
Co-Learning Empirical Games and World Models

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Game-based decision-making involves reasoning over both world dynamics and
2 strategic interactions among the agents. Typically, empirical models capturing these
3 respective aspects are learned and used separately. We investigate the potential gain
4 from co-learning these elements: a world model for dynamics and an empirical
5 game for strategic interactions. Empirical games drive world models toward a
6 broader consideration of possible game dynamics induced by a diversity of strategy
7 profiles. Conversely, world models guide empirical games to efficiently discover
8 new strategies through planning. We demonstrate these benefits first independently,
9 then in combination as realized by a new algorithm, Dyna-PSRO, that co-learns
10 an empirical game and a world model. When compared to PSRO—a baseline
11 empirical-game building algorithm, Dyna-PSRO is found to compute lower regret
12 solutions on partially observable general-sum games. In our experiments, Dyna-
13 PSRO also requires substantially fewer experiences than PSRO, a key algorithmic
14 advantage for settings where collecting player-game interaction data is a cost-
15 limiting factor.

16 1 Introduction

17 Even seemingly simple games can actually embody a level of complexity rendering them intractable
18 to direct reasoning. This complexity stems from the interplay of two sources: dynamics of the
19 game environment, and strategic interactions among the game’s players. As an alternative to direct
20 reasoning, models have been developed to facilitate reasoning over these distinct aspects of the game.
21 *Empirical games* capture strategic interactions in the form of payoff estimates for joint policies [80].
22 *World models* represent a game’s transition dynamics and reward signal directly [69, 19]. Whereas
23 each of these forms of model have been found useful for game reasoning, typical use in prior work
24 has focused on one or the other, learned and employed in isolation from its natural counterpart.

25 Co-learning both models presents an opportunity to leverage their complementary strengths as a
26 means to improve each other. World models predict successor states and rewards given a game’s
27 current state and action(s). However, their performance depends on coverage of their training data,
28 which is limited by the range of strategies considered during learning. Empirical games can inform
29 training of world models by suggesting a diverse set of salient strategies, based on game-theoretic
30 reasoning [80]. These strategies can expose the world model to a broader range of relevant dynamics.
31 Moreover, as empirical games are estimated through simulation of strategy profiles, this same
32 simulation data can be reused as training data for the world model.

33 Strategic diversity through empirical games, however, comes at a cost. In the popular framework
34 of Policy-Space Response Oracles (PSRO) [38], empirical normal-form game models are built
35 iteratively, at each step expanding a restricted strategy set by computing best-response policies to
36 the current game’s solution. As computing an exact best-response is generally intractable, PSRO
37 uses Deep Reinforcement Learning (DRL) to compute approximate response policies. However,
38 each application of DRL can be considerably resource-intensive, necessitating the generation of

39 a vast amount of gameplays for learning. Whether gameplays, or experiences, are generated via
 40 simulation [48] or from real-world interactions [24], their collection poses a major limiting factor in
 41 DRL and by extension PSRO. World models present one avenue to reduce this cost by transferring
 42 previously learned game dynamics across response computations.

43 We investigate the mutual benefits of co-learning
 44 a world model and an empirical game by first
 45 verifying the potential contributions of each
 46 component independently. We then show how
 47 to realize the combined effects in a new algo-
 48 rithm, *Dyna-PSRO*, that co-learns a world model
 49 and an empirical game (illustrated in Figure 1).
 50 *Dyna-PSRO* extends PSRO to learn a world
 51 model concurrently with empirical game expan-
 52 sion, and applies this world model to reduce the
 53 computational cost of computing new policies.
 54 This is implemented by a Dyna-based reinforce-
 55 ment learner [67, 68] that integrates planning,
 56 acting, and learning in parallel. *Dyna-PSRO*
 57 is evaluated against PSRO on a collection of
 58 partially observable general-sum games. In our
 59 experiments, *Dyna-PSRO* found lower-regret
 60 solutions while requiring substantially fewer cu-
 61 mulative experiences.

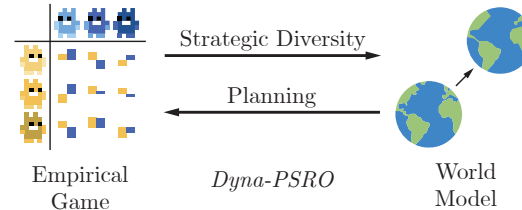


Figure 1: *Dyna-PSRO* co-learns a world model and empirical game. Empirical games offer world models strategically diverse game dynamics. World models offer empirical games more efficient strategy discovery through planning.

62 The main points of novelty of this paper are as follows: (1) empirically demonstrate that world models
 63 benefit from the strategic diversity induced by an empirical game; (2) empirically demonstrate that a
 64 world model can be effectively transferred and used in planning with new other-players. The major
 65 contribution of this work is a new algorithm, *Dyna-PSRO*, that co-learns an empirical game and
 66 world model finding a stronger solution at less cost than the baseline, PSRO.

67 2 Related Work

68 **Empirical Game Theoretic Analysis (EGTA).** The core idea of EGTA [80] is to reason over
 69 approximate game models (*empirical games*) estimated by simulation over a restricted strategy set.
 70 This basic approach was first demonstrated by Walsh et al. [77], in a study of pricing and bidding
 71 games. Phelps et al. [51] introduced the idea of extending a strategy set automatically through
 72 optimization, employing genetic search over a policy space. Schwartzman & Wellman [58] proposed
 73 using RL to derive new strategies that are approximate best responses (BRs) to the current empirical
 74 game’s Nash equilibrium. The general question of which strategies to add to an empirical game
 75 has been termed the *strategy exploration problem* [31]. PSRO [38] generalized the target for BR
 76 beyond NE, and introduced DRL for BR computation in empirical games. Many further variants and
 77 extensions of EGTA have been proposed, for example those using structured game representations
 78 such as extensive-form [43, 34]. Some prior work has considered transfer learning across BR
 79 computations in EGTA, specifically by reusing elements of policies and value functions [64, 65].

80 **Model-Based Reinforcement Learning (MBRL).** *Model-Based* RL algorithms construct or use
 81 a model of the environment (henceforth, *world model*) in the process of learning a policy or value
 82 function [69]. World models may either predict successor observations directly (e.g., at pixel
 83 level [76, 79]), or in a learned latent space [18, 17]. The world models can be either used for
 84 *background planning* by rolling out model-predicted trajectories to train a policy, or by *decision-*
 85 *time planning* where the world model is used to evaluate the current state by planning into the
 86 future. Talvitie [71] demonstrated that even in small Markov decision processes (MDP) [52], model-
 87 prediction errors tend to compound—rendering long-term planning at the abstraction of observations
 88 ineffective. A follow-up study demonstrated that for imperfect models, short-term planning was
 89 no better than repeatedly training on previously collected real experiences; however, medium-term
 90 planning offered advantages even with an imperfect model [27]. Parallel studies hypothesized that
 91 these errors are a result of insufficient data for that transition to be learned [36, 8]. To remedy
 92 the data insufficiency, ensembles of world models were proposed to account for world model

93 uncertainty [8, 36, 84], and another line of inquiry used world model uncertainty to guide exploration
 94 in state-action space [3, 59]. This study extends this problem into the multiagent setting, where
 95 now other-agents may preclude transitions from occurring. The proposed remedy is to leverage the
 96 strategy exploration process of building an empirical game to guide data generation.

97 **Multiagent Reinforcement Learning (MARL).** Previous research intersecting MARL and MBRL
 98 has primarily focused on modeling the opponent, particularly in scenarios where the opponent is fixed
 99 and well-defined. Within specific game sub-classes, like cooperative games and two-player zero-sum
 100 games, it has been theoretically shown that opponent modeling reduces the sample complexity of
 101 RL [73, 85]. Opponent models can either explicitly [46, 15] or implicitly [4, 29] model the behavior
 102 of the opponent. Additionally, these models can either construct a single model of opponent behavior,
 103 or learn a set of models [12, 21]. While opponent modeling details are beyond the scope of this
 104 study, readers can refer to Albrecht & Stone’s survey [1] for a comprehensive review on this subject.
 105 Instead, we consider the case where the learner has explicit access to the opponent’s policy during
 106 training, as is the case in empirical-game building. A natural example is that of Self-Play, where all
 107 agents play the same policy; therefore, a world model can be learned used to evaluate the quality of
 108 actions with Monte-Carlo Tree Search [60, 62, 72, 56]. Li et al. [41] expands on this by building a
 109 population of candidate opponent policies through PSRO to augment the search procedure. Krupnik
 110 et al. [35] demonstrated that a generative world model could be useful in multi-step opponent-action
 111 prediction. Sun et al. [66] examined modeling stateful game dynamics from observations when
 112 the agents’ policies are stationary. Chockalingam et al. [11] explored learning world models for
 113 homogeneous agents with a centralized controller in a cooperative game. World models may also be
 114 shared by independent reinforcement learners in cooperative games [81, 86].

115 3 Co-Learning Benefits

116 We begin by specifying exactly what we mean by world model and empirical game. This requires
 117 defining some primitive elements. Let $t \in \mathcal{T}$ denote time in the real game, with $s^t \in \mathcal{S}$ the
 118 **information state** and $h^t \in \mathcal{H}$ the **game state** at time t . The information state $s^t \equiv (m^{\pi,t}, o^t)$
 119 is composed of the **agent’s memory** $m^\pi \in \mathcal{M}^\pi$, or recurrent state, and the current **observation**
 120 $o \in \mathcal{O}$. Subscripts denote a player-specific component s_i , negative subscripts denote all but the
 121 player s_{-i} , and boldface denote the joint of all players \mathbf{s} . The **transition dynamics** $p : \mathcal{H} \times \mathcal{A} \rightarrow$
 122 $\Delta(\mathcal{H}) \times \Delta(\mathcal{R})$ define the game state update and reward signal. The agent experiences **transitions**, or
 123 **experiences**, $(s^t, a^t, r^{t+1}, s^{t+1})$ of the game; where, sequences of transitions are called **trajectories** τ
 124 and trajectories ending in a terminal game state are **episodes**.

125 At the start of an episode, all players sample their current **policy** π from their **strategy** $\sigma : \Pi \rightarrow$
 126 $[0, 1]$, where Π is the **policy space** and Σ is the corresponding **strategy space**. A **utility function**
 127 $U : \mathbf{\Pi} \rightarrow \mathbb{R}^n$ defines the payoffs/returns (i.e., cumulative reward) for each of n players. The tuple
 128 $\Gamma \equiv (\mathbf{\Pi}, U, n)$ defines a **normal-form game** (NFG) based on these elements. We represent empirical
 129 games in normal form. An **empirical normal-form game** (ENFG) $\hat{\Gamma} \equiv (\hat{\mathbf{\Pi}}, \hat{U}, n)$ models a game
 130 with a **restricted strategy set** $\hat{\mathbf{\Pi}}$ and an estimated payoff function \hat{U} . An empirical game is typically
 131 built by alternating between game reasoning and strategy exploration. During the game reasoning
 132 phase, the empirical game is solved based on a solution concept predefined by the modeler. The
 133 strategy exploration step uses this solution to generate new policies to add to the empirical game. One
 134 common heuristic is to generate new policies that best-respond to the current solution [45, 57]. As
 135 exact best-responses typically cannot be computed, RL or DRL are employed to derive approximate
 136 best-responses [38].

137 An **agent world model** w represents dynamics in terms of information available to the agent. Specifi-
 138 cally, w maps information states and actions to observations and rewards, $w : \mathcal{O} \times \mathcal{A} \times \mathcal{M}^w \rightarrow \mathcal{O} \times \mathcal{R}$,
 139 where $m^w \in \mathcal{M}^w$ is the **world model’s memory**, or recurrent state. For simplicity, in this work, we
 140 assume the agent learns and uses a deterministic world model, irrespective of stochasticity that may be
 141 present in the true game. Specific implementation details for this work are provided in Appendix C.2.

142 Until now, we have implicitly assumed the need for distinct models. However, if a single model could
 143 serve both functions, co-learning two separate models would not be needed. Empirical games, in
 144 general, cannot replace a world model as they entirely abstract away any concept of game dynamics.
 145 Conversely, world models have the potential to substitute for the payoff estimations in empirical
 146 games by estimating payoffs as rollouts with the world model. We explore this possibility in an

147 auxiliary experiment included in Appendix E.4, but our findings indicate that this substitution is
148 impractical. Due to compounding of model-prediction errors, the payoff estimates and entailed game
149 solutions were quite inaccurate.

150 Having defined the models and established the need for their separate instantiations, we can proceed
151 to evaluate the claims of beneficial co-learning. Our first experiment shows that the strategic diversity
152 embodied in an empirical game yields diverse game dynamics, resulting in the training of a more
153 performant world model. The second set of experiments demonstrates that a world model can help
154 reduce the computational cost of policy construction in an empirical game.




155 3.1 Strategic Diversity

156 A world model is trained to predict successor observations and rewards, from the current observations
157 and actions, using a supervised learning signal. Ideally, the training data would cover all possible
158 transitions. This is not feasible, so instead draws are conventionally taken from a dataset generated
159 from play of a *behavioral strategy*. Performance of the world model is then measured against a *target*
160 *strategy*. Differences between the behavioral and target strategies present challenges in learning an
161 effective world model.

162 We call the probability of drawing a state-action pair under some strategy its *reach probability*. From
163 this, we define a strategy’s *strategic diversity* as the distribution induced from reach probabilities.
164 across the full state-action space. These terms allow us to observe two challenges for learning world
165 models. First, the diversity of the behavioral strategy ought to *cover* the target strategy’s diversity.
166 Otherwise, transitions will be absent from the training data. It is possible to supplement coverage of
167 the absent transitions if they can be generalized from covered data; however, this cannot be generally
168 guaranteed. Second, the *closer* the diversities are, the more accurate the learning objective will be.
169 An extended formal argument of these challenges is provided in Appendix C.3.

170 If the target strategy were known, we could readily construct the ideal training data for the world
171 model. However the target is generally not known at the outset; indeed determining this target is the
172 ultimate purpose of empirical game reasoning. The evolving empirical game essentially reflects a
173 search for the target. Serendipitously, construction of this empirical game entails generation of data
174 that captures elements of likely targets. This data can be reused for world model training without
175 incurring any additional data collection cost.

176 **Game.** We evaluate the claims of independent co-learning benefits within the context of a *commons*
177 *game* called “Harvest”. In Harvest, players move around an orchard picking apples. The challenging
178 commons element is that apple regrowth rate is proportional to nearby apples, so that socially optimum
179 behavior would entail managed harvesting. Self-interested agents capture only part of the benefit of
180 optimal growth, thus non-cooperative equilibria tend to exhibit collective over-harvesting. The game
181 has established roots in human-behavioral studies [30] and in agent-based modeling of emergent
182 behavior [53, 40, 39]. For our initial experiments, we use a symmetric two-player version of the game,
183 where in-game entities are represented categorically [28]. Each player has a 10×10 viewbox within
184 their field of vision. The possible actions include moving in the four cardinal directions, rotating
185 either way, tagging, or remaining idle. A successful tag temporarily removes the other player from
186 the game, but can only be done to other nearby players. Players receive a reward of 1 for each apple
187 picked. More detailed information and visualizations are available in Appendix D.1.

188 **Experiment.** To test the effects of strategic diversity, we train a suite of world models that differ
189 in the diversity of their training data. The datasets are constructed from the play of three policies:
190 a random baseline policy, and two PSRO-generated policies. The PSRO policies were arbitrarily
191 sampled from an approximate solution produced by a run of PSRO. We sampled an additional
192 policy from PSRO for evaluating the generalization capacity of the world models. These policies
193 are then subsampled and used to train seven world models. The world models are referred to by
194 icons  that depict the symmetric strategy profiles used to train them in the normal-form. Strategy
195 profiles included in the training data of the world models are shaded black. For instance, the first
196 (random) policy , or the first and third policies . Each world model’s dataset contains 1 million
197 total transitions, collected uniformly from each distinct strategy profile (symmetric profiles are not
198 re-sampled). The world models are then evaluated on accuracy and recall for their predictions of both

199 observation and reward for both players. The world models are optimized with a weighted-average
 200 cross-entropy objective. Additional details are in Appendix C.2.

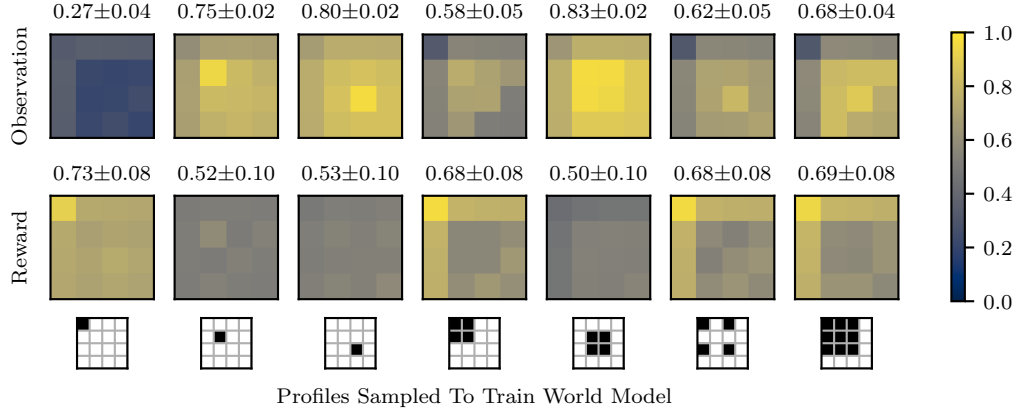


Figure 2: World model accuracy across strategy profiles. Each heatmap portrays a world model’s accuracy over 16 strategy profiles. The meta x-axis corresponds to the profiles used to train the world model (as black cells). Above each heatmap is the model’s average accuracy.

201 **Results.** Figure 2 presents each world model’s per-profile accuracy, as well as its average
 202 over all profiles. Inclusion of the random policy corresponds to decreases in observation
 203 prediction accuracy: $\begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ $0.75 \pm 0.02 \rightarrow \begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ 0.58 ± 0.05 , $\begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ $0.80 \pm 0.02 \rightarrow \begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ 0.62 ± 0.05 ,
 204 and $\begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ $0.83 \pm 0.02 \rightarrow \begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ 0.68 ± 0.04 . Figure 13 (Appendix E.1) contains the world
 205 model’s per-profile recall. Inclusion of the random policy corresponds to increases in reward
 206 1 recall: $\begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ $0.25 \pm 0.07 \rightarrow \begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ 0.37 ± 0.11 , $\begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ $0.25 \pm 0.07 \rightarrow \begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ 0.36 ± 0.11 , and
 207 $\begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ $0.26 \pm 0.07 \rightarrow \begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ 0.37 ± 0.11 .

208 **Discussion.** The PSRO policies offer the most strategically salient view of the game’s dynamics.
 209 Consequently, the world model $\begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ trained with these policies yields the highest observation accuracy.
 210 However, this world model performs poorly on reward accuracy, scoring only 0.50 ± 0.10 . In
 211 comparison, the model trained on the random policy $\begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ scores 0.73 ± 0.08 . This seemingly
 212 counterintuitive result can be attributed to a significant class imbalance in rewards. $\begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ predicts only
 213 the most common class, no reward, which gives the illusion of higher performance. In contrast, the
 214 remaining world models attempt to predict rewarding states, which reduces their overall accuracy.
 215 Therefore, we should compare the world models based on their ability to recall rewards. When we
 216 examine $\begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ again, we find that it also struggles to recall rewards, scoring only 0.26 ± 0.07 . However,
 217 when the random policy is included in the training data ($\begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$), the recall improves to 0.37 ± 0.11 . This
 218 improvement is also due to the same class imbalance. The PSRO policies are highly competitive,
 219 tending to over-harvest. This limits the proportion of rewarding experiences. Including the random
 220 policy enhances the diversity of rewards in this instance, as its coplayer can demonstrate successful
 221 harvesting. Given the importance of accurately predicting both observations and rewards for effective
 222 planning, $\begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$ appears to be the most promising option. However, the strong performance of $\begin{smallmatrix} \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \\ \blacksquare & \blacksquare & \blacksquare \end{smallmatrix}$
 223 suggests future work on algorithms that can benefit solely from observation predictions. Overall,
 224 these results support the claim that strategic diversity enhances the training of world models.

225 3.2 Response Calculations


226 Empirical games are built by iteratively calculating and incorporating responses to the current
 227 solution. However, direct computation of these responses is often infeasible, so RL or DRL is used
 228 to approximate the response. This process of approximating a single response policy using RL is
 229 computationally intensive, posing a significant constraint in empirical game modeling when executed
 230 repeatedly. World models present an opportunity to address this issue. A world model can serve as a
 231 medium for transferring previously learned knowledge about the game’s dynamics. Therefore, the
 232 dynamics need not be relearned, reducing the computational cost associated with response calculation.

233 Exercising a world model for transfer is achieved through a process called *planning*. Planning is
 234 any procedure that takes a world model and produces or improves a policy. In the context of games,
 235 planning can optionally take into account the existence of coplayers. This consideration can reduce
 236 experiential variance caused by unobserved confounders (i.e., the coplayers). However, coplayer
 237 modeling errors may introduce further errors in the planning procedure [21].

238 Planning alongside empirical-game construction allows us to side-step this issue as we have direct
 239 access to the policies of all players during training. This allows us to circumvent the challenge
 240 of building accurate agent models. Instead, the policies of coplayers can be directly queried and
 241 used alongside a world model, leading to more accurate planning. In this section, we empirically
 242 demonstrate the effectiveness of two methods that decrease the cost of response calculation by
 243 integrating planning with a world model and other agent policies.

244 3.2.1 Background Planning

245 The first type of planning that is investigated is *background planning*, popularized by the Dyna
 246 architecture [67]. In background planning, agents interact with the world model to produce *planned*
 247 *experiences*¹. The planned experiences are then used by a model-free reinforcement learning
 248 algorithm as if they were *real experiences* (experiences generated from the real game). Background
 249 planning enables learners to generate experiences of states they are not currently in.

250 **Experiment.** To assess whether planned experiences are effective for training a policy in the actual
 251 game, we compute two response policies. The first response policy, serving as our baseline, learns
 252 exclusively from real experiences. The second response policy, referred to as the planner, is trained
 253 using a two-step procedure. Initially, the planner is exclusively trained on planned experiences. After
 254 10 000 updates, it then transitions to learning solely from real experiences. Policies are trained using
 255 IMPALA [14], with further details available in Appendix C.1. The planner employs the  world
 256 model from Section 3.1, and the opponent plays the previously held-out policy. In this and subsequent
 257 experiments, the cost of methods is measured by the number of experiences they require with the
 258 actual game. This is because, experience collection is often the bottleneck when applying RL-based
 259 methods [48, 24]. Throughout the remainder of this work, each experience represents a trajectory of
 260 20 transitions, facilitating the training of recurrent policies.

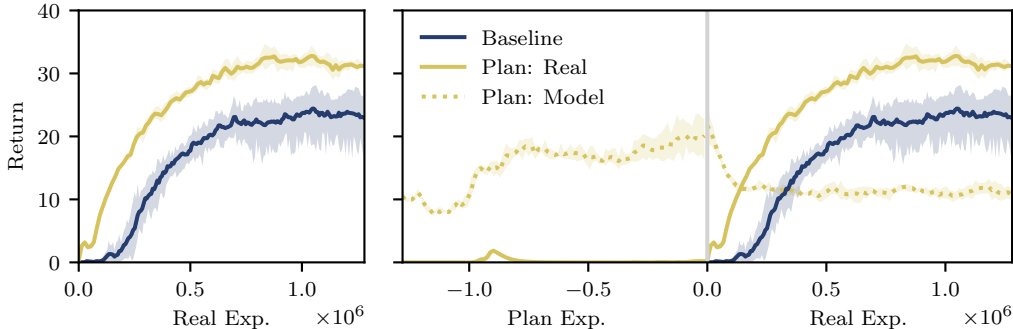


Figure 3: Effects of background planning on response learning. Left: Return curves measured by the number of real experiences used. Right: Return curves measured by usage of both real and planned experiences. The planner’s return is measured against the real game and the world model. (5 seeds, with 95 % bootstrapped CI).

261 **Results.** Figure 3 presents the results of the background planning experiment. The methods are
 262 compared based on their final return, utilizing an equivalent amount of real experiences. The baseline
 263 yields a return of 23.00 ± 4.01 , whereas the planner yields a return of 31.17 ± 0.25 .

264 **Discussion.** In this experiment, the planner converges to a stronger policy, and makes earlier gains
 265 in performance than the baseline. Despite this, there is a significant gap in the planner’s learning

¹Other names include “imaginary”, “simulated”, or “hallucinated” experiences.

266 curves, which are reported with respect to both the world model and real game. This gap arises due
 267 to accumulated model-prediction errors, causing the trajectories to deviate from the true state space.
 268 Nevertheless, the planner effectively learns to interact with the world model during planning, and
 269 this behavior shows positive transfer into the real game, as evidenced by the planner’s rapid learning.
 270 The exact magnitude of benefit will vary across coplayers’ policies, games, and world models. In
 271 Figure 14 (Appendix E.2), we repeat the same experiment with the poorly performing \boxtimes world
 272 model, and observe a marginal benefit (26.05 ± 1.32). The key take-away is that background planning
 273 tends to lead towards learning benefits, and not generally hamper learning.

274 3.2.2 Decision-Time Planning

275 The second main way that a world model is used is to inform action selection at *decision time*
 276 [*planning*] (*DT*). In this case, the agent evaluates the quality of actions by comparing the value of
 277 the model’s predicted successor state for all candidate actions. Action evaluation can also occur
 278 recursively, allowing the agent to consider successor states further into the future. Overall, this
 279 process should enable the learner to select better actions earlier in training, thereby reducing the
 280 amount of experiences needed to compute a response. A potential flaw with decision-time planning
 281 is that the agent’s learned value function may not be well-defined on model-predicted successor
 282 states [71]. To remedy this issue, the value function should also be trained on model-predicted states.

283 **Experiment.** To evaluate the impact the decision-time planning, we perform an experiment similar
 284 to the background planning experiment (Section 3.2.1). However, in this experiment, we evaluate
 285 the quality of four types of decision-time planners that perform one-step three-action search. The
 286 planners differ in their ablations of background planning types: (1) *warm-start background*
 287 *planning (BG: W)* learning from planned experiences before any real experiences, and (2) *concurrent*
 288 *background planning (BG: C)* where after BG: W, learning proceeds simultaneously on both planned
 289 and real experiences. The intuition behind BG: C is that the agent can complement its learning
 290 process by incorporating planned experiences that align with its current behavior, offsetting the
 291 reliance on costly real experiences. Extended experimental details are provided in Appendix C.

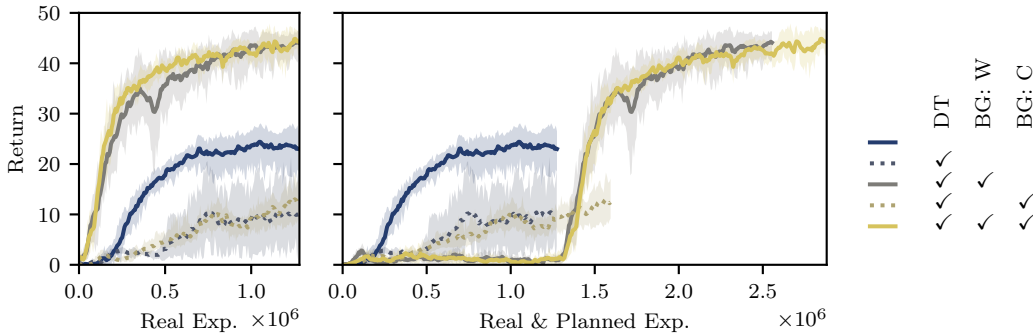


Figure 4: Effects of decision-time planning on response learning. Four planners using decision-time planning (DT) are shown in combinations with warm-start background planning (BG: W) and concurrent background planning (BG: C). (5 seeds, with 95 % bootstrapped CI).

292 **Results.** The results for this experiment are shown in Figure 4. The baseline policy receives a final
 293 return of 23.00 ± 4.01 . The planners that do not include BG: W, perform worse, with final returns of
 294 9.98 ± 7.60 (DT) and 12.42 ± 3.97 (DT & BG: C). The planners that perform BG: W outperform the
 295 baseline, with final returns of 44.11 ± 2.81 (DT & BG: W) and 44.31 ± 2.56 (DT, BG: W, & BG: C).

296 **Discussion.** Our results suggest that the addition of BG: W provides sizable benefits: 9.98 ± 7.60
 297 (DT) \rightarrow 44.11 ± 2.81 (DT & BG:W) and 12.42 ± 3.97 (DT & BG: C) \rightarrow 44.31 ± 2.56 (DT, BG: W,
 298 & BG: C). We postulate that this is because it informs the policy’s value function on model-predictive
 299 states early into training. This allows that the learner is able to more effectively search earlier into
 300 training. BG: C appears to offer minor stability and variance improvements throughout the training

301 procedure; however, it does not have a measurable difference in final performance. This result
302 suggests using planning methods in combination to reap their respective advantages.

303 However, we caution against focusing on the magnitude of improvement found within this experiment.
304 As the margin of benefit depends on many factors including the world model accuracy, the opponent
305 policy, and the game. To exemplify, similar to the background planning section, we repeat the same
306 experiment with the poorly performing \blacksquare world model. The results of this ancillary experiment are
307 in Figure 15 (Appendix E.3). The trend of BG: W providing benefits was reinforced: 6.29 ± 5.12
308 (DT) $\rightarrow 20.98 \pm 9.76$ (DT & BG: W) and 3.64 ± 0.26 (DT & BG: C) $\rightarrow 33.07 \pm 7.67$ (DT, BG: W,
309 & BG: C). However, the addition of BG: C now measurably improved performance 20.98 ± 9.76
310 (DT & BG: W) $\rightarrow 33.07 \pm 7.67$ (DT, BG: W, & BG: C). The main outcome of these experiments
311 is the observation that multi-faceted planning is unlikely to harm a response calculation, and has a
312 potentially large benefit when applied effectively. These results support the claim that world models
313 offer the potential to improve response calculation through decision-time planning.

314 4 Dyna-PSRO

315 In this section we introduce Dyna-PSRO, *Dyna*-Policy-Space Response Oracles, an approximate
316 game-solving algorithm that builds on the PSRO [38] framework. Dyna-PSRO employs co-learning
317 to combine the benefits of world models and empirical games.

318 Dyna-PSRO is defined by two significant alterations to the original PSRO algorithm. First, it trains
319 a world model in parallel with all the typical PSRO routines (i.e., game reasoning and response
320 calculation). We collect training data for the world model from both the episodes used to estimate the
321 empirical game’s payoffs, and the episodes that are generated during response learning and evaluation.
322 This approach ensures that the world model is informed by a diversity of data from a salient set of
323 strategy profiles. By reusing data from empirical game development, training the world model incurs
324 no additional cost for data collection.

325 The second modification introduced by Dyna-PSRO pertains to the way response policies are learned.
326 Dyna-PSRO adopts a Dyna-based reinforcement learner [67, 68, 70] that integrates simultaneous plan-
327 ning, learning, and acting. Consequently, the learner concurrently processes experiences generated
328 from decision-time planning, background planning, and direct game interaction. These experiences,
329 regardless of their origin, are then learned from using the IMPALA [14] update rule. For all accounts
330 of planning, the learner uses the single world model that is trained within Dyna-PSRO. This allows
331 game knowledge accrued from previous response calculations to be transferred and used to reduce
332 the cost of the current and future response calculations. Pseudocode and additional details for both
333 PSRO and Dyna-PSRO are provided in Appendix C.4.

334 **Games.** Dyna-PSRO is evaluated on three games. The first is the harvest commons game used in the
335 experiments described above, denoted “Harvest: Categorical”. The other two games come from the
336 MeltingPot [39] evaluation suite and feature rich image-based observations. “Harvest: RGB” is their
337 version of the same commons harvest game (details in Appendix D.2). “Running With Scissors” is a
338 temporally extended version of rock-paper-scissors (details in Appendix D.3). World model training
339 and implementation details for each game are in Appendix C.2, likewise, policies in Appendix C.1.

340 **Experiment.** Dyna-PSRO’s performance is measured by the quality of the solution it produces
341 when compared against the world-model-free baseline PSRO. The two methods are evaluated on
342 SumRegret (sometimes called *Nash convergence*), which measures the regret across all players
343 $\text{SumRegret}(\sigma, \bar{\Pi}) = \sum_{i \in n} \max_{\pi_i \in \bar{\Pi}_i} \hat{U}_i(\pi_i, \sigma_{-i}) - \hat{U}_i(\sigma_i, \sigma_{-i})$, where σ is the method’s solution
344 and $\bar{\Pi} \subseteq \Pi$ denotes the deviation set. We define deviation sets based on policies generated across
345 methods (i.e., regret is with respect to the *combined game*): $\bar{\Pi} \equiv \bigcup_{\text{method}} \hat{\Pi}^{\text{method}}$, for all methods for
346 a particular seed (detailed in Appendix C.5) [2]. We measure SumRegret for intermediate solutions,
347 and report it as a function of the cumulative number of real experiences employed in the respective
348 methods.

349 **Results.** Figure 5 presents the results for this experiment. For Harvest: Categorical, Dyna-PSRO
350 found a no regret solution within the combined-game in $3.2e6$ experiences. Whereas, PSRO achieves
351 a solution of at best 5.45 ± 1.62 within $2e7$ experiences. In Harvest: RGB, Dyna-PSRO reaches a

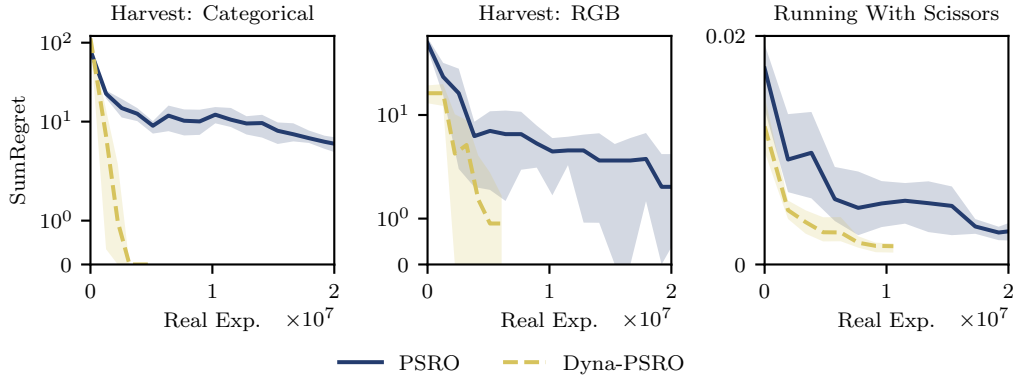


Figure 5: PSRO compared against Dyna-PSRO. (5 seeds, with 95 % bootstrapped CI).

352 solution with 0.89 ± 0.74 regret at $5.12e6$ experiences. At the same time, PSRO had found a solution
 353 with 6.42 ± 4.73 regret, and at the end of its run had 2.50 ± 2.24 regret. In the final game, RWS,
 354 Dyna-PSRO has $2e-3 \pm 5e-4$ regret at $1.06e7$ experiences, and at a similar point ($9.6e6$ experiences),
 355 PSRO has $6.68e-3 \pm 2.51e-3$. At the end of the run, PSRO achieves a regret $3.50e-3 \pm 7.36e-4$.

356 **Discussion.** The results indicate that across all games, Dyna-PSRO consistently outperforms PSRO
 357 by achieving a superior solution. Furthermore, this improved performance is realized while consuming
 358 fewer real-game experiences. For instance, in the case of Harvest: Categorical, the application of
 359 the world model for decision-time planning enables the computation of an effective policy after only
 360 a few iterations. On the other hand, we observe a trend of accruing marginal gains in other games,
 361 suggesting that the benefits are likely attributed to the transfer of knowledge about the game dynamics.
 362 In Harvest: Categorical and Running With Scissors, Dyna-PSRO also had lower variance than PSRO.

363 5 Limitations

364 Although our experiments demonstrate benefits for co-learning world models and empirical games,
 365 there are several areas for potential improvement. The world models used in this study necessitated
 366 observational data from all players for training, and assumed a simultaneous-action game. Future
 367 research could consider relaxing these assumptions to accommodate different interaction protocols,
 368 a larger number of players, and incomplete data perspectives. Furthermore, our world models
 369 functioned directly on agent observations, which made them computationally costly to query. If
 370 the generation of experiences is the major limiting factor, as assumed in this study, this approach is
 371 acceptable. Nevertheless, reducing computational demands through methods like latent world models
 372 presents a promising avenue for future research. Lastly, the evaluation of solution concepts could
 373 also be improved. While combined-game regret employs all available estimates in approximating
 374 regret, its inherent inaccuracies may lead to misinterpretations of relative performance.

375 6 Conclusion

376 This study showed the mutual benefit of co-learning a world model and empirical game. First, we
 377 demonstrated that empirical games provide strategically diverse training data that could inform a more
 378 robust world model. We then showed that world models can reduce the computational cost, measured
 379 in experiences, of response calculations through planning. These two benefits were combined and
 380 realized in a new algorithm, Dyna-PSRO. In our experiments, Dyna-PSRO computed lower-regret
 381 solutions than PSRO on several partially observable general-sum games. Dyna-PSRO also required
 382 substantially fewer experiences than PSRO, a key algorithmic advantage for settings where collecting
 383 experiences is a cost-limiting factor.

384 **References**

- 385 [1] Stefano V Albrecht and Peter Stone. Autonomous agents modelling other agents: A compre-
386 hensive survey and open problems. *Artificial Intelligence*, 258:66–95, 2018.
- 387 [2] David Balduzzi, Karl Tuyls, Julien Pérolat, and Thore Graepel. Re-evaluating evaluation. In
388 *32nd Conference on Neural Information Processing Systems*, 2018.
- 389 [3] Philip Ball, Jack Parker-Holder, Aldo Pacchiano, Krzysztof Choromanski, and Stephen Roberts.
390 Ready policy one: World building through active learning. In *37th International Conference of*
391 *Machine Learning*, 2020.
- 392 [4] Nolan Bard, Michael Johanson, Neil Burch, and Michael Bowling. Online implicit agent
393 modelling. In *12th International Conference on Autonomous Agents and Multiagent Systems*,
394 2013.
- 395 [5] Samy Bengio, Oriol Vinyals, Navdeep Jaitly, and Noam Shazeer. Scheduled sampling for
396 sequence prediction with recurrent neural networks. In *28th Conference on Neural Information*
397 *Processing Systems*, pages 1171–1179, 2015.
- 398 [6] James Bradbury, Roy Frostig, Peter Hawkins, Matthew James Johnson, Chris Leary, Dougal
399 Maclaurin, George Necula, Adam Paszke, Jake VanderPlas, Skye Wanderman-Milne, and Qiao
400 Zhang. JAX: composable transformations of Python+NumPy programs, 2018.
- 401 [7] George W Brown. Iterative solution of games by fictitious play. In *Activity analysis of production*
402 *and allocation*, volume 13, pages 374–376, 1951.
- 403 [8] Jacob Buckman, Danijar Hafner, George Tucker, Eugene Brevdo, and Honglak Lee. Sample-
404 efficient reinforcement learning with stochastic ensemble value expansion. In *22nd Conference*
405 *on Neural Information Processing Systems*, 2018.
- 406 [9] Albin Cassirer, Gabriel Barth-Maron, Eugene Brevdo, Sabela Ramos, Toby Boyd, Thibault
407 Sottiaux, and Manuel Kroiss. Reverb: A framework for experience replay, 2021.
- 408 [10] Silvia Chiappa, Sébastien Racaniere, Daan Wierstra, and Shakir Mohamed. Recurrent environ-
409 ment simulators. In *5th International Conference on Learning Representations*, 2017.
- 410 [11] Valliappa Chockingam, Tegg Taekyong Sung, Feryal Behbanai, Rishab Gargeya, Amlesh
411 Sivanantham, and Aleksandra Malysheva. Extending world models for multi-agent reinforce-
412 ment learning in malmö. In *Joint AIIDE 2018 Workshops co-located with the 14th AAAI*
413 *conference on artificial intelligence and interactive digital entertainment*, 2018.
- 414 [12] Brian Collins. Combining opponent modeling and model-based reinforcement learning in a
415 two-player competitive game. Master’s thesis, University of Edinburgh, 2007.
- 416 [13] B. Curtis Eaves. The linear complementarity problem. *Management Science*, 17(9):612–634,
417 1971.
- 418 [14] Lasse Espeholt, Hubert Soyer, Remi Munos, Karen Simonyan, Volodymyr Mnih, Tom Ward,
419 Yotam Doron, Vlad Firoiu, Tim Harley, Iain Dunning, Shane Legg, and Koray Kavukcuoglu.
420 IMPALA: Scalable distributed deep-RL with importance weighted actor-learner architectures.
421 In *35th International Conference on Machine Learning*, 2018.
- 422 [15] Jakob N Foerster, Richard Y Chen, Maruan Al-Shedivat, Shimon Whiteson, Pieter Abbeel, and
423 Igor Mordatch. Learning with opponent-learning awareness. In *17th International Conference*
424 *on Autonomous Agents and MultiAgent Systems*, 2018.
- 425 [16] Kunihiko Fukushima. Cognitron: A self-organizing multilayered neural network. *Biological*
426 *Cybernetics*, 20:121–136, 1975.
- 427 [17] Carles Gelada, Saurabh Kumar, Jacob Buckman, Ofir Nachum, and Marc G. Bellemare. Deep-
428 MDP: Learning continuous latent space models for representation learning. In *36th International*
429 *Conference on Machine Learning*, volume 97, pages 2170–2179, 2019.

- 430 [18] David Ha and Jürgen Schmidhuber. Recurrent world models facilitate policy evolution. In *31st*
431 *Conference on Neural Information Processing Systems*, 2018.
- 432 [19] David Ha and Jürgen Schmidhuber. World models. In *arXiv preprint arXiv:1803.10122*, 2018.
- 433 [20] Danijar Hafner, Timothy Lillicrap, Mohammad Norouzi, and Jimmy Ba. Mastering atari with
434 discrete world models. In *9th International Conference on Learning Representations*, 2021.
- 435 [21] He He, Jordan Boyd-Graber, Kevin Kwok, and Hal Daumé III. Opponent modeling in deep
436 reinforcement learning. In *33rd International Conference on Machine Learning*, 2016.
- 437 [22] Tom Hennigan, Trevor Cai, Tamara Norman, and Igor Babuschkin. Haiku: Sonnet for JAX,
438 2020.
- 439 [23] Pablo Hernandez-Leal, Michael Kaisers, Tim Baarslag, and Enrique Munoz de Cote. A
440 survey of learning in multiagent environments: Dealing with non-stationarity. *arXiv preprint*
441 *arXiv:1707.09183*, 2017.
- 442 [24] Todd Hester and Peter Stone. Texploré: Real-time sample-efficient reinforcement learning for
443 robots. In *Machine Learning for Robotics (MLR)*, 2012.
- 444 [25] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*,
445 9(8):1735–1780, 1997.
- 446 [26] Matthew W. Hoffman, Bobak Shahriari, John Aslanides, Gabriel Barth-Maron, Nikola Momchev,
447 Danila Sinopalnikov, Piotr Stańczyk, Sabela Ramos, Anton Raichuk, Damien Vincent, Léonard
448 Hussenot, Robert Dadashi, Gabriel Dulac-Arnold, Manu Orsini, Alexis Jacq, Johan Ferret, Nino
449 Vieillard, Seyed Kamyar Seyed Ghasemipour, Sertan Girgin, Olivier Pietquin, Feryal Behbahani,
450 Tamara Norman, Abbas Abdolmaleki, Albin Cassirer, Fan Yang, Kate Baumli, Sarah Henderson,
451 Abe Friesen, Ruba Haroun, Alex Novikov, Sergio Gómez Colmenarejo, Serkan Cabi, Caglar
452 Gulcehre, Tom Le Paine, Srivatsan Srinivasan, Andrew Cowie, Ziyu Wang, Bilal Piot, and
453 Nando de Freitas. Acme: A research framework for distributed reinforcement learning. *arXiv*
454 *preprint arXiv:2006.00979*, 2020.
- 455 [27] G. Zacharias Holland, Erin Talvitie, and Michael Bowling. The effect of planning shape on
456 dyna-style planning in high-dimensional state spaces. In *FAIM workshop “Prediction and*
457 *Generative Modeling in Reinforcement Learning”*, 2018.
- 458 [28] HumanCompatibleAI. <https://github.com/HumanCompatibleAI/multi-agent>, 2019.
- 459 [29] Pararawendy Indarjo. Deep state-space models in multi-agent systems. Master’s thesis, Leiden
460 University, 2019.
- 461 [30] Marco A. Janssen, Robert Holahan, Allen Lee, and Elinor Ostrom. Lab experiments for the
462 study of social-ecological systems. *Science*, 328(5978):613–617, 2010.
- 463 [31] Patrick R. Jordan, L. Julian Schvartzman, and Michael P. Wellman. Strategy exploration in
464 empirical games. In *9th International Conference on Autonomous Agents and Multi-Agent*
465 *Systems*, pages 1131–1138, 2010.
- 466 [32] Gabriel Kalweit and Joschka Boedecker. Uncertainty-driven imagination for continuous deep
467 reinforcement learning. In *1st Conference on Robot Learning*, pages 195–206, 2017.
- 468 [33] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd*
469 *International Conference for Learning Representations*, 2015.
- 470 [34] Christine Konicki, Mithun Chakraborty, and Michael P. Wellman. Exploiting extensive-form
471 structure in empirical game-theoretic analysis. In *Web and Internet Economics: 18th Interna-*
472 *tional Conference*, 2022.
- 473 [35] Orr Krupnik, Igor Mordatch, and Aviv Tamar. Multi-agent reinforcement learning with multi-
474 step generative models. In *4th Conference on Robot Learning*, pages 776–790, 2020.

- 475 [36] Thanard Kurutach, Ignasi Clavera, Yan Duan, Aviv Tamar, and Pieter Abbeel. Model-ensemble
476 trust-region policy optimization. In *6th International Conference on Learning Representations*,
477 2018.
- 478 [37] Marc Lanctot, Edward Lockhart, Jean-Baptiste Lespiau, Vinicius Zambaldi, Satyaki Upadhyay,
479 Julien Pérolat, Sriram Srinivasan, Finbarr Timbers, Karl Tuyls, Shayegan Omidshafiei, Daniel
480 Hennes, Dustin Morrill, Paul Muller, Timo Ewalds, Ryan Faulkner, János Kramár, Bart De
481 Vylder, Brennan Saeta, James Bradbury, David Ding, Sebastian Borgeaud, Matthew Lai,
482 Julian Schrittwieser, Thomas Anthony, Edward Hughes, Ivo Danihelka, and Jonah Ryan-Davis.
483 OpenSpiel: A framework for reinforcement learning in games. *CoRR*, abs/1908.09453, 2019.
- 484 [38] Marc Lanctot, Vinicius Zambaldi, Audrūnas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien
485 Pérolat, David Silver, and Thore Graepel. A unified game-theoretic approach to multiagent
486 reinforcement learning. In *31st Conference on Neural Information Processing Systems*, page
487 4193–4206, 2017.
- 488 [39] Joel Z. Leibo, Edgar Duéñez-Guzmán, Alexander Sasha Vezhnevets, John P. Agapiou, Peter
489 Sunehag, Raphael Koster, Jayd Matyas, Charles Beattie, Igor Mordatch, and Thore Graepel.
490 Scalable evaluation of multi-agent reinforcement learning with melting pot. PMLR, 2021.
- 491 [40] Joel Z. Leibo, Vinicius Zambaldi, Marc Lanctot, Janusz Marecki, and Thore Graepel. Multi-
492 agent reinforcement learning in sequential social dilemmas. In *16th International Conference
493 on Autonomous Agents and Multiagent Systems*, 2017.
- 494 [41] Zun Li, Marc Lanctot, Kevin McKee, Luke Marris, Ian Gemp, Daniel Hennes, Paul Muller,
495 Kate Larson, Yoram Bachrach, and Michael P. Wellman. Search-improved game-theoretic
496 multiagent reinforcement learning in general and negotiation games (extended abstract). In
497 *32nd International Conference on Autonomous Agents and Multiagent Systems*, AAMAS, 2023.
- 498 [42] Michael L. Littman. Markov games as a framework for multi-agent reinforcement learning. In
499 *11th International Conference on Machine Learning*, pages 157–163, 1994.
- 500 [43] Stephen McAleer, John Lanier, Kevin Wang, Pierre Baldi, and Roy Fox. XDO: A double oracle
501 algorithm for extensive-form games. In *35th Conference on Neural Information Processing
502 Systems*, 2021.
- 503 [44] Richard D. McKelvey, Andrew M. McLennan, and Theodore L. Turocy. Gambit: Software
504 tools for game theory. <http://www.gambit-project.org/>, 2016.
- 505 [45] H. Brendan McMahan, Geoffrey J Gordon, and Avrim Blum. Planning in the presence of cost
506 functions controlled by an adversary. In *20th International Conference on Machine Learning*,
507 pages 536–543, 2003.
- 508 [46] Richard Mealing and Jonathan L Shapiro. Opponent modeling by expectation–maximization and
509 sequence prediction in simplified poker. In *IEEE Transactions on Computational Intelligence
510 and AI in Games*, volume 9, pages 11–24, 2015.
- 511 [47] Vicent Michalski, Roland Memisevic, and Kishore Konda. Modeling deep temporal depen-
512 dencies with recurrent grammar cells. In *27th Conference on Neural Information Processing
513 Systems*, 2014.
- 514 [48] Johan S. Obando-Ceron and Pablo Samuel Castro. Revisiting rainbow: Promoting more
515 insightful and inclusive deep reinforcement learning research. In *38th International Conference
516 on Machine Learning*, 2021.
- 517 [49] Junhyuk Oh, Xiaoxiao Guo, Honglak Lee, Richard Lewis, and Satinder Singh. Action-
518 conditional video prediction using deep networks in atari games. In *28th Conference on
519 Neural Information Processing Systems*, 2015.
- 520 [50] Junhyuk Oh, Satinder Singh, and Honglak Lee. Value prediction network. In *30th Conference
521 on Neural Information Processing Systems*, pages 6118–6128, 2017.

- 522 [51] S. Phelps, M. Marcinkiewicz, and S. Parsons. A novel method for automatic strategy acquisition
523 in N -player non-zero-sum games. In *Fifth International Joint Conference on Autonomous*
524 *Agents and Multiagent Systems*, page 705–712, 2006.
- 525 [52] Martin L Puterman. *Markov Decision Processes: Discrete Stochastic Dynamic Programming*.
526 John Wiley & Sons, Inc., 1994.
- 527 [53] Julien Pérolat, Joel Z. Leibo, Vinicius Zambaldi, Charles Beattie, Karl Tuyls, and Thore Graepel.
528 A multi-agent reinforcement learning model of common-pool resource appropriation. In *31st*
529 *Conference on Neural Information Processing Systems*, 2017.
- 530 [54] Stéphane Ross and J. Andrew Bagnell. Reinforcement and imitation learning via interactive
531 no-regret learning. *CoRR*, abs/1406.5979, 2014.
- 532 [55] Stéphane Ross, Geoffrey J. Gordon, and J. Andrew Bagnell. A reduction of imitation learning
533 and structured prediction to no-regret online learning. In *14th International Conference on*
534 *Artificial Intelligence and Statistics*, 2011.
- 535 [56] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre,
536 Simon Schmitt, Arthur Guez, Edward Lockhart, Demis Hassabis, Thore Graepel, Timothy
537 Lillicrap, and David Silver. Mastering atari, go, chess and shogi by planning with a learned
538 model. *Nature*, 588:604–609, 2020.
- 539 [57] L. Julian Schvartzman and Michael P. Wellman. Exploring large strategy spaces in empirical
540 game modeling. In *AAMAS-09 Workshop on Agent-Mediated Electronic Commerce*, 2009.
- 541 [58] L. Julian Schvartzman and Michael P. Wellman. Stronger CDA strategies through empiri-
542 cal game-theoretic analysis and reinforcement learning. In *8th International Conference on*
543 *Autonomous Agents and Multi-Agent Systems*, pages 249–256, 2009.
- 544 [59] Ramanan Sekar, Oleh Rybkin, Kostas Daniilidis, Pieter Abbeel, Danijar Hafner, and Deepak
545 Pathak. Planning to explore via self-supervised world models. In *37th International Conference*
546 *of Machine Learning*, pages 8583–8592, 2020.
- 547 [60] David Silver, Aja Huang, Chris J. Maddison, Arthur Guez, Laurent Sifre, George van den Driess-
548 che, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, Sander
549 Dieleman, Dominik Grewe, John Nham, Nal Kalchbrenner, Ilya Sutskever, Timothy Lillicrap,
550 Madeleine Leach, Koray Kavukcuoglu, Thore Graepel, and Demis Hassabis. Mastering the
551 game of Go with deep neural networks and tree search. *Nature*, 529:484–489, 2016.
- 552 [61] David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George van den
553 Driessche, Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al.
554 Mastering the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–
555 489, 2016.
- 556 [62] David Silver, Julian Schrittwieser, Karen Simonyan, Ioannis Antonoglou, Aja Huang, Arthur
557 Guez, Thomas Hubert, Lucas Baker, Matthew Lai, Adrian Bolton, et al. Mastering the game of
558 go without human knowledge. *Nature*, 550(7676):354–359, 2017.
- 559 [63] David Silver, Hado van Hasselt, Matteo Hessel, Tom Schaul, Arthur Guez, Tim Harley, Gabriel
560 Dulac-Arnold, David Reichert, Neil Rabinowitz, Andre Barreto, and Thomas Degris. The
561 predictron: End-to-end learning and planning. In *34th International Conference on Machine*
562 *Learning*, volume 70, pages 3191–3199, 2017.
- 563 [64] Max Olan Smith, Thomas Anthony, Yongzhao Wang, and Michael P. Wellman. Learning to
564 play against any mixture of opponents. *CoRR*, 2020.
- 565 [65] Max Olan Smith, Thomas Anthony, and Michael P. Wellman. Iterative empirical game solving
566 via single policy best response. In *9th International Conference on Learning Representations*,
567 2021.
- 568 [66] Chen Sun, Per Karlsson, Jiajun Wu, Joshua B Tenenbaum, and Kevin Murphy. Stochastic
569 prediction of multi-agent interactions from partial observations. In *7th International Conference*
570 *on Learning Representations*, 2019.

- 571 [67] Richard S. Sutton. Integrated architectures for learning, planning, and reacting based on
572 approximating dynamic programming. In *7th International Workshop on Machine Learning*,
573 pages 216–224. Morgan Kaufmann, 1990.
- 574 [68] Richard S. Sutton. Dyna, an integrated architecture for learning, planning, and reacting. In
575 *SIGART Bulletin*, volume 2, pages 160–163. ACM, 1991.
- 576 [69] Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. The MIT
577 Press, 2018.
- 578 [70] Richard S Sutton, Csaba Szepesvári, Alborz Geramifard, and Michael P. Bowling. Dyna-style
579 planning with linear function approximation and prioritized sweeping. In *28th Conference on*
580 *Uncertainty in Artificial Intelligence*, 2012.
- 581 [71] Erin Talvitie. Model regularization for stable sample rollouts. In *30th Conference on Uncertainty*
582 *in Artificial Intelligence*, 2014.
- 583 [72] Gerald Tesauro. Temporal difference learning and td-gammon. *Communications of the ACM*,
584 38(3):58–68, 1995.
- 585 [73] Zheng Tian, Ying Wen, Zhichen Gong, Faiz Punakkath, Shihao Zou, and Jun Wang. A
586 regularized opponent model with maximum entropy objective. In *International Joint Conference*
587 *on Artificial Intelligence*, 2019.
- 588 [74] Karl Tuyls, Julien Pérolat, Marc Lanctot, Edward Hughes, Richard Everett, Joel Z. Leibo, Csaba
589 Szepesvári, and Thore Graepel. Bounds and dynamics for empirical game theoretic analysis.
590 *Autonomous Agents and Multi-Agent Systems*, 34(7), 2020.
- 591 [75] Yevgeniy Vorobeychik. Probabilistic analysis of simulation-based games. *ACM Transactions*
592 *on Modeling and Computer Simulation*, 20(3), 2010.
- 593 [76] Niklas Wahlström, Thomas B. Schön, and Marc Peter Deisenroth. From pixels to torques:
594 Policy learning with deep dynamical models. *arXiv preprint arXiv:1502.02251*, 2015.
- 595 [77] William Walsh, Rajarshi Das, Gerald Tesauro, and Jeffrey Kephart. Analyzing complex strategic
596 interactions in multi-agent systems. In *AAAI-02 Workshop on Game Theoretic and Decision*
597 *Theoretic Agents*, 2002.
- 598 [78] Rose E Wang, Chase Kew, Dennis Lee, Edward Lee, Brian Andrew Ichter, Tingnan Zhang, Jie
599 Tan, and Aleksandra Faust. Model-based reinforcement learning for decentralized multiagent
600 rendezvous. In *Conference on Robot Learning*, 2020.
- 601 [79] Manuel Watter, Jost Tobias Springenberg, Joschka Boedecker, and Martin Riedmiller. Embed to
602 control: A locally linear latent dynamics model for control from raw images. In *28th Conference*
603 *on Neural Information Processing Systems*, pages 2746–2754, 2015.
- 604 [80] Michael P. Wellman. Methods for empirical game-theoretic analysis. In *21st National Confer-*
605 *ence on Artificial Intelligence*, page 1552–1555, 2006.
- 606 [81] Daniël Willemsen, Mario Coppola, and Guido CHE de Croon. MAMBPO: Sample-efficient
607 multi-robot reinforcement learning using learned world models. In *IEEE/RSJ International*
608 *Conference on Intelligent Robots and Systems*, 2021.
- 609 [82] Ronald J. Williams and David Zipser. A learning algorithm for continually running fully
610 recurrent neural networks. *Neural Computation*, 1(2), 1989.
- 611 [83] Fan Yang, Gabriel Barth-Maron, Piotr Stańczyk, Matthew Hoffman, Siqi Liu, Manuel Kroiss,
612 Aedan Pope, and Alban Rustemi. Launchpad: A programming model for distributed machine
613 learning research. *arXiv preprint arXiv:2106.04516*, 2021.
- 614 [84] Tianhe Yu, Garrett Thomas, Lantao Yu, Stefano Ermon, James Zou, Sergey Levine, Chelsea
615 Finn, and Tengyu Ma. MOPO: Model-based offline policy optimization. In *33rd Conference on*
616 *Neural Information Processing Systems*, 2020.

- 617 [85] Kaiqing Zhang, Sham Kakade, Tamer Basar, and Lin Yang. Model-based multi-agent rl in
618 zero-sum markov games with near-optimal sample complexity. In *33rd Conference on Neural*
619 *Information Processing Systems*, 2020.
- 620 [86] Qizhen Zhang, Chris Lu, Animesh Garg, and Jakob Foerster. Centralized model and exploration
621 policy for multi-agent RL. In *21st International Conference on Autonomous Agents and*
622 *Multiagent Systems*, 2022.