arxiv.org/abs/2107.08241>.
- [Pla22] A. Plaat. *Deep reinforcement learning, a textbook*. Berlin, Germany: Springer, Jan. 2022. URL: <https://link.springer.com/10.1007/978-981-19-0638-1>.
- [PLG23] B. Prystawski, M. Y. Li, and N. D. Goodman. “Why think step-by-step? Reasoning emerges from the locality of experience”. In: *NIPS*. Vol. abs/2304.03843. Apr. 2023, pp. 70926–70947. URL: https://proceedings.neurips.cc/paper_files/paper/2023/hash/e0af79ad53a336b4c4b4f7e2a68eb609-Abstract-Conference.html.
- [PMB22] K. Paster, S. McIlraith, and J. Ba. “You can’t count on luck: Why decision transformers and RvS fail in stochastic environments”. In: *NIPS*. May 2022. URL: <http://arxiv.org/abs/2205.15967>.
- [Pom89] D. Pomerleau. “ALVINN: An Autonomous Land Vehicle in a Neural Network”. In: *NIPS*. 1989, pp. 305–313.
- [Por+20] R. Portelas, C. Colas, L. Weng, K. Hofmann, and P.-Y. Oudeyer. “Automatic Curriculum Learning for deep RL: A short survey”. In: *IJCAI*. Mar. 2020. URL: <https://arxiv.org/abs/2003.04664>.
- [Pow19] W. B. Powell. “From Reinforcement Learning to Optimal Control: A unified framework for sequential decisions”. In: *arXiv [cs.AI]* (Dec. 2019). URL: <http://arxiv.org/abs/1912.03513>.
- [Pow22] W. B. Powell. *Reinforcement Learning and Stochastic Optimization: A Unified Framework for Sequential Decisions*. en. 1st ed. Wiley, Mar. 2022. URL: <https://www.amazon.com/Reinforcement-Learning-Stochastic-Optimization-Sequential/dp/1119815037>.
- [PR12] W. B. Powell and I. O. Ryzhov. *Optimal Learning*. Wiley Series in Probability and Statistics. <http://optimalllearning.princeton.edu/>. Hoboken, NJ: Wiley-Blackwell, Mar. 2012. URL: <https://castle.princeton.edu/wp-content/uploads/2019/02/Powell-OptimalLearningWileyMarch112018.pdf>.
- [PS07] J. Peters and S. Schaal. “Reinforcement Learning by Reward-Weighted Regression for Operational Space Control”. In: *ICML*. 2007, pp. 745–750.
- [PSS00] D. Precup, R. S. Sutton, and S. P. Singh. “Eligibility Traces for Off-Policy Policy Evaluation”. In: *ICML*. ICML ’00. Morgan Kaufmann Publishers Inc., 2000, pp. 759–766. URL: <http://dl.acm.org/citation.cfm?id=645529.658134>.

- [PT87] C. Papadimitriou and J. Tsitsiklis. “The complexity of Markov decision processes”. In: *Mathematics of Operations Research* 12.3 (1987), pp. 441–450.
- [Pte+24] M. Pternea, P. Singh, A. Chakraborty, Y. Oruganti, M. Milletari, S. Bapat, and K. Jiang. “The RL/LLM taxonomy tree: Reviewing synergies between Reinforcement Learning and Large Language Models”. In: *JAIR* (Feb. 2024). URL: <https://www.jair.org/index.php/jair/article/view/15960>.
- [Pur+25] I. Puri, S. Sudalairaj, G. Xu, K. Xu, and A. Srivastava. “A probabilistic inference approach to inference-time scaling of LLMs using particle-based Monte Carlo methods”. In: *arXiv [cs.LG]* (Feb. 2025). URL: <http://arxiv.org/abs/2502.01618>.
- [Put94] M. L. Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*. Wiley, 1994.
- [PW94] J. Peng and R. J. Williams. “Incremental Multi-Step Q-Learning”. In: *Machine Learning Proceedings*. Elsevier, Jan. 1994, pp. 226–232. URL: <http://dx.doi.org/10.1016/B978-1-55860-335-6.50035-0>.
- [QPC21] J. Queeney, I. C. Paschalidis, and C. G. Cassandras. “Generalized Proximal Policy Optimization with Sample Reuse”. In: *NIPS*. Oct. 2021.
- [QPC24] J. Queeney, I. C. Paschalidis, and C. G. Cassandras. “Generalized Policy Improvement algorithms with theoretically supported sample reuse”. In: *IEEE Trans. Automat. Contr.* (2024). URL: <http://arxiv.org/abs/2206.13714>.
- [Rac+20] S. Racaniere, A. K. Lampinen, A. Santoro, D. P. Reichert, V. Firoiu, and T. P. Lillicrap. “Automated curricula through setter-solver interactions”. In: *ICLR*. 2020. URL: <https://arxiv.org/abs/1909.12892>.
- [Raf+23] R. Rafailov, A. Sharma, E. Mitchell, S. Ermon, C. D. Manning, and C. Finn. “Direct Preference Optimization: Your language model is secretly a reward model”. In: *arXiv [cs.LG]* (May 2023). URL: <http://arxiv.org/abs/2305.18290>.
- [Raf+24] R. Rafailov et al. “D5RL: Diverse datasets for data-driven deep reinforcement learning”. In: *RLC*. Aug. 2024. URL: <https://arxiv.org/abs/2408.08441>.
- [Raj+17] A. Rajeswaran, K. Lowrey, E. Todorov, and S. Kakade. “Towards generalization and simplicity in continuous control”. In: *NIPS*. Mar. 2017.
- [Rao10] A. V. Rao. “A Survey of Numerical Methods for Optimal Control”. In: *Adv. Astronaut. Sci.* 135.1 (2010). URL: <http://dx.doi.org/>.
- [Ras+18] T. Rashid, M. Samvelyan, C. S. de Witt, G. Farquhar, J. Foerster, and S. Whiteson. “QMIX: Monotonic value function factorisation for deep multi-agent reinforcement learning”. In: *ICML*. Mar. 2018. URL: <https://arxiv.org/abs/1803.11485>.
- [RB12] S. Ross and J. A. Bagnell. “Agnostic system identification for model-based reinforcement learning”. In: *ICML*. Mar. 2012.
- [RB99] R. P. Rao and D. H. Ballard. “Predictive coding in the visual cortex: a functional interpretation of some extra-classical receptive-field effects”. en. In: *Nat. Neurosci.* 2.1 (1999), pp. 79–87. URL: <http://dx.doi.org/10.1038/4580>.
- [RCdP07] S. Ross, B. Chaib-draa, and J. Pineau. “Bayes-Adaptive POMDPs”. In: *NIPS* 20 (2007). URL: https://proceedings.neurips.cc/paper_files/paper/2007/file/3b3dbaf68507998acd6a5a5254ab2d76-Paper.pdf.
- [Rec19] B. Recht. “A Tour of Reinforcement Learning: The View from Continuous Control”. In: *Annual Review of Control, Robotics, and Autonomous Systems* 2 (2019), pp. 253–279. URL: <http://arxiv.org/abs/1806.09460>.
- [Ree+22] S. Reed et al. “A Generalist Agent”. In: *TMLR* (May 2022). URL: <https://arxiv.org/abs/2205.06175>.

- [Ren+24] A. Z. Ren, J. Lidard, L. L. Ankile, A. Simeonov, P. Agrawal, A. Majumdar, B. Burchfiel, H. Dai, and M. Simchowitz. “Diffusion Policy Policy Optimization”. In: *arXiv [cs.RO]* (Aug. 2024). URL: <http://arxiv.org/abs/2409.00588>.
- [Ren+25] X. Ren, C. Forder, Q. Zang, A. Tahir, R. J. Georgio, S. Deb, P. Carroll, O. Gurcan, and Z. Guo. “Anemol: A semi-centralized multi-agent system based on agent-to-agent communication MCP server from Coral Protocol”. In: *arXiv [cs.MA]* (Aug. 2025). URL: <http://arxiv.org/abs/2508.17068>.
- [RFP15] I. O. Ryzhov, P. I. Frazier, and W. B. Powell. “A new optimal stepsize for approximate dynamic programming”. en. In: *IEEE Trans. Automat. Contr.* 60.3 (Mar. 2015), pp. 743–758. URL: <https://castle.princeton.edu/Papers/Ryzhov-OptimalStepsizeforADPFeb242015.pdf>.
- [RG66] A. Rapoport and M. Guyer. “A Taxonomy of 2 X 2 Games”. In: *General System: Yearbook of the Society for General Systems Research* 11 (1966), pp. 203–214.
- [RGB11] S. Ross, G. J. Gordon, and J. A. Bagnell. “A reduction of imitation learning and structured prediction to no-regret online learning”. In: *AISTATS*. 2011.
- [RHH23] J. Robine, M. Höftmann, and S. Harmeling. “A simple framework for self-supervised learning of sample-efficient world models”. In: *NIPS SSL Workshop*. 2023, pp. 17–18. URL: https://sslneurips23.github.io/paper_pdfs/paper_44.pdf.
- [RHH25] J. Robine, M. Hoftmann, and S. Harmeling. “Simple, Good, Fast: Self-Supervised World Models Free of Baggage”. In: *The Thirteenth International Conference on Learning Representations*. 2025. URL: <https://openreview.net/pdf?id=yFGR36PLDJ>.
- [Ric+25] J. Richens, D. Abel, A. Bellot, and T. Everitt. “General agents need world models”. In: *ICML*. June 2025. URL: <https://arxiv.org/abs/2506.01622>.
- [Rie05] M. Riedmiller. “Neural fitted Q iteration – first experiences with a data efficient neural reinforcement learning method”. en. In: *ECML. Lecture notes in computer science*. 2005, pp. 317–328. URL: https://link.springer.com/chapter/10.1007/11564096_32.
- [Rin21] M. Ring. “Representing knowledge as predictions (and state as knowledge)”. In: *arXiv [cs.AI]* (Dec. 2021). URL: <http://arxiv.org/abs/2112.06336>.
- [RJ22] A. Rao and T. Jelvis. *Foundations of Reinforcement Learning with Applications in Finance*. Chapman and Hall/ CRC, 2022. URL: <https://github.com/TikhonJelvis/RL-book>.
- [RK04] R. Rubinstein and D. Kroese. *The Cross-Entropy Method: A Unified Approach to Combinatorial Optimization, Monte-Carlo Simulation, and Machine Learning*. Springer-Verlag, 2004.
- [RLT18] M. Riemer, M. Liu, and G. Tesauro. “Learning Abstract Options”. In: *NIPS* 31 (2018). URL: https://proceedings.neurips.cc/paper_files/paper/2018/file/cdf28f8b7d14ab02d12a2329d71e4079-Paper.pdf.
- [RMD22] J. B. Rawlings, D. Q. Mayne, and M. M. Diehl. *Model Predictive Control: Theory, Computation, and Design (2nd ed)*. en. Nob Hill Publishing, LLC, Sept. 2022. URL: <https://sites.engineering.ucsb.edu/~jbrow/mpc/MPC-book-2nd-edition-1st-printing.pdf>.
- [RMK20] A. Rajeswaran, I. Mordatch, and V. Kumar. “A game theoretic framework for model based reinforcement learning”. In: *ICML*. 2020.
- [RN94] G. A. Rummery and M. Niranjan. *On-Line Q-Learning Using Connectionist Systems*. Tech. rep. Cambridge Univ. Engineering Dept., 1994. URL: <http://dx.doi.org/>.
- [Roc+24] J. Rocamonde, V. Montesinos, E. Nava, E. Perez, and D. Lindner. “Vision-language models are zero-shot reward models for reinforcement learning”. In: *ICLR*. 2024. URL: <https://arxiv.org/abs/2310.12921>.
- [RP+24] B. Romera-Paredes et al. “Mathematical discoveries from program search with large language models”. In: *Nature* (2024).

- [RR14] D. Russo and B. V. Roy. “Learning to Optimize via Posterior Sampling”. In: *Math. Oper. Res.* 39.4 (2014), pp. 1221–1243.
- [RTV12] K. Rawlik, M. Toussaint, and S. Vijayakumar. “On stochastic optimal control and reinforcement learning by approximate inference”. In: *Robotics: Science and Systems VIII*. Robotics: Science and Systems Foundation, 2012. URL: https://blogs.cuit.columbia.edu/zp2130/files/2019/03/On_Stochastic_Optimal_Control_and_Reinforcement_Learning_by_Approximate_Inference.pdf.
- [Rub97] R. Y. Rubinstein. “Optimization of computer simulation models with rare events”. In: *Eur. J. Oper. Res.* 99.1 (1997), pp. 89–112. URL: <http://www.sciencedirect.com/science/article/pii/S0377221796003852>.
- [Rud+25] M. Rudolph, N. Lichtle, S. Mohammadpour, A. Bayen, J. Z. Kolter, A. Zhang, G. Farina, E. Vinitzky, and S. Sokota. “Reevaluating policy gradient methods for imperfect-information games”. In: *arXiv [cs.LG]* (Feb. 2025). URL: <http://arxiv.org/abs/2502.08938>.
- [Rus+18] D. J. Russo, B. Van Roy, A. Kazerouni, I. Osband, and Z. Wen. “A Tutorial on Thompson Sampling”. In: *Foundations and Trends in Machine Learning* 11.1 (2018), pp. 1–96. URL: <http://dx.doi.org/10.1561/22000000070>.
- [Rus19] S. Russell. *Human Compatible: Artificial Intelligence and the Problem of Control*. en. Kindle. Viking, 2019. URL: https://www.amazon.com/Human-Compatible-Artificial-Intelligence-Problem-ebook/dp/B07N5J5FTS/ref=zg_bs_3887_4?_encoding=UTF8&psc=1&refRID=0JE0ST011W4K15PTFZAT.
- [RW91] S. Russell and E. Wefald. “Principles of metareasoning”. en. In: *Artif. Intell.* 49.1-3 (May 1991), pp. 361–395. URL: [http://dx.doi.org/10.1016/0004-3702\(91\)90015-C](http://dx.doi.org/10.1016/0004-3702(91)90015-C).
- [Ryu+20] M. Ryu, Y. Chow, R. Anderson, C. Tjandraatmadja, and C. Boutilier. “CAQL: Continuous Action Q-Learning”. In: *ICLR*. 2020. URL: <https://openreview.net/forum?id=BkxXe0Etwr>.
- [Saj+21] N. Sajid, P. J. Ball, T. Parr, and K. J. Friston. “Active Inference: Demystified and Compared”. en. In: *Neural Comput.* 33.3 (Mar. 2021), pp. 674–712. URL: https://web.archive.org/web/20210628163715id_/https://discovery.ucl.ac.uk/id/eprint/10119277/1/Friston_neco_a_01357.pdf.
- [Sal+17] T. Salimans, J. Ho, X. Chen, and I. Sutskever. “Evolution Strategies as a Scalable Alternative to Reinforcement Learning”. In: (2017). arXiv: 1703.03864 [stat.ML]. URL: <http://arxiv.org/abs/1703.03864>.
- [Sal+23] T. Salvatori, A. Mali, C. L. Buckley, T. Lukasiewicz, R. P. N. Rao, K. Friston, and A. Ororbia. “Brain-inspired computational intelligence via predictive coding”. In: *arXiv [cs.AI]* (Aug. 2023). URL: <http://arxiv.org/abs/2308.07870>.
- [Sal+24] T. Salvatori, Y. Song, Y. Yordanov, B. Millidge, L. Sha, C. Emde, Z. Xu, R. Bogacz, and T. Lukasiewicz. “A Stable, Fast, and Fully Automatic Learning Algorithm for Predictive Coding Networks”. In: *ICLR*. Oct. 2024. URL: <https://openreview.net/pdf?id=RyUvzda8GH>.
- [Sam+19] M. Samvelyan, T. Rashid, C. S. de Witt, G. Farquhar, N. Nardelli, T. G. J. Rudner, C.-M. Hung, P. H. S. Torr, J. Foerster, and S. Whiteson. “The StarCraft Multi-Agent Challenge”. In: *arXiv [cs.LG]* (Feb. 2019). URL: <http://arxiv.org/abs/1902.04043>.
- [Sar+25] B. Sarkar et al. “Evolution Strategies at the Hyperscale”. In: *arXiv [cs.LG]* (Nov. 2025). URL: <http://arxiv.org/abs/2511.16652>.
- [SB08] O. Simsek and A. G. Barto. “Skill characterization based on betweenness”. In: (2008).
- [SB18] R. Sutton and A. Barto. *Reinforcement learning: an introduction (2nd edn)*. MIT Press, 2018.
- [Sch10] J. Schmidhuber. “Formal Theory of Creativity, Fun, and Intrinsic Motivation”. In: *IEEE Trans. Autonomous Mental Development* 2 (2010). URL: <http://people.idsia.ch/~juergen/ieeecreative.pdf>.

- [Sch12] L. Schulz. “The origins of inquiry: inductive inference and exploration in early childhood”. en. In: *Trends Cogn. Sci.* 16.7 (July 2012), pp. 382–389. URL: <https://www.cell.com/action/showAbstract?pii=S1364661312001301>.
- [Sch+15a] T. Schaul, D. Horgan, K. Gregor, and D. Silver. “Universal Value Function Approximators”. en. In: *ICML*. PMLR, June 2015, pp. 1312–1320. URL: <https://proceedings.mlr.press/v37/schaul15.html>.
- [Sch+15b] J. Schulman, S. Levine, P. Moritz, M. I. Jordan, and P. Abbeel. “Trust Region Policy Optimization”. In: *ICML*. 2015. URL: <http://arxiv.org/abs/1502.05477>.
- [Sch+16a] T. Schaul, J. Quan, I. Antonoglou, and D. Silver. “Prioritized Experience Replay”. In: *ICLR*. 2016. URL: <http://arxiv.org/abs/1511.05952>.
- [Sch+16b] J. Schulman, P. Moritz, S. Levine, M. Jordan, and P. Abbeel. “High-Dimensional Continuous Control Using Generalized Advantage Estimation”. In: *ICLR*. 2016. URL: <http://arxiv.org/abs/1506.02438>.
- [Sch+17] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. “Proximal Policy Optimization Algorithms”. In: (2017). arXiv: [1707.06347](https://arxiv.org/abs/1707.06347) [cs.LG]. URL: <http://arxiv.org/abs/1707.06347>.
- [Sch+19] M. Schlegel, W. Chung, D. Graves, J. Qian, and M. White. “Importance Resampling for Off-policy Prediction”. In: *NIPS*. June 2019. URL: <https://arxiv.org/abs/1906.04328>.
- [Sch19] J. Schmidhuber. “Reinforcement learning Upside Down: Don’t predict rewards – just map them to actions”. In: *arXiv [cs.AI]* (Dec. 2019). URL: <http://arxiv.org/abs/1912.02875>.
- [Sch+20] J. Schrittwieser et al. “Mastering Atari, Go, Chess and Shogi by Planning with a Learned Model”. In: *Nature* (2020). URL: <http://arxiv.org/abs/1911.08265>.
- [Sch+21a] M. Schmid et al. “Student of Games: A unified learning algorithm for both perfect and imperfect information games”. In: *Sci. Adv.* (Dec. 2021). URL: <https://www.science.org/doi/10.1126/sciadv.adg3256>.
- [Sch+21b] J. Schrittwieser, T. K. Hubert, A. Mandhane, M. Barekatain, I. Antonoglou, and D. Silver. “Online and Offline Reinforcement Learning by Planning with a Learned Model”. In: *NIPS*. Nov. 2021. URL: <https://openreview.net/forum?id=HKtsGW-lNbw>.
- [Sch+21c] M. Schwarzer, A. Anand, R. Goel, R. Devon Hjelm, A. Courville, and P. Bachman. “Data-Efficient Reinforcement Learning with Self-Predictive Representations”. In: *ICLR*. 2021. URL: <https://openreview.net/pdf?id=uCQfPZwRaUu>.
- [Sch+23a] I. Schubert, J. Zhang, J. Bruce, S. Behtle, E. Parisotto, M. Riedmiller, J. T. Springenberg, A. Byravan, L. Hasenclever, and N. Heess. “A Generalist Dynamics Model for Control”. In: *arXiv [cs.AI]* (May 2023). URL: <http://arxiv.org/abs/2305.10912>.
- [Sch+23b] M. Schwarzer, J. Obando-Ceron, A. Courville, M. Bellemare, R. Agarwal, and P. S. Castro. “Bigger, Better, Faster: Human-level Atari with human-level efficiency”. In: *ICML*. May 2023. URL: <http://arxiv.org/abs/2305.19452>.
- [Sch+24] M. Schneider, R. Krug, N. Vaskevicius, L. Palmieri, and J. Boedecker. “The Surprising Ineffectiveness of Pre-Trained Visual Representations for Model-Based Reinforcement Learning”. In: *NIPS*. Nov. 2024. URL: <https://openreview.net/pdf?id=LvAy07mCxU>.
- [Sch+25] T. Schmied, T. Adler, V. P. Patil, M. Beck, K. Pöppel, J. Brandstetter, G. Klambauer, R. Pascanu, and S. Hochreiter. “A Large Recurrent Action Model: xLSTM enables Fast Inference for Robotics Tasks”. In: *ICML*. June 2025. URL: <https://openreview.net/forum?id=J2JxJOP1LI>.
- [Sco10] S. Scott. “A modern Bayesian look at the multi-armed bandit”. In: *Applied Stochastic Models in Business and Industry* 26 (2010), pp. 639–658.

- [Sei+16] H. van Seijen, A Rupam Mahmood, P. M. Pilarski, M. C. Machado, and R. S. Sutton. “True Online Temporal-Difference Learning”. In: *JMLR* (2016). URL: <http://jmlr.org/papers/volume17/15-599/15-599.pdf>.
- [Sek+20] R. Sekar, O. Rybkin, K. Daniilidis, P. Abbeel, D. Hafner, and D. Pathak. “Planning to explore via self-supervised world models”. In: *ICML*. May 2020. URL: <https://arxiv.org/abs/2005.05960>.
- [Sey+22] T. Seyde, P. Werner, W. Schwarting, I. Gilitschenski, M. Riedmiller, D. Rus, and M. Wulfmeier. “Solving Continuous Control via Q-learning”. In: *ICLR*. Sept. 2022. URL: <https://openreview.net/pdf?id=U5X0GxAgccS>.
- [Sha+20] R. Shah, P. Freire, N. Alex, R. Freedman, D. Krasheninnikov, L. Chan, M. D. Dennis, P. Abbeel, A. Dragan, and S. Russell. “Benefits of Assistance over Reward Learning”. In: *NIPS Workshop*. 2020. URL: <https://aima.cs.berkeley.edu/~russell/papers/neurips20ws-assistance.pdf>.
- [Sha+24] Z. Shao, P. Wang, Q. Zhu, R. Xu, J. Song, M. Zhang, Y. K. Li, Y. Wu, and D. Guo. “DeepSeek-Math: Pushing the limits of mathematical reasoning in open language models”. In: *arXiv [cs.CL]* (Feb. 2024). URL: <http://arxiv.org/abs/2402.03300>.
- [Sha+25] Y. Shang et al. “A survey of embodied world models”. In: (Sept. 2025). URL: <https://www.researchgate.net/doi/10.13140/RG.2.2.28036.39041>.
- [Sha53a] L. Shapley. “A value for n-person games”. en. In: *Classics in Game Theory*. Princeton University Press, 1953, pp. 69–79. URL: <https://www.degruyterbrill.com/document/doi/10.1515/9781400829156-012/html>.
- [Sha53b] L. S. Shapley. “Stochastic games”. en. In: *Proc. Natl. Acad. Sci. U. S. A.* 39.10 (Oct. 1953), pp. 1095–1100. URL: <https://pnas.org/doi/full/10.1073/pnas.39.10.1095>.
- [SHS20] S. Schmitt, M. Hessel, and K. Simonyan. “Off-Policy Actor-Critic with Shared Experience Replay”. en. In: *ICML*. PMLR, Nov. 2020, pp. 8545–8554. URL: <https://proceedings.mlr.press/v119/schmitt20a.html>.
- [Sie+20] N. Siegel, J. T. Springenberg, F. Berkenkamp, A. Abdolmaleki, M. Neunert, T. Lampe, R. Hafner, N. Heess, and M. Riedmiller. “Keep Doing What Worked: Behavior Modelling Priors for Offline Reinforcement Learning”. In: *ICLR*. 2020. URL: <https://openreview.net/pdf?id=rke7geHtwH>.
- [Sil+14] D. Silver, G. Lever, N. Heess, T. Degris, D. Wierstra, and M. Riedmiller. “Deterministic Policy Gradient Algorithms”. In: *ICML*. ICML’14. JMLR.org, 2014, pp. I–387–I–395. URL: <http://dl.acm.org/citation.cfm?id=3044805.3044850>.
- [Sil+16] D. Silver et al. “Mastering the game of Go with deep neural networks and tree search”. en. In: *Nature* 529.7587 (2016), pp. 484–489. URL: <http://dx.doi.org/10.1038/nature16961>.
- [Sil+17a] D. Silver et al. “Mastering the game of Go without human knowledge”. en. In: *Nature* 550.7676 (2017), pp. 354–359. URL: <http://dx.doi.org/10.1038/nature24270>.
- [Sil+17b] D. Silver et al. “The predictron: end-to-end learning and planning”. In: *ICML*. 2017. URL: <https://openreview.net/pdf?id=BkJscIcgl>.
- [Sil18] D. Silver. *Lecture 9L Exploration and Exploitation*. 2018. URL: http://www0.cs.ucl.ac.uk/staff/d.silver/web/Teaching_files/XX.pdf.
- [Sil+18] D. Silver et al. “A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play”. en. In: *Science* 362.6419 (2018), pp. 1140–1144. URL: <http://dx.doi.org/10.1126/science.aar6404>.
- [Sil+21] D. Silver, S. Singh, D. Precup, and R. S. Sutton. “Reward is enough”. en. In: *Artif. Intell.* 299.103535 (Oct. 2021), p. 103535. URL: <https://www.sciencedirect.com/science/article/pii/S0004370221000862>.

- [Sin+00] S. Singh, T. Jaakkola, M. L. Littman, and C. Szepesvári. “Convergence Results for Single-Step On-Policy Reinforcement-Learning Algorithms”. In: *MLJ* 38.3 (2000), pp. 287–308. URL: <https://doi.org/10.1023/A:1007678930559>.
- [SK18] Z. Sunberg and M. Kochenderfer. “Online algorithms for POMDPs with continuous state, action, and observation spaces”. In: *ICAPS*. 2018. URL: <https://arxiv.org/abs/1709.06196>.
- [Ska+22] J. Skalse, N. H. R. Howe, D. Krashennnikov, and D. Krueger. “Defining and characterizing reward hacking”. In: *NIPS*. Sept. 2022.
- [SKM00] S. P. Singh, M. J. Kearns, and Y. Mansour. “Nash Convergence of Gradient Dynamics in General-Sum Games”. en. In: *UAI*. June 2000, pp. 541–548. URL: <https://dl.acm.org/doi/10.5555/647234.719924>.
- [SKM18] S. Schwöbel, S. Kiebel, and D. Marković. “Active Inference, Belief Propagation, and the Bethe Approximation”. en. In: *Neural Comput.* 30.9 (2018), pp. 2530–2567. URL: http://dx.doi.org/10.1162/neco_a_01108.
- [SL25] J. Schulman and T. M. Lab. “LoRA Without Regret”. In: *Thinking Machines Lab: Connectionism* (2025). <https://thinkingmachines.ai/blog/lora/>.
- [SLB08] Y. Shoham and K. Leyton-Brown. *Multiagent systems: Algorithmic, game-theoretic, and logical foundations*. Cambridge, England: Cambridge University Press, Dec. 2008. URL: <https://www.masfoundations.org/>.
- [SLG23] N. Shinn, B. Labash, and A. Gopinath. “Reflexion: an autonomous agent with dynamic memory and self-reflection”. In: *NIPS*. Mar. 2023. URL: https://proceedings.neurips.cc/paper_files/paper/2023/hash/1b44b878bb782e6954cd888628510e90-Abstract-Conference.html.
- [Sli19] A. Slivkins. “Introduction to Multi-Armed Bandits”. In: *Foundations and Trends in Machine Learning* (2019). URL: <http://arxiv.org/abs/1904.07272>.
- [Smi+23] F. B. Smith, A. Kirsch, S. Farquhar, Y. Gal, A. Foster, and T. Rainforth. “Prediction-Oriented Bayesian Active Learning”. In: *AISTATS*. Apr. 2023. URL: <http://arxiv.org/abs/2304.08151>.
- [Sne+24] C. Snell, J. Lee, K. Xu, and A. Kumar. “Scaling LLM test-time compute optimally can be more effective than scaling model parameters”. In: *arXiv [cs.LG]* (Aug. 2024). URL: <http://arxiv.org/abs/2408.03314>.
- [Sok+21] S. Sokota, E. Lockhart, F. Timbers, E. Davoodi, R. D’Orazio, N. Burch, M. Schmid, M. Bowling, and M. Lanctot. “Solving common-payoff games with approximate policy iteration”. In: *AAAI*. Jan. 2021. URL: <https://arxiv.org/abs/2101.04237>.
- [Sok+22] S. Sokota, R. D’Orazio, J. Zico Kolter, N. Loizou, M. Lanctot, I. Mitliagkas, N. Brown, and C. Kroer. “A Unified Approach to Reinforcement Learning, Quantal Response Equilibria, and Two-Player Zero-Sum Games”. In: *ICLR*. Sept. 2022. URL: <https://openreview.net/forum?id=DpE5UYUQzZH>.
- [Sok+23] S. Sokota, G. Farina, D. J. Wu, H. Hu, K. A. Wang, J. Z. Kolter, and N. Brown. “The update-equivalence framework for decision-time planning”. In: *arXiv [cs.AI]* (Apr. 2023). URL: <http://arxiv.org/abs/2304.13138>.
- [Sol64] R. J. Solomonoff. “A formal theory of inductive inference. Part I”. In: *Information and Control* 7.1 (Mar. 1964), pp. 1–22. URL: <https://www.sciencedirect.com/science/article/pii/S0019995864902232>.
- [Som+25] S. Somvanshi, M. M. Islam, M. S. Mimi, S. B. B. Pollock, G. Chhetri, and S. Das. “From S4 to Mamba: A comprehensive survey on Structured State Space Models”. In: *arXiv [cs.LG]* (Mar. 2025). URL: <http://arxiv.org/abs/2503.18970>.

- [Son+23] S Sontakke, J. Zhang, S. M. R. Arnold, K. Pertsch, E. Biyik, D. Sadigh, C. Finn, and L. Itti. “RoboCLIP: One demonstration is enough to learn robot policies”. In: *NIPS*. Vol. abs/2310.07899. Oct. 2023, pp. 55681–55693. URL: https://proceedings.neurips.cc/paper_files/paper/2023/hash/ae54ce310476218f26dd48c1626d5187-Abstract-Conference.html.
- [Son+25] H. Song, J. Jiang, W. Tian, Z. Chen, Y. Wu, J. Zhao, Y. Min, W. X. Zhao, L. Fang, and J.-R. Wen. “R1-Searcher++: Incentivizing the Dynamic Knowledge Acquisition of LLMs via reinforcement learning”. In: *arXiv [cs.CL]* (May 2025). URL: <http://arxiv.org/abs/2505.17005>.
- [Son98] E. D. Sontag. *Mathematical Control Theory: Deterministic Finite Dimensional Systems*. 2nd. Vol. 6. Texts in Applied Mathematics. Springer, 1998.
- [Spi+24] B. A. Spiegel, Z. Yang, W. Jurayj, B. Bachmann, S. Tellex, and G. Konidaris. “Informing Reinforcement Learning Agents by Grounding Language to Markov Decision Processes”. In: *Workshop on Training Agents with Foundation Models at RLC 2024*. Aug. 2024. URL: <https://openreview.net/pdf?id=uFm9e4Ly26>.
- [Spr17] M. W. Spratling. “A review of predictive coding algorithms”. en. In: *Brain Cogn.* 112 (2017), pp. 92–97. URL: <http://dx.doi.org/10.1016/j.bandc.2015.11.003>.
- [SPS99] R. S. Sutton, D. Precup, and S. Singh. “Between MDPs and semi-MDPs: A framework for temporal abstraction in reinforcement learning”. In: *Artif. Intell.* 112.1 (Aug. 1999), pp. 181–211. URL: <http://www.sciencedirect.com/science/article/pii/S0004370299000521>.
- [Sri+18] S. Srinivasan, M. Lanctot, V. Zambaldi, J. Pérolat, K. Tuyls, R. Munos, and M. Bowling. “Actor-Critic Policy Optimization in Partially Observable Multiagent Environments”. In: *NIPS*. 2018.
- [SS21] D. Schmidt and T. Schmied. “Fast and Data-Efficient Training of Rainbow: an Experimental Study on Atari”. In: *Deep RL Workshop NeurIPS 2021*. Dec. 2021. URL: <https://openreview.net/pdf?id=GvM7A3cv63M>.
- [SS25] D. Silver and R. S. Sutton. *Welcome to the era of experience*. 2025. URL: <https://storage.googleapis.com/deepmind-media/Era-of-Experience%20/The%20Era%20of%20Experience%20Paper.pdf>.
- [SSE23] A. Shih, D. Sadigh, and S. Ermon. “Long Horizon Temperature Scaling”. In: *arXiv [cs.LG]* (Feb. 2023). URL: <http://arxiv.org/abs/2302.03686>.
- [SSM08] R. S. Sutton, C. Szepesvári, and H. R. Maei. “A convergent O(n) algorithm for off-policy temporal-difference learning with linear function approximation”. en. In: *NIPS*. NIPS’08. Red Hook, NY, USA: Curran Associates Inc., Dec. 2008, pp. 1609–1616. URL: https://proceedings.neurips.cc/paper_files/paper/2008/file/e0c641195b27425bb056ac56f8953d24-Paper.pdf.
- [SSTH22] S. Schmitt, J. Shawe-Taylor, and H. van Hasselt. “Chaining value functions for off-policy learning”. en. In: *AAAI*. Vol. 36. Association for the Advancement of Artificial Intelligence (AAAI), June 2022, pp. 8187–8195. URL: <https://ojs.aaai.org/index.php/AAAI/article/view/20792>.
- [SSTVH23] S. Schmitt, J. Shawe-Taylor, and H. Van Hasselt. “Exploration via Epistemic Value Estimation”. en. In: *AAAI* 37.8 (June 2023), pp. 9742–9751. URL: <https://ojs.aaai.org/index.php/AAAI/article/view/26164>.
- [Str00] M. Strens. “A Bayesian Framework for Reinforcement Learning”. In: *ICML*. 2000.
- [Sub+22] J. Subramanian, A. Sinha, R. Seraj, and A. Mahajan. “Approximate information state for approximate planning and reinforcement learning in partially observed systems”. In: *JMLR* 23.12 (2022), pp. 1–83. URL: <http://jmlr.org/papers/v23/20-1165.html>.
- [Suk+18] S. Sukhbaatar, E. Denton, A. Szlam, and R. Fergus. “Learning goal embeddings via self-play for hierarchical reinforcement learning”. In: *arXiv [cs.LG]* (Nov. 2018). URL: <http://arxiv.org/abs/1811.09083>.

- [Sum+24] T. R. Sumers, S. Yao, K. Narasimhan, and T. L. Griffiths. “Cognitive Architectures for Language Agents”. In: *Trans. on Machine Learning Research* (2024). URL: <https://openreview.net/forum?id=1i6ZCvflQJ>.
- [Sun+17] P. Sunehag et al. “Value-decomposition networks for cooperative multi-agent learning”. In: *arXiv [cs.AI]* (June 2017). URL: <http://arxiv.org/abs/1706.05296>.
- [Sun+24] F.-Y. Sun, S. I. Harini, A. Yi, Y. Zhou, A. Zook, J. Tremblay, L. Cross, J. Wu, and N. Haber. “FactorSim: Generative Simulation via Factorized Representation”. In: *NIPS*. Nov. 2024. URL: <https://openreview.net/forum?id=wBzvYh3PRA>.
- [Sun+25] H. Sun, Y. Wu, Y. Cheng, and X. Chu. “Game theory meets large language models: A systematic survey”. In: *arXiv [cs.AI]* (Feb. 2025). URL: <http://arxiv.org/abs/2502.09053>.
- [Sut04] R. Sutton. “The reward hypothesis”. In: (2004). URL: <http://incompleteideas.net/rlai.cs.ualberta.ca/RLAI/rewardhypothesis.html>.
- [Sut+08] R. S. Sutton, C. Szepesvari, A. Geramifard, and M. P. Bowling. “Dyna-style planning with linear function approximation and prioritized sweeping”. In: *UAI*. 2008.
- [Sut+11] R. Sutton, J. Modayil, M. Delp, T. Degris, P. Pilarski, A. White, and D. Precup. “Horde: a scalable real-time architecture for learning knowledge from unsupervised sensorimotor interaction”. In: *Adapt Agent Multi-agent Syst* (May 2011), pp. 761–768. URL: <https://www.semanticscholar.org/paper/Horde%3A-a-scalable-real-time-architecture-for-from-Sutton-Modayil/50e9a441f56124b7b969e6537b66469a0e1aa707>.
- [Sut15] R. Sutton. *Introduction to RL with function approximation*. NIPS Tutorial. 2015. URL: <http://media.nips.cc/Conferences/2015/tutorialslides/SuttonIntroRL-nips-2015-tutorial.pdf>.
- [Sut22] R. S. Sutton. “The quest for a common model of the intelligent decision maker”. In: *RLDM*. Feb. 2022. URL: <https://arxiv.org/abs/2202.13252>.
- [Sut+23] R. S. Sutton, M. C. Machado, G. Zacharias Holland, D. Szepesvari, F. Timbers, B. Tanner, and A. White. “Reward-Respecting Subtasks for Model-Based Reinforcement Learning”. In: *Artificial Intelligence Journal* (2023). URL: <http://arxiv.org/abs/2202.03466>.
- [Sut88] R. Sutton. “Learning to predict by the methods of temporal differences”. In: *Machine Learning* 3.1 (1988), pp. 9–44.
- [Sut90] R. S. Sutton. “Integrated Architectures for Learning, Planning, and Reacting Based on Approximating Dynamic Programming”. In: *ICML*. Ed. by B. Porter and R. Mooney. Morgan Kaufmann, 1990, pp. 216–224. URL: <http://www.sciencedirect.com/science/article/pii/B9781558601413500304>.
- [Sut95] R. S. Sutton. “TD models: Modeling the world at a mixture of time scales”. en. In: *ICML*. Jan. 1995, pp. 531–539. URL: <https://www.sciencedirect.com/science/article/abs/pii/B9781558603776500724>.
- [Sut96] R. S. Sutton. “Generalization in Reinforcement Learning: Successful Examples Using Sparse Coarse Coding”. In: *NIPS*. Ed. by D. S. Touretzky, M. C. Mozer, and M. E. Hasselmo. MIT Press, 1996, pp. 1038–1044. URL: <http://papers.nips.cc/paper/1109-generalization-in-reinforcement-learning-successful-examples-using-sparse-coarse-coding.pdf>.
- [Sut+99] R. Sutton, D. McAllester, S. Singh, and Y. Mansour. “Policy Gradient Methods for Reinforcement Learning with Function Approximation”. In: *NIPS*. 1999.
- [SV10] D. Silver and J. Veness. “Monte-Carlo Planning in Large POMDPs”. In: *Advances in Neural Information Processing Systems*. Vol. 23. 2010. URL: <https://proceedings.neurips.cc/paper/2010/hash/edf6e1afcf9246bb0d40eb4d8027d90f-Abstract.html>.
- [SW06] J. E. Smith and R. L. Winkler. “The Optimizer’s Curse: Skepticism and Postdecision Surprise in Decision Analysis”. In: *Manage. Sci.* 52.3 (2006), pp. 311–322.

- [Sze10] C. Szepesvari. *Algorithms for Reinforcement Learning*. Morgan Claypool, 2010.
- [SZL24] R. Shwartz-Ziv and Y. LeCun. “To Compress or Not to Compress- Self-Supervised Learning and Information Theory: A Review”. In: *Entropy* (2024). URL: <http://arxiv.org/abs/2304.09355>.
- [Ta25] A. Tekparmak and andthattoo. *Python is all you need? introducing dria-agent-a*. 2025. URL: <https://huggingface.co/blog/andthattoo/dria-agent-a>.
- [Tam+16] A. Tamar, Y. Wu, G. Thomas, S. Levine, and P. Abbeel. “Value Iteration Networks”. In: *NIPS*. 2016. URL: <http://arxiv.org/abs/1602.02867>.
- [Tan+23] Y. Tang et al. “Understanding Self-Predictive Learning for Reinforcement Learning”. In: *ICML*. 2023. URL: <https://proceedings.mlr.press/v202/tang23d/tang23d.pdf>.
- [Tan+24] H. Tang, K. Hu, J. P. Zhou, S. C. Zhong, W.-L. Zheng, X. Si, and K. Ellis. “Code Repair with LLMs gives an Exploration-Exploitation Tradeoff”. In: *NIPS*. Nov. 2024. URL: <https://openreview.net/pdf?id=o863gX6DxA>.
- [TCG21] Y. Tian, X. Chen, and S. Ganguli. “Understanding self-supervised Learning Dynamics without Contrastive Pairs”. In: *ICML*. Feb. 2021. URL: <http://arxiv.org/abs/2102.06810>.
- [Ten02] R. B. A. Tennenholtz. “R-max – A General Polynomial Time Algorithm for Near-Optimal Reinforcement Learning”. In: *JMLR* 3 (2002), pp. 213–231. URL: <http://www.ai.mit.edu/projects/jmlr/papers/volume3/brafman02a/source/brafman02a.pdf>.
- [TG96] G. Tesauro and G. R. Galperin. “On-line Policy Improvement using Monte-Carlo Search”. In: *NIPS*. 1996. URL: <https://arxiv.org/abs/2501.05407>.
- [Tha+22] S. Thakoor, M. Rowland, D. Borsa, W. Dabney, R. Munos, and A. Barreto. “Generalised Policy Improvement with Geometric Policy Composition”. en. In: *ICML*. PMLR, June 2022, pp. 21272–21307. URL: <https://proceedings.mlr.press/v162/thakoor22a.html>.
- [Tho33] W. R. Thompson. “On the Likelihood that One Unknown Probability Exceeds Another in View of the Evidence of Two Samples”. In: *Biometrika* 25.3/4 (1933), pp. 285–294.
- [Tim+20] F. Timbers, N. Bard, E. Lockhart, M. Lanctot, M. Schmid, N. Burch, J. Schrittwieser, T. Hubert, and M. Bowling. “Approximate exploitability: Learning a best response in large games”. In: *arXiv [cs.LG]* (Apr. 2020). URL: <http://arxiv.org/abs/2004.09677>.
- [Tio+24] S. Tiomkin, I. Nemenman, D. Polani, and N. Tishby. “Intrinsic motivation in dynamical control systems”. en. In: *PRX Life* 2.3 (Aug. 2024), p. 033009. URL: <http://dx.doi.org/10.1103/PRXLife.2.033009>.
- [TKE24] H. Tang, D. Y. Key, and K. Ellis. “WorldCoder, a Model-Based LLM Agent: Building World Models by Writing Code and Interacting with the Environment”. In: *NIPS*. Nov. 2024. URL: <https://openreview.net/pdf?id=QGJSXhVaL>.
- [TL05] E. Todorov and W. Li. “A Generalized Iterative LQG Method for Locally-optimal Feedback Control of Constrained Nonlinear Stochastic Systems”. In: *ACC*. 2005, pp. 300–306.
- [TLO23] J. Tarbouriech, T. Lattimore, and B. O’Donoghue. “Probabilistic Inference in Reinforcement Learning Done Right”. In: *NIPS*. Nov. 2023. URL: <https://openreview.net/pdf?id=9yQ2aaArDn>.
- [TM25] Y. Tang and R. Munos. “On a few pitfalls in KL divergence gradient estimation for RL”. In: *arXiv [cs.LG]* (June 2025). URL: <http://arxiv.org/abs/2506.09477>.
- [TMM19] C. Tessler, D. J. Mankowitz, and S. Mannor. “Reward Constrained Policy Optimization”. In: *ICLR*. 2019. URL: <https://openreview.net/pdf?id=SkfrvsA9FX>.
- [TO21] A. Touati and Y. Ollivier. “Learning one representation to optimize all rewards”. In: *NIPS*. Mar. 2021. URL: https://openreview.net/pdf?id=q_eWErV46er.
- [Tom+20] M. Tomar, L. Shani, Y. Efroni, and M. Ghavamzadeh. “Mirror descent policy optimization”. In: *arXiv [cs.LG]* (May 2020). URL: <http://arxiv.org/abs/2005.09814>.

- [Tom+22] T. Tomilin, T. Dai, M. Fang, and M. Pechenizkiy. “LevDoom: A benchmark for generalization on level difficulty in reinforcement learning”. In: *2022 IEEE Conference on Games (CoG)*. IEEE, Aug. 2022. URL: https://ieee-cog.org/2022/assets/papers/paper_30.pdf.
- [Tom+23] M. Tomar, U. A. Mishra, A. Zhang, and M. E. Taylor. “Learning Representations for Pixel-based Control: What Matters and Why?”. In: *Transactions on Machine Learning Research* (2023). URL: <https://openreview.net/pdf?id=wIXHG8LZ2w>.
- [Tom+24] M. Tomar, P. Hansen-Estruch, P. Bachman, A. Lamb, J. Langford, M. E. Taylor, and S. Levine. “Video Occupancy Models”. In: *arXiv [cs.CV]* (June 2024). URL: <http://arxiv.org/abs/2407.09533>.
- [Tou09] M. Toussaint. “Robot Rrajectory Optimization using Approximate Inference”. In: *ICML*. 2009, pp. 1049–1056.
- [Tou14] M. Toussaint. *Bandits, Global Optimization, Active Learning, and Bayesian RL – understanding the common ground*. Autonomous Learning Summer School. 2014. URL: <https://www.user.tu-berlin.de/mtoussai/teaching/14-BanditsOptimizationActiveLearningBayesianRL.pdf>.
- [TR97] J. Tsitsiklis and B. V. Roy. “An analysis of temporal-difference learning with function approximation”. In: *IEEE Trans. on Automatic Control* 42.5 (1997), pp. 674–690.
- [Tre+23] L. Treven, J. Hübötter, B. Sukhija, F. Dörfler, and A. Krause. “Efficient exploration in continuous-time model-based reinforcement learning”. In: *NIPS*. Oct. 2023. URL: <https://arxiv.org/abs/2310.19848>.
- [TRO23] A. Touati, J. Rapin, and Y. Ollivier. “Does Zero-Shot Reinforcement Learning Exist?”. In: *ICLR*. 2023. URL: https://openreview.net/forum?id=MYEap_OcQI.
- [TS06] M. Toussaint and A. Storkey. “Probabilistic inference for solving discrete and continuous state Markov Decision Processes”. In: *ICML*. 2006, pp. 945–952.
- [TS11] E. Talvitie and S. Singh. “Learning to make predictions in partially observable environments without a generative model”. en. In: *JAIR* 42 (Sept. 2011), pp. 353–392. URL: <https://jair.org/index.php/jair/article/view/10729>.
- [Tsc+20] A. Tschantz, B. Millidge, A. K. Seth, and C. L. Buckley. “Reinforcement learning through active inference”. In: *ICLR workshop on “Bridging AI and Cognitive Science”*. Feb. 2020.
- [Tsc+23] A. Tschantz, B. Millidge, A. K. Seth, and C. L. Buckley. “Hybrid predictive coding: Inferring, fast and slow”. en. In: *PLoS Comput. Biol.* 19.8 (Aug. 2023), e1011280. URL: <https://journals.plos.org/ploscompbiol/article/file?id=10.1371/journal.pcbi.1011280&type=printable>.
- [Tsi+17] P. A. Tsividis, T. Pouncy, J. L. Xu, J. B. Tenenbaum, and S. J. Gershman. “Human Learning in Atari”. en. In: *AAAI Spring Symposium Series*. 2017. URL: <https://www.aaai.org/ocs/index.php/SSS/SSS17/paper/viewPaper/15280>.
- [TVR97] J. N. Tsitsiklis and B. Van Roy. “An analysis of temporal-difference learning with function approximation”. en. In: *IEEE Trans. Automat. Contr.* 42.5 (May 1997), pp. 674–690. URL: <https://ieeexplore.ieee.org/abstract/document/580874>.
- [TWM25] Y. Tang, S. Wang, and R. Munos. “Learning to chain-of-thought with Jensen’s evidence lower bound”. In: *arXiv [cs.LG]* (Mar. 2025). URL: <http://arxiv.org/abs/2503.19618>.
- [Ueh+24] M. Uehara, Y. Zhao, T. Biancalani, and S. Levine. “Understanding reinforcement learning-based fine-tuning of diffusion models: A tutorial and review”. In: *arXiv [cs.LG]* (July 2024). URL: <http://arxiv.org/abs/2407.13734>.
- [Uel+25] T. Uelwer, J. Robine, S. S. Wagner, M. Höftmann, E. Upschulte, S. Konietzny, M. Behrendt, and S. Harmeling. “A survey on self-supervised methods for visual representation learning”. en. In: *Mach. Learn.* 114.4 (Apr. 2025), pp. 1–56. URL: <http://dx.doi.org/10.1007/s10994-024-06708-7>.

- [Unk24] Unknown. “Beyond The Rainbow: High Performance Deep Reinforcement Learning On A Desktop PC”. In: (Oct. 2024). URL: <https://openreview.net/pdf?id=0ydseYDKRi>.
- [VA+25] H. Van Assel, M. Ibrahim, T. Biancalani, A. Regev, and R. Balestriero. “Joint embedding vs reconstruction: Provable benefits of latent space prediction for Self Supervised Learning”. In: *NIPS*. Oct. 2025. URL: <https://arxiv.org/abs/2505.12477>.
- [Val00] H. Valpola. “Bayesian Ensemble Learning for Nonlinear Factor Analysis”. PhD thesis. Helsinki University of Technology, 2000. URL: https://users.ics.aalto.fi/harri/thesis/valpola_thesis.ps.gz.
- [van+18] H. van Hasselt, Y. Doron, F. Strub, M. Hessel, N. Sonnerat, and J. Modayil. *Deep Reinforcement Learning and the Deadly Triad*. arXiv:1812.02648. 2018.
- [Vas+21] S. Vaswani, O. Bachem, S. Totaro, R. Mueller, S. Garg, M. Geist, M. C. Machado, P. S. Castro, and N. L. Roux. “A general class of surrogate functions for stable and efficient reinforcement learning”. In: *arXiv [cs.LG]* (Aug. 2021). URL: <http://arxiv.org/abs/2108.05828>.
- [VBW15] S. S. Villar, J. Bowden, and J. Wason. “Multi-armed Bandit Models for the Optimal Design of Clinical Trials: Benefits and Challenges”. en. In: *Stat. Sci.* 30.2 (2015), pp. 199–215. URL: <http://dx.doi.org/10.1214/14-STS504>.
- [Vee+19] V. Veeriah, M. Hessel, Z. Xu, J. Rajendran, R. L. Lewis, J. Oh, H. P. van Hasselt, D. Silver, and S. Singh. “Discovery of Useful Questions as Auxiliary Tasks”. In: *NIPS*. Vol. 32. 2019. URL: https://proceedings.neurips.cc/paper_files/paper/2019/file/10ff0b5e85e5b85cc3095d431d8c08b4-Paper.pdf.
- [Ven+24] D. Venuto, S. N. Islam, M. Klissarov, D. Precup, S. Yang, and A. Anand. “Code as reward: Empowering reinforcement learning with VLMs”. In: *ICML*. Feb. 2024. URL: <https://openreview.net/forum?id=6P88DMUDvH>.
- [Vez+17] A. S. Vezhnevets, S. Osindero, T. Schaul, N. Heess, M. Jaderberg, D. Silver, and K. Kavukcuoglu. “FeUdal Networks for Hierarchical Reinforcement Learning”. en. In: *ICML*. PMLR, July 2017, pp. 3540–3549. URL: <https://proceedings.mlr.press/v70/vezhnevets17a.html>.
- [Vil+22] A. R. Villaflor, Z. Huang, S. Pande, J. M. Dolan, and J. Schneider. “Addressing Optimism Bias in Sequence Modeling for Reinforcement Learning”. en. In: *ICML*. PMLR, June 2022, pp. 22270–22283. URL: <https://proceedings.mlr.press/v162/villaflor22a.html>.
- [Vin+19] O. Vinyals et al. “Grandmaster level in StarCraft II using multi-agent reinforcement learning”. en. In: *Nature* 575.7782 (Nov. 2019), pp. 350–354. URL: <http://dx.doi.org/10.1038/s41586-019-1724-z>.
- [VPG20] N. Vieillard, O. Pietquin, and M. Geist. “Munchausen Reinforcement Learning”. In: *NIPS*. Vol. 33. 2020, pp. 4235–4246. URL: https://proceedings.neurips.cc/paper_files/paper/2020/file/2c6a0bae0f071cbbf0bb3d5b11d90a82-Paper.pdf.
- [Vri+25] B. de Vries et al. “Expected Free Energy-based planning as variational inference”. In: *arXiv [stat.ML]* (Apr. 2025). URL: <http://arxiv.org/abs/2504.14898>.
- [WA25] R. Wang and P. Ammanabrolu. “A practitioner’s guide to multi-turn agentic reinforcement learning”. In: *arXiv [cs.LG]* (Oct. 2025). URL: <http://arxiv.org/abs/2510.01132>.
- [Wag+19] N. Wagener, C.-A. Cheng, J. Sacks, and B. Boots. “An online learning approach to model predictive control”. In: *Robotics: Science and Systems*. Feb. 2019. URL: <https://arxiv.org/abs/1902.08967>.
- [Wan+16] Z. Wang, T. Schaul, M. Hessel, H. van Hasselt, M. Lanctot, and N. de Freitas. “Dueling Network Architectures for Deep Reinforcement Learning”. In: *ICML*. 2016. URL: <http://proceedings.mlr.press/v48/wangf16.pdf>.
- [Wan+19] T. Wang, X. Bao, I. Clavera, J. Hoang, Y. Wen, E. Langlois, S. Zhang, G. Zhang, P. Abbeel, and J. Ba. “Benchmarking Model-Based Reinforcement Learning”. In: *arXiv [cs.LG]* (July 2019). URL: <http://arxiv.org/abs/1907.02057>.

- [Wan+22] T. Wang, S. S. Du, A. Torralba, P. Isola, A. Zhang, and Y. Tian. “Denoised MDPs: Learning World Models Better Than the World Itself”. In: *ICML*. June 2022. URL: <http://arxiv.org/abs/2206.15477>.
- [Wan+23] T. T. Wang et al. “Adversarial Policies Beat Superhuman Go AIs”. In: *ICML*. 2023. URL: <https://arxiv.org/abs/2211.00241>.
- [Wan+24a] G. Wang, Y. Xie, Y. Jiang, A. Mandlekar, C. Xiao, Y. Zhu, L. Fan, and A. Anandkumar. “Voyager: An Open-Ended Embodied Agent with Large Language Models”. In: *TMLR* (2024). URL: <https://openreview.net/forum?id=ehfRiFOR3a>.
- [Wan+24b] S. Wang, S. Liu, W. Ye, J. You, and Y. Gao. “EfficientZero V2: Mastering discrete and continuous control with limited data”. In: *arXiv [cs.LG]* (Mar. 2024). URL: <http://arxiv.org/abs/2403.00564>.
- [Wan+24c] X. Wang, Y. Chen, L. Yuan, Y. Zhang, Y. Li, H. Peng, and H. Ji. “Executable Code Actions Elicit Better LLM Agents”. In: *ICML*. June 2024. URL: <https://openreview.net/forum?id=jJ9BoXAfFa>.
- [Wan+25a] Z. Wan et al. “ReMA: Learning to meta-think for LLMs with multi-Agent Reinforcement Learning”. In: *arXiv [cs.AI]* (Mar. 2025). URL: <http://arxiv.org/abs/2503.09501>.
- [Wan+25b] W. Wang, I. Dusparic, Y. Shi, K. Zhang, and V. Cahill. “Drama: Mamba-enabled model-based reinforcement learning is sample and parameter efficient”. In: *arXiv [cs.LG]* (2025). URL: <http://arxiv.org/abs/2410.08893>.
- [Wan+25c] Z. Wang et al. “RAGEN: Understanding self-evolution in LLM agents via multi-turn reinforcement learning”. In: *arXiv [cs.LG]* (Apr. 2025). URL: <http://arxiv.org/abs/2504.20073>.
- [WAT17] G. Williams, A. Aldrich, and E. A. Theodorou. “Model Predictive Path Integral Control: From Theory to Parallel Computation”. In: *J. Guid. Control Dyn.* 40.2 (Feb. 2017), pp. 344–357. URL: <https://doi.org/10.2514/1.G001921>.
- [Wat+21] J. Watson, H. Abdulsamad, R. Findeisen, and J. Peters. “Stochastic Control through Approximate Bayesian Input Inference”. In: *arxiv* (2021). URL: <http://arxiv.org/abs/2105.07693>.
- [Wau+15] K. Waugh, D. Morrill, J. Bagnell, and M. Bowling. “Solving games with functional regret estimation”. en. In: *AAAI* 29.1 (Feb. 2015). URL: <https://ojs.aaai.org/index.php/AAAI/article/view/9445>.
- [WCM24] C. Wang, Y. Chen, and K. Murphy. “Model-based Policy Optimization under Approximate Bayesian Inference”. en. In: *AISTATS*. PMLR, Apr. 2024, pp. 3250–3258. URL: <https://proceedings.mlr.press/v238/wang24g.html>.
- [WD92] C. Watkins and P. Dayan. “Q-learning”. In: *Machine Learning* 8.3 (1992), pp. 279–292.
- [Web+17] T. Weber et al. “Imagination-Augmented Agents for Deep Reinforcement Learning”. In: *NIPS*. 2017. URL: <http://arxiv.org/abs/1707.06203>.
- [Wei+22] J. Wei, X. Wang, D. Schuurmans, M. Bosma, E. Chi, Q. Le, and D. Zhou. “Chain of Thought Prompting Elicits Reasoning in Large Language Models”. In: *arXiv [cs.CL]* (Jan. 2022). URL: <http://arxiv.org/abs/2201.11903>.
- [Wei+24] R. Wei, N. Lambert, A. McDonald, A. Garcia, and R. Calandra. “A unified view on solving objective mismatch in model-based Reinforcement Learning”. In: *Trans. on Machine Learning Research* (2024). URL: <https://openreview.net/forum?id=tQVZgvXhZb>.
- [Wen18a] L. Weng. “A (Long) Peek into Reinforcement Learning”. In: *lilianweng.github.io* (2018). URL: <https://lilianweng.github.io/posts/2018-02-19-rl-overview/>.
- [Wen18b] L. Weng. “Policy Gradient Algorithms”. In: *lilianweng.github.io* (2018). URL: <https://lilianweng.github.io/posts/2018-04-08-policy-gradient/>.
- [Wen+25] Y. Wen, Y. Ye, Y. Zhang, D. Yang, and H. Zhu. “Real-time reasoning agents in evolving environments”. In: *arXiv [cs.AI]* (Nov. 2025). URL: <http://arxiv.org/abs/2511.04898>.

- [Whi15] M White Adam. “Developing a predictive approach to knowledge”. PhD thesis. University of Alberta Libraries, 2015. URL: <http://dx.doi.org/10.7939/R3FF3M75H>.
- [WHT19] Y. Wang, H. He, and X. Tan. “Truly Proximal Policy Optimization”. In: *UAI*. 2019. URL: <http://auai.org/uai2019/proceedings/papers/21.pdf>.
- [WHZ23] Z. Wang, J. J. Hunt, and M. Zhou. “Diffusion Policies as an Expressive Policy Class for Offline Reinforcement Learning”. In: *ICLR*. 2023. URL: <https://openreview.net/pdf?id=AHvFDPi-FA>.
- [Wie03] E Wiewiora. “Potential-Based Shaping and Q-Value Initialization are Equivalent”. In: *JAIR*. 2003. URL: <https://jair.org/index.php/jair/article/view/10338>.
- [Wil+17] G. Williams, N. Wagener, B. Goldfain, P. Drews, J. M. Rehg, B. Boots, and E. A. Theodorou. “Information theoretic MPC for model-based reinforcement learning”. In: *ICRA*. IEEE, May 2017, pp. 1714–1721. URL: <https://ieeexplore.ieee.org/document/7989202>.
- [Wil92] R. J. Williams. “Simple statistical gradient-following algorithms for connectionist reinforcement learning”. In: *MLJ* 8.3-4 (1992), pp. 229–256.
- [WIP20] J. Watson, A. Imohiosen, and J. Peters. “Active Inference or Control as Inference? A Unifying View”. In: *International Workshop on Active Inference*. 2020. URL: <http://arxiv.org/abs/2010.00262>.
- [Wit+20] C. S. de Witt, T. Gupta, D. Makoviichuk, V. Makoviychuk, P. H. S. Torr, M. Sun, and S. Whiteson. “Is independent learning all you need in the StarCraft multi-agent challenge?” In: *arXiv [cs.AI]* (Nov. 2020). URL: <http://arxiv.org/abs/2011.09533>.
- [WNS21] Y. Wan, A. Naik, and R. S. Sutton. “Learning and planning in average-reward Markov decision processes”. In: *ICML*. 2021. URL: <https://arxiv.org/abs/2006.16318>.
- [Won+22] A. Wong, T. Bäck, A. V. Kononova, and A. Plaat. “Deep multiagent reinforcement learning: challenges and directions”. en. In: *Artif. Intell. Rev.* 56.6 (Oct. 2022), pp. 5023–5056. URL: <https://link.springer.com/article/10.1007/s10462-022-10299-x>.
- [Won+23] L. Wong, J. Mao, P. Sharma, Z. S. Siegel, J. Feng, N. Korneev, J. B. Tenenbaum, and J. Andreas. “Learning adaptive planning representations with natural language guidance”. In: *arXiv [cs.AI]* (Dec. 2023). URL: <http://arxiv.org/abs/2312.08566>.
- [WSF25] K. A. Wang, J. Shi, and E. B. Fox. “Test-time regression: a unifying framework for designing sequence models with associative memory”. In: *arXiv [cs.LG]* (Jan. 2025). URL: <http://arxiv.org/abs/2501.12352>.
- [Wu+17] Y. Wu, E. Mansimov, S. Liao, R. Grosse, and J. Ba. “Scalable trust-region method for deep reinforcement learning using Kronecker-factored approximation”. In: *NIPS*. 2017. URL: <https://arxiv.org/abs/1708.05144>.
- [Wu+21] Y. Wu, S. Zhai, N. Srivastava, J. Susskind, J. Zhang, R. Salakhutdinov, and H. Goh. “Uncertainty Weighted Actor-critic for offline Reinforcement Learning”. In: *ICML*. May 2021. URL: <https://arxiv.org/abs/2105.08140>.
- [Wu+22] P. Wu, A. Escontrela, D. Hafner, K. Goldberg, and P. Abbeel. “DayDreamer: World Models for Physical Robot Learning”. In: (June 2022). arXiv: [2206.14176 \[cs.R0\]](https://arxiv.org/abs/2206.14176). URL: [http://arxiv.org/abs/2206.14176](https://arxiv.org/abs/2206.14176).
- [Wu+23] G. Wu, W. Fang, J. Wang, P. Ge, J. Cao, Y. Ping, and P. Gou. “Dyna-PPO reinforcement learning with Gaussian process for the continuous action decision-making in autonomous driving”. en. In: *Appl. Intell.* 53.13 (July 2023), pp. 16893–16907. URL: <https://link.springer.com/article/10.1007/s10489-022-04354-x>.
- [Wur+22] P. R. Wurman et al. “Outracing champion Gran Turismo drivers with deep reinforcement learning”. en. In: *Nature* 602.7896 (Feb. 2022), pp. 223–228. URL: https://www.researchgate.net/publication/358484368_Outracing_champion_Gran_Turismo_drivers_with_deep_reinforcement_learning.

- [Xi+25] Z. Xi et al. “BAPO: Stabilizing off-policy reinforcement learning for LLMs via BALanced Policy Optimization with Adaptive Clipping”. In: *arXiv [cs.LG]* (Oct. 2025). URL: <http://arxiv.org/abs/2510.18927>.
- [Xie+25] Z. Xie, Q. Zhang, F. Yang, M. Hutter, and R. Xu. “Simple Policy Optimization”. In: *ICML*. June 2025. URL: <https://openreview.net/forum?id=SG8Yx1FyeU>.
- [XL24] S. Xie and Z. Li. “Implicit bias of AdamW: ℓ_∞ norm constrained optimization”. In: *arXiv [cs.LG]* (Apr. 2024). URL: <http://arxiv.org/abs/2404.04454>.
- [Xu+17] C. Xu, T. Qin, G. Wang, and T.-Y. Liu. “Reinforcement learning for learning rate control”. In: *arXiv [cs.LG]* (May 2017). URL: <http://arxiv.org/abs/1705.11159>.
- [Xu+22] T. Xu, Z. Yang, Z. Wang, and Y. Liang. “A Unifying Framework of Off-Policy General Value Function Evaluation”. In: *NIPS*. Oct. 2022. URL: https://openreview.net/pdf?id=LdKdbHw3A_6.
- [Xu+25] F. Xu et al. “Towards large reasoning models: A survey of reinforced reasoning with Large Language Models”. In: *arXiv [cs.AI]* (Jan. 2025). URL: <http://arxiv.org/abs/2501.09686>.
- [XZ25] T. Xiao and J. Zhu. “Foundations of large language models”. In: *arXiv [cs.CL]* (Jan. 2025). URL: <http://arxiv.org/abs/2501.09223>.
- [Yab+25] Z. Yabo, Z. Yihan, L. Qingyun, H. Zhen, H. Kavin, and Z. Wangmeng. “Tool-R1: Sample-efficient reinforcement learning for agentic tool use”. In: *arXiv [cs.LG]* (Sept. 2025). URL: <http://arxiv.org/abs/2509.12867>.
- [Yan+23] M. Yang, D. Schuurmans, P. Abbeel, and O. Nachum. “Dichotomy of control: Separating what you can control from what you cannot”. In: *ICLR*. Vol. abs/2210.13435. 2023. URL: <https://github.com/google-research/google-research/tree/>.
- [Yan+24] S. Yang, Y. Du, S. K. S. Ghasemipour, J. Tompson, L. P. Kaelbling, D. Schuurmans, and P. Abbeel. “Learning Interactive Real-World Simulators”. In: *ICLR*. 2024. URL: <https://openreview.net/pdf?id=sFyTZEqmUY>.
- [Yan25] W. Yan. “Don’t build multi-agents”. In: (2025). URL: <https://cognition.ai/blog/dont-build-multi-agents>.
- [Yan+25] N. Yang, J. Gao, F. Gao, Y. Wu, and C. Yu. “Fine-tuning Diffusion Policies with backpropagation through diffusion timesteps”. In: *arXiv [cs.LG]* (Sept. 2025). URL: <http://arxiv.org/abs/2505.10482>.
- [Yao+22] S. Yao, J. Zhao, D. Yu, N. Du, I. Shafran, K. R. Narasimhan, and Y. Cao. “ReAct: Synergizing Reasoning and Acting in Language Models”. In: *ICLR*. Sept. 2022. URL: https://openreview.net/pdf?id=WE_vluYUL-X.
- [Ye+21] W. Ye, S. Liu, T. Kurutach, P. Abbeel, and Y. Gao. “Mastering Atari games with limited data”. In: *NIPS*. Oct. 2021.
- [Yiu+25] E. Yiu, K. Allen, S. Ginosar, and A. Gopnik. “Empowerment Gain and Causal Model Construction: Children and adults are sensitive to controllability and variability in their causal interventions”. en. In: *PsyArXiv* (June 2025). URL: https://osf.io/preprints/psyarxiv/ept4n_v1.
- [YKSR23] T. Yamagata, A. Khalil, and R. Santos-Rodriguez. “Q-learning Decision Transformer: Leveraging Dynamic Programming for conditional sequence modelling in offline RL”. In: *ICML*. 2023. URL: <https://arxiv.org/abs/2209.03993>.
- [Yu17] H. Yu. “On convergence of some gradient-based temporal-differences algorithms for off-policy learning”. In: *arXiv [cs.LG]* (Dec. 2017). URL: <http://arxiv.org/abs/1712.09652>.
- [Yu+20] T. Yu, G. Thomas, L. Yu, S. Ermon, J. Y. Zou, S. Levine, C. Finn, and T. Ma. “MOPO: Model-based Offline Policy Optimization”. In: *NIPS*. Vol. 33. 2020, pp. 14129–14142. URL: https://proceedings.neurips.cc/paper_files/paper/2020/hash/a322852ce0df73e204b7e67cbbef0d0a-Abstract.html.

- [Yu+22] C. Yu, A. Velu, E. Vinitisky, J. Gao, Y. Wang, A. Bayen, and Y. Wu. “The surprising effectiveness of PPO in cooperative, multi-agent games”. In: *NeurIPS 2022 Datasets and Benchmarks*. 2022. URL: <https://arxiv.org/abs/2103.01955>.
- [Yu+23] C. Yu, N. Burgess, M. Sahani, and S. Gershman. “Successor-Predecessor Intrinsic Exploration”. In: *NIPS*. Vol. abs/2305.15277. Curran Associates, Inc., May 2023, pp. 73021–73038. URL: https://proceedings.neurips.cc/paper_files/paper/2023/hash/e6f2b968c4ee8ba260cd7077e39590dd-Abstract-Conference.html.
- [Yu+25] Q. Yu et al. “DAPO: An open-source LLM reinforcement learning system at scale”. In: *arXiv [cs.LG]* (Mar. 2025). URL: <http://arxiv.org/abs/2503.14476>.
- [Yua22] M. Yuan. “Intrinsically-motivated reinforcement learning: A brief introduction”. In: *arXiv [cs.LG]* (Mar. 2022). URL: <http://arxiv.org/abs/2203.02298>.
- [Yua+25] M. Yuan, R. C. Castanyer, B. Li, X. Jin, G. Berseth, and W. Zeng. “RLeXplore: Accelerating research in intrinsically-motivated reinforcement learning”. In: *Trans. on Machine Learning Research* (May 2025). URL: <http://arxiv.org/abs/2405.19548>.
- [Yue+25] Y. Yue, Z. Chen, R. Lu, A. Zhao, Z. Wang, Y. Yue, S. Song, and G. Huang. “Does reinforcement Learning really incentivize reasoning capacity in LLMs beyond the base model?” In: *arXiv [cs.AI]* (Apr. 2025). URL: <http://arxiv.org/abs/2504.13837>.
- [YW20] Y. Yang and J. Wang. “An overview of multi-agent reinforcement learning from game theoretical perspective”. In: *arXiv [cs.MA]* (Nov. 2020). URL: <http://arxiv.org/abs/2011.00583>.
- [YWW25] M. Yin, M. Wang, and Y.-X. Wang. “On the statistical complexity for offline and low-adaptive reinforcement learning with structures”. In: *arXiv [cs.LG]* (Jan. 2025). URL: <http://arxiv.org/abs/2501.02089>.
- [YZ22] Y. Yang and P. Zhai. “Click-through rate prediction in online advertising: A literature review”. In: *Inf. Process. Manag.* 59.2 (2022), p. 102853. URL: <https://www.sciencedirect.com/science/article/pii/S0306457321003241>.
- [ZABD10] B. D. Ziebart, J. Andrew Bagnell, and A. K. Dey. “Modeling Interaction via the Principle of Maximum Causal Entropy”. In: *ICML*. 2010. URL: <https://www.cs.uic.edu/pub/Ziebart/Publications/maximum-causal-entropy.pdf>.
- [Zbo+21] J. Zbontar, L. Jing, I. Misra, Y. LeCun, and S. Deny. “Barlow Twins: Self-Supervised Learning via Redundancy Reduction”. In: Mar. 2021. URL: <https://arxiv.org/abs/2103.03230>.
- [Zel+24] E. Zelikman, G. Harik, Y. Shao, V. Jayasiri, N. Haber, and N. D. Goodman. “Quiet-STaR: Language Models Can Teach Themselves to Think Before Speaking”. In: *arXiv [cs.CL]* (Mar. 2024). URL: <http://arxiv.org/abs/2403.09629>.
- [Zen+25] S. Zeng, L. Viano, C. Li, J. Li, V. Cevher, M. Wulfmeier, S. Ermon, A. Garcia, and M. Hong. “Aligning large language models with human feedback: Mathematical foundations and algorithm design”. In: *Techrxiv* (May 2025). URL: <https://www.authorea.com/doi/full/10.36227/techrxiv.174784525.51683948?commit=a5839728003f40ea2bcac936bca54b11d7b1edff>.
- [Zha+18] R. Zhang, C. Chen, C. Li, and L. Carin. “Policy Optimization as Wasserstein Gradient Flows”. In: *ICML*. July 2018, pp. 5737–5746. URL: <https://proceedings.mlr.press/v80/zhang18a.html>.
- [Zha+19] S. Zhang, B. Liu, H. Yao, and S. Whiteson. “Provably convergent two-timescale off-policy actor-critic with function approximation”. In: *ICML* 119 (Nov. 2019). Ed. by H. D. Iii and A. Singh, pp. 11204–11213. URL: <https://proceedings.mlr.press/v119/zhang20s/zhang20s.pdf>.
- [Zha+21] A. Zhang, R. T. McAllister, R. Calandra, Y. Gal, and S. Levine. “Learning Invariant Representations for Reinforcement Learning without Reconstruction”. In: *ICLR*. 2021. URL: <https://openreview.net/pdf?id=-2FCwDKRREu>.

- [Zha+23a] J. Zhang, J. T. Springenberg, A. Byravan, L. Hasenclever, A. Abdolmaleki, D. Rao, N. Heess, and M. Riedmiller. “Leveraging Jumpy Models for Planning and Fast Learning in Robotic Domains”. In: *arXiv [cs.RO]* (Feb. 2023). URL: <http://arxiv.org/abs/2302.12617>.
- [Zha+23b] W. Zhang, G. Wang, J. Sun, Y. Yuan, and G. Huang. “STORM: Efficient Stochastic Transformer based world models for reinforcement learning”. In: *arXiv [cs.LG]* (Oct. 2023). URL: <http://arxiv.org/abs/2310.09615>.
- [Zha+24a] Y. Zhai et al. “Fine-tuning large vision-language models as decision-making agents via reinforcement learning”. In: *NIPS* (May 2024). URL: <http://arxiv.org/abs/2405.10292>.
- [Zha+24b] J. Zhang, J. Lehman, K. Stanley, and J. Clune. “OMNI: Open-endedness via Models of human Notions of Interestingness”. In: *ICLR*. 2024. URL: <https://openreview.net/forum?id=AgM3MzT99c>.
- [Zha+24c] Z. Zhang, X. Bo, C. Ma, R. Li, X. Chen, Q. Dai, J. Zhu, Z. Dong, and J.-R. Wen. “A survey on the memory mechanism of large language model based agents”. In: *arXiv [cs.AI]* (Apr. 2024). URL: <http://arxiv.org/abs/2404.13501>.
- [Zha+24d] S. Zhao, R. Brekelmans, A. Makhzani, and R. Grosse. “Probabilistic inference in language models via twisted Sequential Monte Carlo”. In: *ICML*. Apr. 2024. URL: <https://arxiv.org/abs/2404.17546>.
- [Zha+24e] S. Zhao, R. Brekelmans, A. Makhzani, and R. B. Grosse. “Probabilistic Inference in Language Models via Twisted Sequential Monte Carlo”. In: *ICML*. June 2024. URL: <https://openreview.net/pdf?id=frA0NNBS1n>.
- [Zha+25a] G. Zhang et al. “The landscape of agentic reinforcement learning for LLMs: A survey”. In: *arXiv [cs.AI]* (Sept. 2025). URL: <http://arxiv.org/abs/2509.02547>.
- [Zha+25b] K. Zhang et al. “A survey of Reinforcement Learning for large reasoning models”. In: *arXiv [cs.CL]* (Sept. 2025). URL: <http://arxiv.org/abs/2509.08827>.
- [Zha+25c] K. Zhang et al. “A survey of Reinforcement Learning for large reasoning models”. In: *arXiv [cs.CL]* (Sept. 2025). URL: <http://arxiv.org/abs/2509.08827>.
- [Zha+25d] Q. Zhang et al. “Agentic Context Engineering: Evolving contexts for self-improving language models”. In: *arXiv [cs.LG]* (Oct. 2025). URL: <http://arxiv.org/abs/2510.04618>.
- [Zha+25e] Y. Zhang, Y. Liu, H. Yuan, Y. Yuan, Q. Gu, and A. C.-C. Yao. “On the design of KL-Regularized Policy Gradient algorithms for LLM reasoning”. In: *arXiv [cs.LG]* (Sept. 2025). URL: <http://arxiv.org/abs/2505.17508>.
- [Zha+25f] A. Zhao et al. “Absolute Zero: Reinforced self-play reasoning with zero data”. In: *arXiv [cs.LG]* (May 2025). URL: <http://arxiv.org/abs/2505.03335>.
- [Zha+25g] W. Zhao, X. Wang, C. Ma, L. Kong, Z. Yang, M. Tuo, X. Shi, Y. Zhai, and X. Cai. “MUA-RL: Multi-turn User-interacting Agent Reinforcement Learning for agentic tool use”. In: *arXiv [cs.AI]* (Aug. 2025). URL: <http://arxiv.org/abs/2508.18669>.
- [Zhe+22a] L. Zheng, T. Fiez, Z. Alumbaugh, B. Chasnov, and L. J. Ratliff. “Stackelberg actor-critic: Game-theoretic reinforcement learning algorithms”. en. In: *AAAI* 36.8 (June 2022), pp. 9217–9224. URL: <https://ojs.aaai.org/index.php/AAAI/article/view/20908>.
- [Zhe+22b] R. Zheng, X. Wang, H. Xu, and F. Huang. “Is Model Ensemble Necessary? Model-based RL via a Single Model with Lipschitz Regularized Value Function”. In: *ICLR*. Sept. 2022. URL: <https://openreview.net/pdf?id=hNyJBk3CwR>.
- [Zhe+23] C. Zheng, P. Ke, Z. Zhang, and M. Huang. “Click: Controllable text generation with sequence likelihood contrastive learning”. In: *ACL*. June 2023. URL: <https://arxiv.org/abs/2306.03350>.

- [Zhe+24a] Q. Zheng, M. Henaff, A. Zhang, A. Grover, and B. Amos. “Online intrinsic rewards for decision making agents from large language model feedback”. In: *arXiv [cs.LG]* (Oct. 2024). URL: <http://arxiv.org/abs/2410.23022>.
- [Zhe+24b] R. Zheng et al. “Toward optimal LLM alignments using two-player games”. In: *arXiv [cs.CL]* (June 2024). URL: <http://arxiv.org/abs/2406.10977>.
- [Zhe+25a] C. Zheng, J. Tuyls, J. Peng, and B. Eysenbach. “Can a MISL Fly? Analysis and Ingredients for Mutual Information Skill Learning”. In: *The Thirteenth International Conference on Learning Representations*. 2025. URL: <https://openreview.net/forum?id=xoIeVdF07U>.
- [Zhe+25b] C. Zheng et al. “Group Sequence Policy Optimization”. In: *arXiv [cs.LG]* (July 2025). URL: <http://arxiv.org/abs/2507.18071>.
- [Zhe+25c] J. Zheng, C. Shi, X. Cai, Q. Li, D. Zhang, C. Li, D. Yu, and Q. Ma. “Lifelong learning of large language model based agents: A roadmap”. In: *arXiv [cs.AI]* (Jan. 2025). URL: <http://arxiv.org/abs/2501.07278>.
- [Zho+24a] G. Zhou, H. Pan, Y. LeCun, and L. Pinto. “DINO-WM: World models on pre-trained visual features enable zero-shot planning”. In: *arXiv [cs.RO]* (Nov. 2024). URL: <http://arxiv.org/abs/2411.04983>.
- [Zho+24b] G. Zhou, S. Swaminathan, R. V. Raju, J. S. Guntupalli, W. Lehrach, J. Ortiz, A. Dedieu, M. Lázaro-Gredilla, and K. Murphy. “Diffusion Model Predictive Control”. In: *arXiv [cs.LG]* (Oct. 2024). URL: <http://arxiv.org/abs/2410.05364>.
- [Zho+25a] X. Zhou, Z. Liu, H. Wang, C. Du, M. Lin, C. Li, L. Wang, and T. Pang. “Variational reasoning for language models”. In: *arXiv [cs.CL]* (Sept. 2025). URL: <http://arxiv.org/abs/2509.22637>.
- [Zho+25b] Y. Zhou, S. Levine, J. Weston, X. Li, and S. Sukhbaatar. “Self-challenging language model agents”. In: *arXiv [cs.AI]* (June 2025). URL: <http://arxiv.org/abs/2506.01716>.
- [Zho+25c] Y. Zhou, S. Wang, X. Dong, X. Jin, Y. Chen, Y. Min, K. Yang, X. Ren, D. Liu, and L. Zhang. “Reasoning like an economist: Post-training on economic problems induces strategic generalization in LLMs”. In: *arXiv [cs.AI]* (May 2025). URL: <http://arxiv.org/abs/2506.00577>.
- [ZHR24] H. Zhu, B. Huang, and S. Russell. “On representation complexity of model-based and model-free reinforcement learning”. In: *ICLR*. 2024.
- [Zhu+24] Z. Zhu et al. “Is Sora a world simulator? A comprehensive survey on General world models and beyond”. In: *arXiv [cs.CV]* (May 2024). URL: <http://arxiv.org/abs/2405.03520>.
- [Zie+08] B. D. Ziebart, A. L. Maas, J. A. Bagnell, and A. K. Dey. “Maximum Entropy Inverse Reinforcement Learning”. In: *AAAI*. 2008, pp. 1433–1438.
- [Zin+07] M. Zinkevich, M. Johanson, M. Bowling, and C. Piccione. “Regret minimization in games with incomplete information”. en. In: *NIPS*. Dec. 2007, pp. 1729–1736. URL: <https://dl.acm.org/doi/10.5555/2981562.2981779>.
- [Zin+21] L. Zintgraf, S. Schulze, C. Lu, L. Feng, M. Igl, K. Shiarlis, Y. Gal, K. Hofmann, and S. Whiteson. “VariBAD: Variational Bayes-Adaptive Deep RL via meta-learning”. In: *J. Mach. Learn. Res.* 22.289 (2021), 289:1–289:39. URL: <https://www.jmlr.org/papers/volume22/21-0657/21-0657.pdf>.
- [Zit+23] B. Zitkovich et al. “RT-2: Vision-Language-Action Models Transfer Web Knowledge to Robotic Control”. en. In: *Conference on Robot Learning*. PMLR, Dec. 2023, pp. 2165–2183. URL: <https://proceedings.mlr.press/v229/zitkovich23a.html>.
- [ZS22] N. Zucchet and J. Sacramento. “Beyond backpropagation: Bilevel optimization through implicit differentiation and equilibrium propagation”. en. In: *Neural Comput.* 34.12 (Nov. 2022), pp. 2309–2346. URL: https://direct.mit.edu/neco/article-pdf/34/12/2309/2057431/neco_a_01547.pdf.

- [ZSE24] C. Zheng, R. Salakhutdinov, and B. Eysenbach. “Contrastive Difference Predictive Coding”. In: *The Twelfth International Conference on Learning Representations*. 2024. URL: <https://openreview.net/pdf?id=0akLDTFR9x>.
- [ZW19] S. Zhang and S. Whiteson. “DAC: The Double Actor-Critic Architecture for Learning Options”. In: *NIPS* 32 (2019). URL: https://proceedings.neurips.cc/paper_files/paper/2019/file/4f284803bd0966cc24fa8683a34afc6e-Paper.pdf.
- [ZWG22] E. Zelikman, Y. Wu, and N. D. Goodman. “STaR: Bootstrapping Reasoning With Reasoning”. In: *NIPS*. Mar. 2022. URL: <https://arxiv.org/abs/2203.14465>.