

628 A Broader Impact

629 There are no direct broader impacts from this work. However, this work promotes the adoption
630 of both empirical games and world models. Potential negative impacts may arise due to errors
631 introduced when compressing the true game into the confines of *any* model, which could lead to
632 negative consequences. World model errors within Dyna-PSRO are transferred across response
633 calculations potentially reinforcing biases about the world. If these biases are not rectified, they could
634 negatively influence policies learned from these models. The strategic diversity component of this
635 work aims to mitigate these potential biases, though it represents only the initial step in addressing
636 this concern. When considering empirical games, inaccuracies within them can lead to the suggestion
637 of flawed solutions. The adoption of these inaccurate solutions could have negative repercussions
638 for practitioners or other stakeholders involved in the game. Vigilance and thorough evaluation are
639 required to prevent these potential issues.

640 B Compute

641 GPUs are used for training world models, and policies within Dyna-PSRO. Two types of GPUs were
642 used throughout this work interchangeably: TITAN X and GTX 1080 Ti. All other computation was
643 completed using CPUs. Each response calculation had additional CPUs corresponding to the number
644 of experience generation arenas described in Appendix C. Experiments were run on internal clusters.

645 C Methods Details

646 In this work, the both the policies and world models are implemented in JAX [6] with Haiku [23].
647 The software is architected using Launchpad [85] with design patterns inspired by ACME [28]. All
648 replay buffers are implemented using Reverb [10]. Gambit [46] is used as a game solver via linear
649 complementarity [14].

650 C.1 Policy Implementation & Training

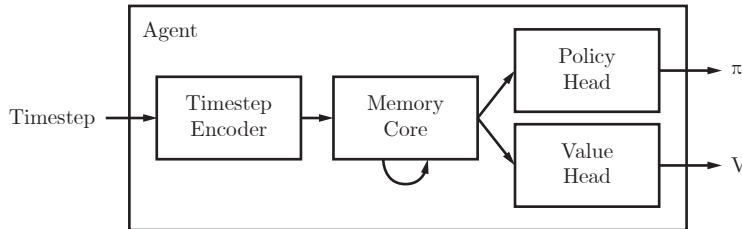


Figure 6: Agent Architecture.

651 All policies follow the general architecture depicted in Figure 6. This consists of four modules:

- 652 • *Timestep Encoder*: Processes all of the current observation’s information into a single
653 embedding vector. The timestep includes the new observation and the policy’s previous
654 action.
- 655 • *Memory Core*: The component of the agent that maintains and update’s the agent’s memory.
- 656 • *Policy Head*: Computes the agent’s policy.
- 657 • *Value Head*: Computes the agent’s state value function.

658 All of the components are simultaneously trained and their joint parameters are $\theta^\pi \in \Theta^\pi$. The
659 policies are trained using the IMPALA algorithm [15]. For the IMPALA loss, the coefficients for
660 each component loss are:

$$\mathcal{L}_{\text{IMPALA}} = \lambda_\pi \cdot \mathcal{L}_\pi + \lambda_V \cdot \mathcal{L}_V + \lambda_{\text{entropy}} \cdot \mathcal{L}_{\text{entropy}},$$

661 with a discount factor of 0.99. The training details for each specific response calculation are itemized
662 below.

663 **Baseline Parameters** The learning rate begins is linearly decayed over 10 000 updates. Each
 664 update is computed from a mini-batch of 128 examples that are generated from 8 arenas² Policy
 665 parameters are synchronized at the beginning of each episode. Each example in the mini-batch is a
 666 sequence of 20 transitions. Moreover, sequences are stored in a replay buffer with a period of 19, to
 667 ensure that the action played at the end of a sequence is trained. Sequences are stored in a replay
 668 buffer with a max capacity of 1 000 000, and are evicted once sampled. Additional hyperparameters
 669 are specified in Table 1.

Table 1: Baseline policy hyperparameters per game.

Hyperparameter	Harvest: Categorical	Harvest: RGB	Running with Scissors
Optimizer	Adam [35]	RMSProp [26]	RMSProp [26]
λ_π	1.0	1.0	1.0
λ_V	0.2	0.5	0.2
λ_{entropy}	0.04	0.01	0.003
Learning Rate Start	6e−6	6e−4	1e−4
Learning Rate Stop	6e−9	6e−9	1e−4
Max Grad Norm	10.0	1.0	0.1
Batch Size	128	128	128

670 Harvest: Categorical module implementations:

- 671 • *Timestep Encoder*: The encoder processes two timestep components: the current observation
 672 and the previous action the policy took. First the observation is passed through a two-layer
 673 fully connected neural network with hidden sizes of [256, 256]. The representation of the
 674 observation is then concatenated with the previous action (represented as a one-hot vector),
 675 and passed together through a second neural network with sizes [256, 256]. All of the
 676 layers have ReLU [17] activations including the final layers of both networks. The final
 677 representation is the output of the timestep encoder.
- 678 • *Memory Core*: A single-layer LSTM [27] with 256 units.
- 679 • *Policy Head*: A single linear layer of size 8.
- 680 • *Value Head*: A single linear layer of size 1.

681 Harvest: RGB and Running with Scissors module implementations:

- 682 • *Timestep Encoder*: The encoder processes two timestep components: the current observation
 683 and the previous action the policy took. The observation is first process by a two-layer
 684 convolutional neural network with ReLU activations [17]. The first layer has 16 channels,
 685 a kernel with shape [8, 8], and a stride of [8, 8]. The second layer has 32 channels, a
 686 kernel shape of [4, 4], and a stride of [1, 1]. The output of this layer is then flattened
 687 and concatenated with a one-hot encoding of the policy’s previous action. The resulting
 688 embedding is then passed through a two-layer fully connected neural network with hidden
 689 sizes of [128, 128], and ReLU activations.
- 690 • *Memory Core*: A single-layer LSTM [27] with 128 units.
- 691 • *Policy Head*: A single linear layer of size 8.
- 692 • *Value Head*: A single linear layer of size 1.

693 **Planning Parameters** The planners have the same hyperparameters as the baseline method, but
 694 with the addition of planning-specific settings. For all planners, an additional 4 arenas are used to
 695 generate planned experiences (for background planning). The additional settings for each version of
 696 planning are as follows:

- 697 • *Warm-Start Background Planning*: An additional 10 000 updates are performed on exclu-
 698 sively planned experiences before play in the real game occurs.

²The term *arena* is used to refer to an experience generation process. This is more commonly referred to as an “actor”; however, this terminology may be confounding with language in RL, Dyna, or multiagent learning.

699 • *Concurrent Background Planning*: Each mini-batch sampled after warm-starting contains
 700 25 % planned experiences, and 75 % real experiences.

701 • *Decision-Time Planning*: In the training arenas (those that have the real game, and are not
 702 used for evaluation), the agent selects actions with a beam-search of width 3 and depth 1.

703 Background planning also requires defining a *search control* procedure [69, 70, 71]. Search control
 704 defines how the agent prioritizes selecting starting states and actions for background planning.
 705 This work considers the simplest search-control method: maintain a buffer of the initial states and
 706 uniformly sample.

707 C.2 World Model Implementation & Training

708 C.2.1 Action-Conditioned Scheduled Sampling

709 As noted by Talvitie [73], rolling out trajectories with an imperfect model tends to result in com-
 710 pounding errors in prediction. Their work suggests training a Markovian world model with previous
 711 predictions (referred to as “hallucinated replay”), to train the model to correct errors. For stateful world
 712 models, as studied in this work, it has been demon-
 713 strated that curricula of n -step future predictions can train a fruitful world model [49, 51, 11]. All of the
 714 preceding work was studying single-agent systems; therefore, they could assume a much more stable
 715 data distribution for training. As a result, these fixed curricula style approaches may prove fatal as
 716 the data distribution may change dramatically throughout training based on the coplayers’ strategies.
 717

Algorithm 1: Action-Conditioned Scheduled Sampling

$m \leftarrow$ Initial recurrent state
for $t \in T$ **do**
 $\mathbf{o} \leftarrow \mathbf{o}^t$ if $\text{Unif}[0, 1] < \epsilon(t)$ else $\hat{\mathbf{o}}^t$
 $\hat{\mathbf{o}}^{t+1}, \hat{\mathbf{r}}^{t+1}, m \leftarrow w(\mathbf{o}, \mathbf{a}^t, m)$

Output: Predicted trajectory $(\hat{\mathbf{o}}^{0:T}, \hat{\mathbf{r}}^{0:T})$

722 Instead, this work adapts the scheduled sampling [5] algorithm as a stochastic curricula, which will
 723 allow both short- and long-term predictions throughout the course of training. Scheduled sampling
 724 is an algorithm for training auto-regressive sequence prediction models where at each predictive
 725 step during training the model input is sampled from either the previous prediction or the ground
 726 truth. Adapting this algorithm for world model rollouts requires biasing each predictive step with
 727 the true actions while sampling between the predicted successor observation and the true successor
 728 observation. Therefore, the predictions will always be biased on true actions, but must learn to handle
 729 model-predicted observation. The sampling follows a schedule $\epsilon : \mathbb{Z} \rightarrow [0, 1]$ that determines the
 730 probability of sampling the true observation over the previous prediction. When ϵ is 1.0, the algorithm
 731 behaves akin to teacher forcing [84] (with the same action-conditional modification); whereas, as it
 732 approaches 0.0 it becomes fully auto-regressive.

733 C.2.2 Implementation

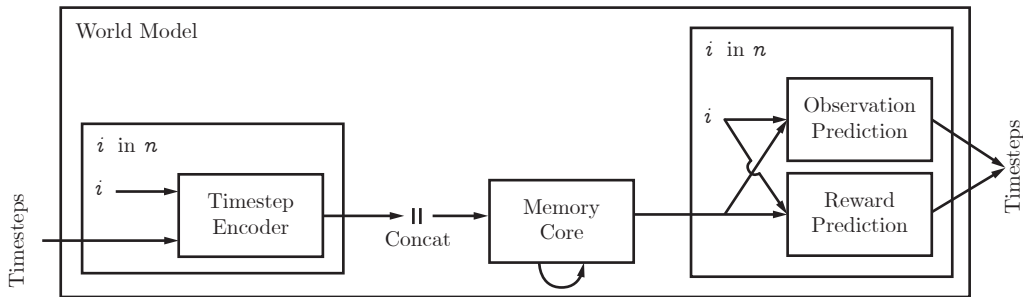


Figure 7: World Model Architecture.

734 The high-level architecture of the world model is illustrated in Figure 7. The world model is composed
 735 of several modules that are quite similar to the policy:

- 736 • *Timestep Encoder*: Processes all of the current observation’s information into a single
737 embedding vector. The timestep includes all new observational data that the agent gains
738 at the current point in time. Different from the agent’s timestep encoder, this encoder also
739 receives the ID that corresponds with the timestep.
- 740 • *Memory Core*: The component of the agent that maintains and update’s the agent’s memory.
741 Different from the agent’s timestep encoder, this memory core receives the representation of
742 each player’s timestep concatenated.
- 743 • *Observation Prediction (Head)*: Predicts the successor observation for each player. As
744 all games considered in this work are gridworld games, the predicted observation is a
745 classification task for each future grid cell (that are within the respective player’s observation
746 window).
- 747 • *Reward Prediction (Head)*: Predicts the reward received for each player. Rewards are treated
748 as categorical values.

749 Note, that the timestep encoder, observation prediction head, and reward prediction head each use
750 the same parameters across each player. Similar to the agent, all components are simultaneously
751 trained and their joint parameters are referred to as $\theta^w \in \Theta^w$. Both observation and reward losses are
752 optimized with a cross entropy objective, and averaged across players. The total world model loss is
753 as follows:

$$\mathcal{L}_w = \lambda_{\text{observation}} \cdot \mathcal{L}_{\text{observation}} + \lambda_{\text{reward}} \cdot \mathcal{L}_{\text{reward}}.$$

754 The implementation of each component is as follows:

- 755 • *Timestep Encoder*: The same as the agent’s timestep encoder, but the player’s ID is also
756 provided alongside the action into the second neural network.
- 757 • *Memory Core*: Identical to the agent.
- 758 • *Observation Prediction (Head)*: The observation prediction is based on the memory core’s
759 output and a one-hot ID of the predicted player’s ID. These inputs are concatenated and fed
760 into an transposed version of the timestep encoder.
- 761 • *Reward Prediction (Head)*: A linear layer of size one. For Harvest: Categorical this output
762 is handled as a discrete prediction; whereas, it is continuous for the other games.

763 A world model is trained for 1 250 000 updates. Each example in the mini-batch is a sequence of 20
764 transitions, where the first 5 timesteps are used to burn-in the memory. Burn-in does not occur for
765 examples where the first 5 transitions are at the beginning of the episode. Moreover, sequences are
766 added into the replay buffer at a period of 14 so that all timesteps show up as prediction targets.

767 The world model is trained using action-conditioned scheduled sampling (Appendix 1, Algorithm 1).
768 The schedule ϵ follows the following schedule:

$$\epsilon(t) = \begin{cases} 1.0 & t < 250000 \\ \frac{4}{3} - \frac{t}{750000} & 250000 \leq t \leq 1000000 \\ 0.0 & t > 1000000. \end{cases}$$

769 This schedule starts out training as a variation of teacher forcing [84], and slowly transitions to fully
770 auto-regressive. Additional hyperparameters are specified in Table 2.

Table 2: World model hyperparameters per game.

Hyperparameter	Harvest: Categorical	Harvest: RGB	Running with Scissors
$\lambda_{\text{observation}}$	1.0	1.0	1.0
λ_{reward}	10.0	0.01	0.01
Optimizer	Adam [35]	Adam [35]	Adam [35]
Learning Rate	3e−4	3e−4	3e−4
Max Grad Norm	10.0	10.0	10.0
Batch Size	32	24	24

771 **C.3 Strategic Diversity**

772 Learning a general world model assumes that the transitions are drawn from the space of all possible
 773 transitions. This is typically not tractable, but instead draws are taken from a dataset generated from
 774 play of a *behavioral strategy* σ . And the performance of the world model is measured under a *target*
 775 *strategy* σ^* , instead of all possible strategies Σ . Differences between σ and σ^* present challenges in
 776 learning an effective world model.

777 We call the probability of drawing a state-action pair s, a under some strategy its *reach probability*
 778 $\eta^{\hat{\sigma}}$ under joint strategy $\hat{\sigma}$. From this, we define *strategic diversity* as the distribution induced from
 779 reach probabilities. These terms allow us to observe two challenges for learning world models.

780 First, the diversity of the behavioral strategy *cover* the target strategy’s diversity:

$$\eta^{\sigma^*}(s, a) \rightarrow \eta^{\sigma}(s, a). \quad (1)$$

781 Otherwise, transitions will be absent from the training data. As an aside, it is possible to construct a
 782 weaker claim for coverage. This is done through making additional assumptions about the generaliza-
 783 tion capacity of a world model across transitions. For example, if transitions are drawn from two
 784 discrete latent variables, unseen combinations of these variables may be generalized if the individual
 785 values are known. However, generalization cannot be generally guaranteed, so we consider coverage.

786 The second challenge is that the *closer* the diversities are, the more accurate the learning objective
 787 will be. In other words, we want

$$\eta^{\sigma^*}(s, a) \approx \eta^{\sigma}(s, a). \quad (2)$$

788 If closeness is not ensured, crucial dynamics knowledge may not be learned as the learning signal is
 789 dominated from unimportant transitions. An example of the issue of closeness can be seen in the
 790 “noisy TV problem,” [9]. This exploration problem poses that novelty-seeking agents may be stuck
 791 forever watching the ever new TV static, and not experiencing practical novelty. In the same vein,
 792 if a world model is trained almost entirely on “noisy TV”-like experiences, and as a rarely on the
 793 few salient experiences, it may never learn. Therefore, we should strive to correct the distribution of
 794 experiences to be informed by a target strategy.

795 By design, empirical-game building algorithms offer a means to construct the target world model
 796 objective. These algorithms require the specification of a solution concept that serves the dual roll
 797 as the target strategy for a world model. Then through an iterative process, the empirical-game
 798 constructs strategies that progressively approach the target. In turn, generating transitions that match
 799 the target world model objective.

800 **Claim 1.** *Dyna-PSRO produces a correct world-model objective η^{σ^*} with a best-response oracle*
 801 *and a correct empirical game for a game with a unique Nash Equilibrium σ^* .*

802 *Proof.* Following McMahan et al. [47], the Double Oracle algorithm will converge to a NE in the
 803 limit of enumerating the full strategy space. Let $\sigma^0, \sigma^1, \dots, \sigma^e$ be the solutions discovered for each
 804 epoch, ending at epoch e . Then a dataset composed of experiences generated by the current empirical
 805 game solution evolves as follows:

$$\eta^{\sigma^0} \rightarrow \eta^{\sigma^1} \rightarrow \dots \rightarrow \eta^{\sigma^e} = \eta^{\sigma^*}. \quad (3)$$

806 □

807 The previous claim contains two strong assumptions: an exact best-response oracle and error-less
 808 empirical game. These assumptions must be made, because PSRO is parameterized by its choice
 809 of response oracle and empirical game model; therefore, PSRO’s convergence must be proven for
 810 each choice. Theoretically PSRO has been shown to converge to an ϵ -NE, where ϵ depends on
 811 the empirical game’s modelling error, to a corresponding NE in the true game [76, 77]. Therefore,
 812 in practice Dyna-PSRO produces $\eta^{\sigma^e} \approx \eta^{\sigma^*}$, which supports the weaker claim that Dyna-PSRO
 813 generally improves the quality of a world model.

814 It is also worth noting the connections between this analysis and MARL regimes that seek to find any
 815 solution the game. In these regimes, the priority is finding *any* performant strategy. This matches the
 816 approach taken by the majority of studies in MARL falling under paradigms such as Independent RL
 817 or Self-Play. Therefore, their target distribution is the best-response to the previous strategy $\eta^{\text{BR}(\sigma^{i-1})}$
 818 and changes in tandem with the strategies. When no best-response can be found, then the current
 819 strategy matches the solution and the dataset correspondingly reflects this.

820 **C.4 Dyna-PSRO**

821 The Dyna-PSRO builds upon PSRO (Algorithm 3) by including the co-learning of a world model.
 822 The high-level pseudocode of Dyna-PSRO is provided in Algorithm 5 and a high-level application
 823 architecture diagram is depicted in Figure 8. There are three main co-routines of Dyna-PSRO:
 824 response computation, world-model learning, and empirical-game simulation. The details of each
 825 routine are first provided; then, how the routines interact with each other is explained.

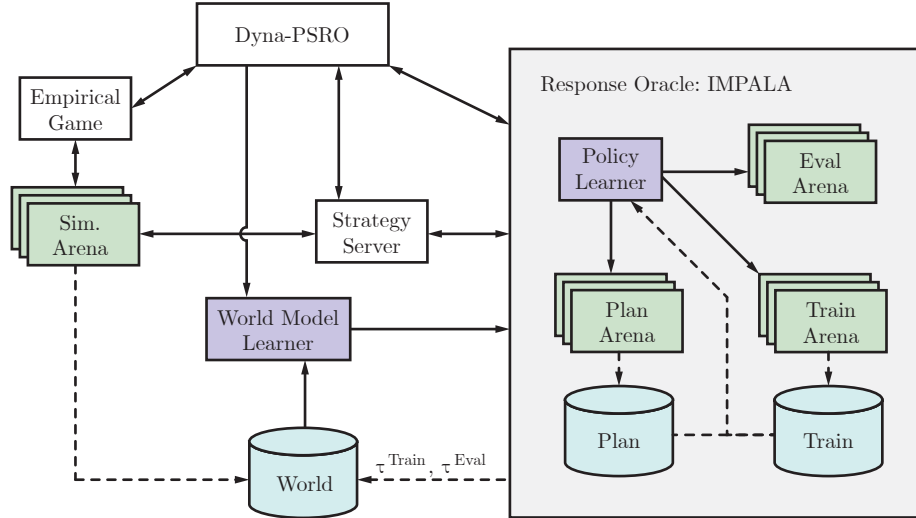


Figure 8: Overview of the major Dyna-PSRO processes.

826 **C.4.1 Empirical Game**

827 The empirical game routine is responsible for maintaining the empirical game, including simulating
 828 new payoffs and game reasoning. New profiles are sent to *simulation (sim.) arenas* for payoff
 829 estimation in parallel. Once all profiles are estimated, the game is solved, and the solution is based
 830 to the main Dyna-PSRO process. In the experiments in this work, the chosen solution is Nash
 831 Equilibrium, and it is solved through the linear complementarity [14] algorithm that is implemented
 832 by Gambit [46].

833 **C.4.2 World Model**

834 The world model routine is responsible for training the world model and serving its parameters. This
 835 routine’s pseudocode is provided in Algorithm 2, and follows mostly the same method details as
 836 the strategic diversity experiment. The difference is that instead of there being a precomputed fixed
 837 dataset, the world model is now trained over a dynamic dataset. The dataset is represented by a
 838 replay buffer that is populated from: (1) trajectories from the simulation arena used for expanding
 839 the empirical game, and (2) trajectories from the training and evaluation arenas from the response
 840 calculation. Notably, all of this data must be generated in the standard PSRO procedure, so it collected
 841 with no additional cost. The world-model learner samples and evicts data randomly from this buffer.

Algorithm 2: World Model Learner

Input: World model w and data buffer \mathcal{B}^w

Input: n no. of updates (default: ∞).

842 **for** $i \in [[n]]$ **do**
 Train w over $\tau \sim \mathcal{B}^w$

Output: w

843 C.4.3 Response Oracle

844 The response oracle uses the IMPALA [15] algorithm to compute an approximate best-response to
 845 the opponent’s strategy according to the current empirical game. IMPALA uses several processes
 846 that generate experiences for the agent to train on. These process are referred to in this work as arenas.
 847 The *train arenas* generate real experiences, and the *plan arenas* generate planned experiences. If the
 848 learner is using decision-time planning they will only use it in the train arenas. A third set of arenas
 849 called *eval arenas* periodically evaluate the performance of the greedy policy and record additional
 850 metrics. The arenas attempt to synchronize all parameters at the start of each episode.

851 The policy learner runs for a fixed number of updates, querying the datastores for experiences to learn
 852 from. The specifics of how each policy learns is described in Appendix C.1.

853 C.4.4 Runtime Procedure

854 A sketch of the respective processes runtime is shown in Figure 9 As in PSRO, the main empirical-
 855 game building loop iterates between response computation and empirical-game simulation.

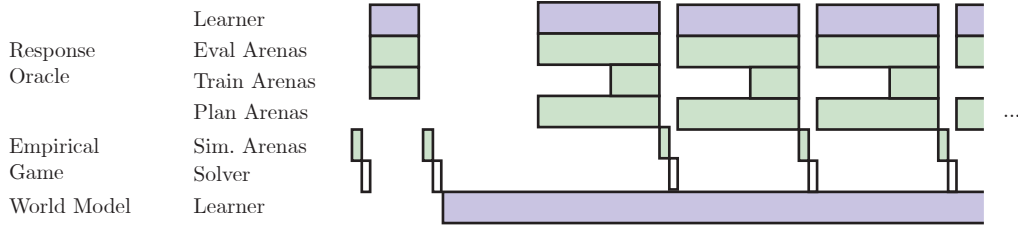


Figure 9: Example Dyna-PSRO runtime. Planning is set to occur after the first epoch. Each players’ response oracle runs in parallel.

856 The runtime is defined by a parameter specifying on which PSRO epoch to begin planning. Before
 857 that epoch, the response oracles do not use planning at all, because the world models are untrained.
 858 All of these policies therefore are trained exclusively on real experiences just like standard PSRO.
 859 However, these experiences are also being used to populate the world model’s replay buffer. Once the
 860 first planning epoch has arrived, computing responses is temporarily paused. The world model is then
 861 given a set number of updates to warm-start its parameters, before being used in response calculation.
 862 Once the world model’s warm-start phase is over, all process proceed concurrently.

863 Throughout this work planning begins on the second epoch. The world models is given 1 million
 864 updates of warm starting.

Algorithm 3: Policy-Space Response Oracles [40]

Input: Initial strategy sets for all players Π^0
 Simulate utilities \hat{U}^{Π^0} for each joint $\pi^0 \in \Pi^0$
 Initialize solution $\sigma_i^{*,0} = \text{Uniform}(\Pi_i^0)$
while epoch e in $\{1, 2, \dots\}$ **do**
 for player $i \in [[n]]$ **do**
 // Algorithm 4.
 $\pi_{i,-}^e = \text{response_oracle}(\sigma_{-i}^{*,e-1})$
 $\Pi_i^e = \Pi_i^{e-1} \cup \{\pi_i^e\}$
 Simulate missing entries in \hat{U}^{Π^e} from Π^e
 Compute a solution $\sigma_i^{*,e}$ from $\hat{\Gamma}^e$
Output: Current solution $\sigma_i^{*,e}$ for player i

Algorithm 4: Response Oracle

Input: Coplayer strategy profile σ_{-i}
Input: Num updates k
 $\pi_i \leftarrow \theta^\pi$
 $\mathcal{B} \leftarrow \{\}$ // Replay Buffer.
for many async episodes **do**
 $\pi_{-i} \sim \sigma_{-i}$
 $\mathcal{B} = \mathcal{B} \cup \{\tau \sim (\pi_i, \pi_{-i})\}$
for $i \in [[k]]$ **do**
 Train π_i over $\tau \sim \mathcal{B}$
Output: π_i, \mathcal{B}

Algorithm 5: Dyna-PSRO

Input: Initial strategy sets for all players Π^0 **Input:** No. of world model head-start updates n_w **Input:** Epoch to begin planning e^{plan} Simulate utilities \hat{U}^{Π^0} for each joint $\pi^0 \in \Pi^0$ Initialize solution $\sigma_i^{*,0} = \text{Uniform}(\Pi_i^0)$ $w \leftarrow \theta^w$ $\mathcal{B}^w \leftarrow \{\}$

// World Model's Replay Buffer.

while epoch e in $\{1, 2, \dots\}$ **do** **for** player $i \in [[n]]$ **do** **if** $e > e^{\text{plan}}$ **then** $\pi_i^e, \tau = \text{async}(\text{planner_oracle}(\sigma_{-i}^{*,e-1}, w))$ // Algorithm 6. **else** $\pi_i^e, \tau = \text{async}(\text{response_oracle}(\sigma_{-i}^{*,e-1}))$ // Algorithm 4. $\mathcal{B}^w = \mathcal{B}^w \cup \{\tau\}$ $\Pi_i^e = \Pi_i^{e-1} \cup \{\pi_i^e\}$ Wait on all futures π^e, τ Simulate missing entries in \hat{U}^{Π^e} from Π^e Add τ from simulation to \mathcal{B}^w Compute a solution $\sigma^{*,e}$ from $\hat{\Gamma}^e$ **if** $e == l$ **then** $w = \text{world_model_learner}(w, n_w)$ // Algorithm 2. $w = \text{async}(\text{world_model_learner}(w))$ // Parameters periodically sync.**Output:** Current solution $\sigma_i^{*,e}$ for player i

Algorithm 6: Planner Oracle

Input: Coplayer strategy profile σ_{-i} **Input:** World model w , real game dynamics p **Input:** Warm-start background planning updates $n^{\text{BG:WS}}$ **Input:** Training updates n **Input:** Concurrent background planning fraction $f^{\text{BG:C}}$ $\pi_i \leftarrow \theta^\pi$ $\mathcal{B}^{\text{plan}} \leftarrow \{\}$

// Replay Buffer with planned experience.

 $\mathcal{B}^{\text{train}} \leftarrow \{\}$

// Replay Buffer with real experience.

// Asynchronously generate data on arenas.

for many async episodes **do**

866

 $\pi_{-i} \sim \sigma_{-i}$ $\mathcal{B}^{\text{plan}} = \mathcal{B}^{\text{plan}} \cup \{\tau \sim (\pi_i, \pi_{-i}, w)\}$ **for** many async episodes **do** $\pi_{-i} \sim \sigma_{-i}$ $\mathcal{B}^{\text{train}} = \mathcal{B}^{\text{train}} \cup \{\tau \sim (\pi_i, \pi_{-i}, p)\}$

// Train the response policy.

for $i \in [[n^{\text{BG:WS}}]]$ **do** Train π_i over $\tau \sim \mathcal{B}^{\text{plan}}$ **for** $i \in [[n]]$ **do** Train π_i over $\tau \sim \{(1.0 - f^{\text{BG:C}}) \cdot \mathcal{B}^{\text{train}}\} \cup \{f^{\text{BG:C}} \cdot \mathcal{B}^{\text{plan}}\}$ **Output:** π_i, \mathcal{B}

867 **C.5 Combined-Game Regret**

868 Combined-game regret is an approximate measure of regret that all available estimates to approximate
 869 the regret within the true game. Intuitively, combined-game regret is the regret of a strategy with
 870 respect to all discovered policies. When comparing empirical-game building algorithms this is
 871 formalized as follows:

$$\text{SumRegret}(\sigma, \bar{\Pi}) = \sum_{i \in \mathcal{N}} \max_{\pi_i \in \bar{\Pi}_i} \hat{U}_i(\pi_i, \sigma_{-i}) - \hat{U}_i(\sigma_i, \sigma_{-i}), \quad \bar{\Pi}_i \equiv \bigcup_{\text{method}} \hat{\Pi}_i^{\text{method}}, \quad (4)$$

872 where $\hat{\Pi}$ is the restricted strategy set from one of the algorithms.

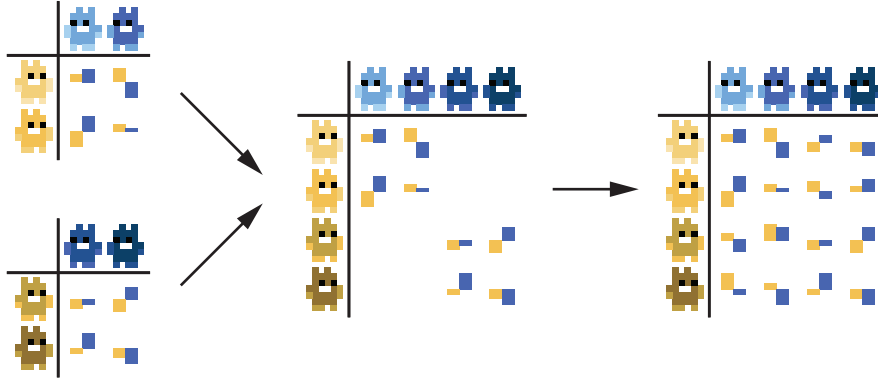


Figure 10: Combined-game construction. Left: Constituent empirical games. Middle: Combination of the strategy sets and payoff functions. Right: Completion of the empirical game by estimating new strategy profile payoffs.

873 The process of constructing a combined-game is illustrated in Figure 10. Where, the strategy sets
 874 (depicted by the toons) across methods are first combined. The new *combined game* that results from
 875 this can be initialized with the payoff estimates from the constituent empirical games. Unestimated
 876 payoffs must then be simulated for the new strategy profiles. Then the complete combined game
 877 can be used to compute the combined-game regret from the solutions computed throughout the
 878 empirical-game building algorithms.

879 **D Games**

880 **D.1 Harvest: Categorical**

881 In Harvest, players move around an orchard picking apples. The challenging commons element is that
882 apple regrowth rate is proportional to nearby apples, so that socially optimum behavior would entail
883 managed harvesting. Self-interested agents capture only part of the benefit of optimal growth, thus
884 non-cooperative equilibria tend to exhibit collective over-harvesting. The game has established roots
885 in human-behavioral studies [32] and in agent-based modeling of emergent behavior [55, 42, 41].

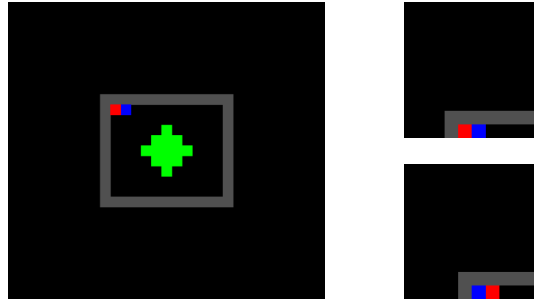


Figure 11: Harvest: Categorical. Left: game state. Right: player observations.

886 For our initial experiments, we use a symmetric two-player version of the game, where in-game
887 entities are represented categorically [30]. This categorical representation facilitates faster experi-
888 mentation and simplifies the interpretation of results. Figure 11 depicts the game state and player
889 observations. Each player has a 10×10 viewbox within their field of vision. The cells of the grid
890 world can be occupied by either agent shown in red and blue, the apples shown in green, or a wall in
891 gray. The possible actions include moving in the four cardinal directions, rotating either way, tagging,
892 or remaining idle. A successful tag temporarily removes the other player from the game, but can only
893 be done to other nearby players. Players receive a reward of 1 for each apple picked. Episodes are
894 limited to 100 timesteps.

895 **D.2 Harvest: RGB**

896 Harvest: RGB is a different implementation of the Harvest game introduced by Harvest: Categorical
897 (Appendix D.1. Harvest: RGB is exactly the harvest implementation from MeltingPot [41] with the
898 same orchard map. A rendering of the game state and observations is shown in Figure 12. The main
899 difference between the Harvest versions is that the observations are $88 \times 88 \times 3$ images of the 11×11
900 viewbox in front of them. There are also minor differences in the implementation of tagging and
901 apple respawn mechanism. Episodes play for 1000 timesteps.

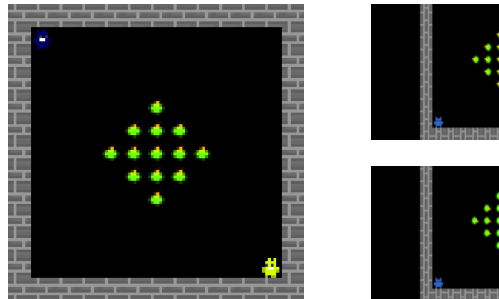


Figure 12: Harvest: RGB. Left: game state. Right: player observations.

902 **D.3 Running With Scissors**

903 Running With Scissors (RWS) is a temporally extended version of rock-paper-scissors (RPS). In it,
904 players collect rock, paper, and scissor items into their inventory. At any point the player has the
905 option to tag their opponent if they're nearby. Then they play RPS corresponding to the distribution
906 of items in their inventories. The agents have the same action space as in the previous games. The
907 observation space is $40 \times 40 \times 3$ image-based viewbox in front of them corresponding to a 6×6 grid
908 around them. A portion of items are placed within the game deterministically, the rest are randomly
909 sampled before play. If neither player tags each other before 1000 timesteps, the players are forced
910 into playing RPS.

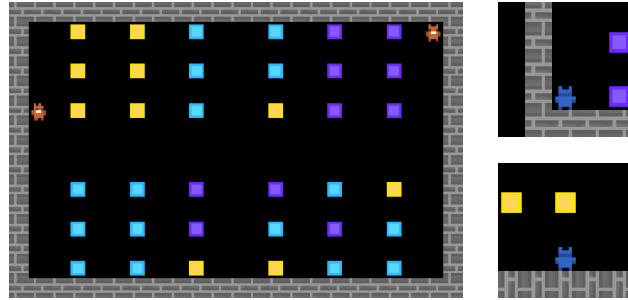


Figure 13: Running With Scissors. Left: game state. Right: player observations.

911 **E Additional Results**

912 **E.1 Strategic Diversity**

913 Figure 14 displays the recall results that correspond to the accuracies portrayed in Figure 2. See
 914 Section 3.1 for a discussion of these results.

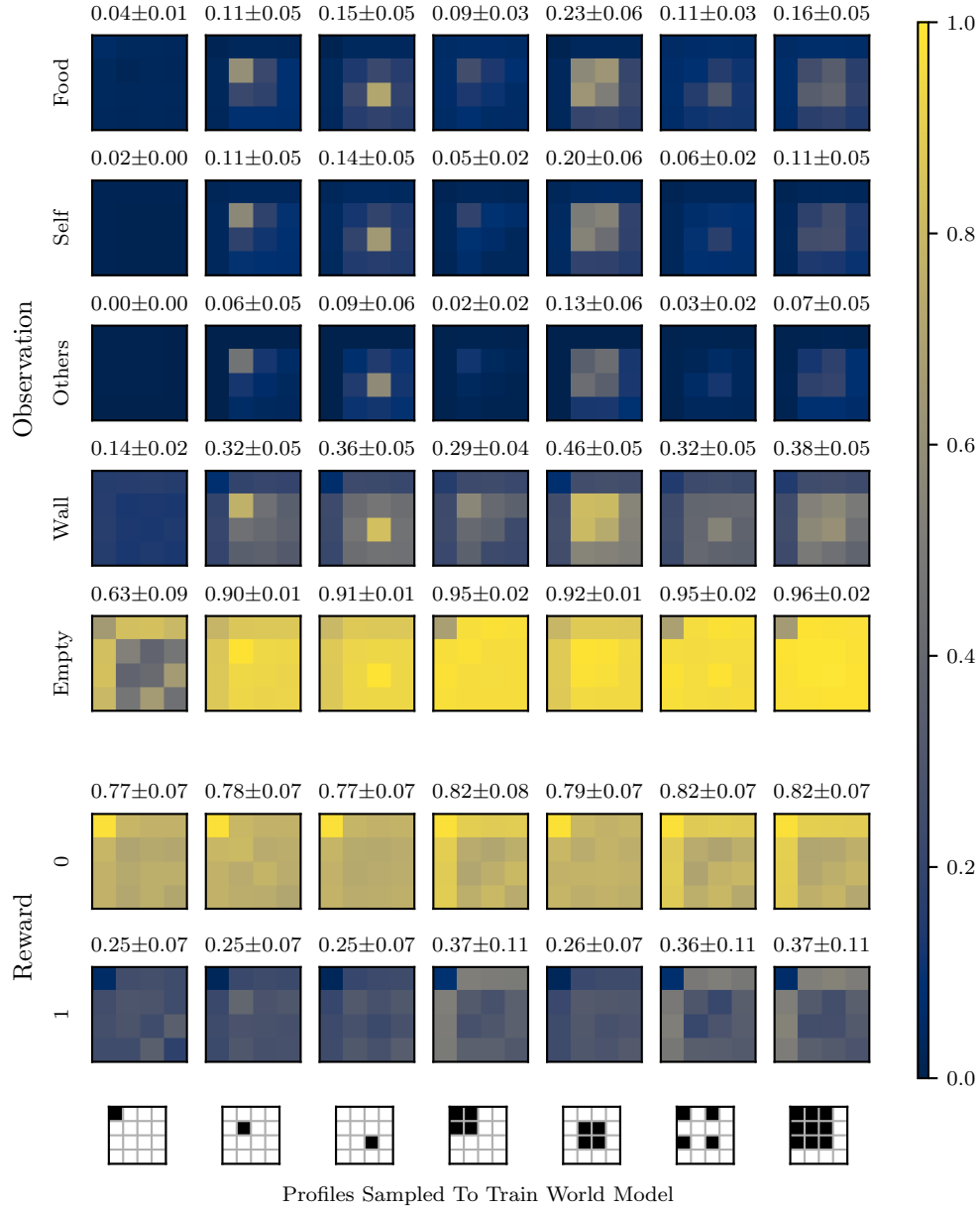


Figure 14: World model recall on Harvest: Categorical.

915 **E.2 Background Planning**

916 Figure 15 shows the results of repeating the background planning experiment with world model \boxtimes .
 917 Besides changing the world model, the methodology is consistent with that described in Section 3.2.1.
 918 This result shows the planner achieving results comparable to the baseline method. Supporting the result
 919 the adoption of planning, as it tends to not negatively impact the learning process.

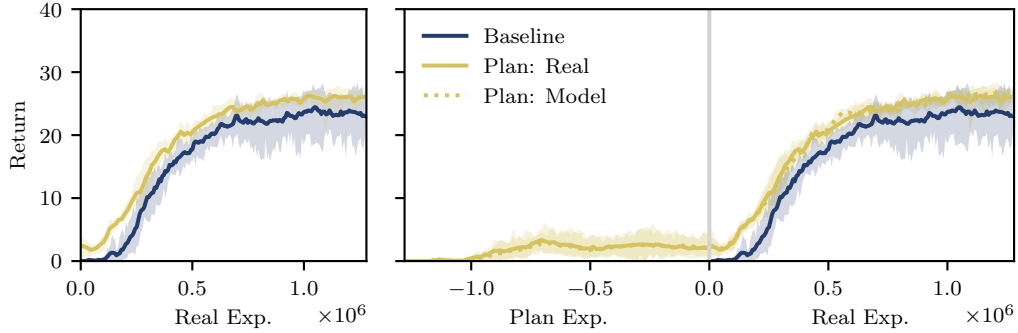


Figure 15: Effects of background planning on response computation using world model \boxtimes . (5 seeds, with 95 % bootstrapped CI).

920 **E.3 Decision-Time Planning**

921 Figure 16 shows the results of repeating the decision-time planning experiment with world model \boxtimes .
 922 Besides changing the world model, the methodology is consistent with that described in Section 3.2.2.
 923 This result further exemplifies the trend shown in Figure 4, where the planners that did not use BG: W
 924 failed to learn an effective policy. The planner that used BG: W achieved performance comparable to
 925 the baseline. Finally, the planner that used both BG: W and BG: C achieves the strongest performance
 926 at 33.07 ± 6.76 . These results support the benefit of BG: W when using DT, and that effective
 927 planning performs as least as good as the baseline.

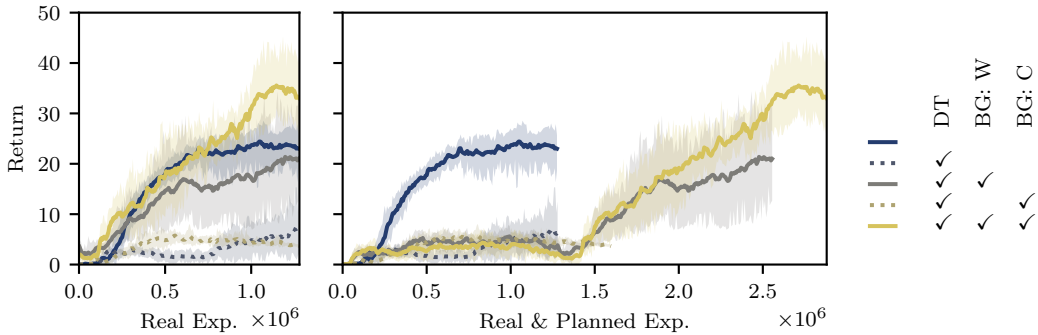


Figure 16: Effects of decision-time planning on response computation using world model \boxtimes . (5 seeds, with 95 % bootstrapped CI).

928 **E.4 World Models as Empirical Games**

929 In this section, we verify the need for separate models.

930 First, consider if an empirical game can substitute for a world model. The majority of previous work
 931 on empirical games represents the model in the normal form. This representation abstracts away any
 932 notion of dynamics within an episode into a choice in policy and the resulting payoff. Since empirical
 933 games currently lack dynamics information completely, this supports the choice of separate models.

934 This is not without any exceptions. If the original game is one-shot and stateless (i.e., an episode is
 935 played through a single action), then a normal-form empirical game is exactly a world model.

936 Now, consider if a world model can substitute for an empirical game. World models predict successor
 937 states and rewards; and thus, can rollout planned trajectories to estimate payoffs. Note, that rolling
 938 out a trajectory a trajectory with a world model is an auto-regressive prediction that tends to result
 939 in compounding errors [73, 29]. Despite this, it is plausible that a world model can substitute as a
 940 high-fidelity empirical game.

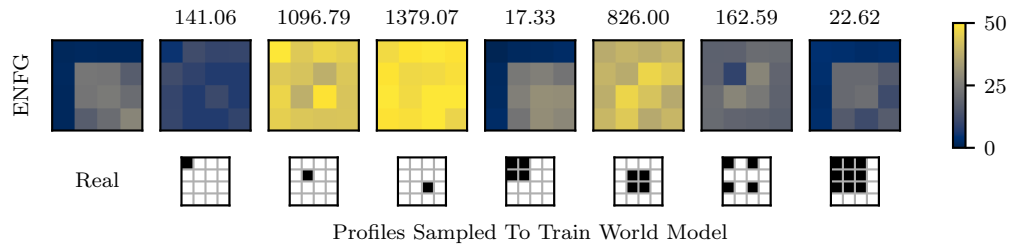


Figure 17: Empirical normal-form games (ENFG) estimated by world model rollouts. The title of each plot is its L2 distance with the real ENFG.

941 Figure 17 compares an empirical game estimated from real game payouts empirical games estimated
 942 with payouts predicted by a world model. In this experiment, the world models are the same that were
 943 used in Section 3.1. In general, the empirical games estimated by world models have large errors
 944 ($L2 > 100$), with several having exceptionally large errors ($L2 > 1000$). These result suggest that this
 945 direction may be possible with future algorithmic improvements; however, currently, the prediction
 946 errors are too large to substitute empirical games with world models. Especially in games with long
 947 time horizons.