CODE WORLD MODELS FOR GENERAL GAME PLAYING

Anonymous authors
Paper under double-blind review

000

001

002 003 004

005

006

008 009

010 011

012

013

014

015

016

017

018

019

021

022

025

026

027

028

029

031

032

033

034

036

038

040

042

043

044

046

048

049

051

052

Abstract

Large Language Models (LLMs) reasoning abilities are increasingly being applied to classical board and card games, but the dominant approachinvolving prompting for direct move generation—has significant drawbacks. It relies on the model's implicit fragile pattern-matching capabilities, leading to frequent illegal moves and strategically shallow play. Here we introduce an alternative approach: We use the LLM to translate natural language rules and game trajectories into a formal, executable world model represented as Python code. This generated model—comprising functions for state transition, legal move enumeration, and termination checks—serves as a verifiable simulation engine for high-performance planning algorithms like Monte Carlo tree search (MCTS). In addition, we prompt the LLM to generate heuristic value functions (to make MCTS more efficient), and inference functions (to estimate hidden states in imperfect information games). Our method offers three distinct advantages compared to directly using the LLM as a policy: (1) Verifiability: The generated CWM serves as a formal specification of the game's rules, allowing planners to algorithmically enumerate valid actions and avoid illegal moves, contingent on the correctness of the synthesized model; (2) Strategic Depth: We combine LLM semantic understanding with the deep search power of classical planners; and (3) Generalization: We direct the LLM to focus on the meta-task of datato-code translation, enabling it to adapt to new games more easily. We evaluate our agent on 10 different games, of which 4 are novel and created for this paper. 5 of the games are fully observed (perfect information), and 5 are partially observed (imperfect information). We find that our method outperforms or matches Gemini 2.5 Pro in 9 out of the 10 considered games.

1 Introduction

Large Language Models (LLMs) have shown impressive abilities at solving various reasoning tasks, and recently have been applied as "agents" which can play classical (often multiplayer) games, like Chess, Go, and even complex imperfect information games like Poker and Bridge. The standard approach is to treat the LLM as a policy, by asking it to pick a move at each step using a prompting strategy based on the trajectory of observations and actions seen so far, plus optional text meta data about the game. This method treats the LLM as an end-to-end "intuitive player", leveraging its vast training data to recognize patterns and select moves that seem promising. However, strategic mastery often requires deep multi-step lookahead, characteristic of a "System 2" deliberation (Kahneman, 2003). While strong play can be achieved through training specialist models (Ruoss et al., 2024; Schultz et al., 2025), direct play from generalist LLMs often lacks deep tactical foresight, despite recent advances in "thinking" (Liao et al., 2025), as we show empirically in this paper. In addition, the LLM as policy approach does not work very well on novel games that are not part of the LLM's training set, as we will also show.

We propose to use LLMs in a different way, namely as induction engines that can leverage their prior knowledge to map a small amount of observed trajectory (game play) data, plus a textual game description, into plausible world models, represented as Python code, using iterative code refinement methods as in Tang et al. (2024a). We call the result of this pro-

cess a "Code World Model" (CWM). In the context of game playing, a CWM consists of a definition of the (possibly latent) state, a function that specifies which moves are legal at each step, a state transition function, an observation function (for latent states), a reward function, and a function that checks for termination. Furthermore, for the challenging case of partially observable games, we introduce a novel paradigm that effectively tasks the LLM with synthesizing a regularized autoencoder: an inference function (the encoder) maps observations to plausible latent histories, and the CWM (the decoder) reconstructs observations from them, with the game's rules and API serving as a strong structural regularizer.

Although there is prior work that uses LLMs to learn symbolic world models (see Sec. 3), and then leverage them for planning, we differ in three main ways. First, we handle the case of partially observed and stochastic worlds (such as Poker), whereas all prior work (to the best of our knowledge) either assumes fully observed and deterministic environments, or (in the case of Curtis et al. (2025)) assumes post-hoc observability; both cases make model learning much easier. Second, in addition to learning a CWM, we ask the LLM to generate heuristic value functions, which significantly improves the performance of our search-based policies, such as MCTS and Information Set MCTS (Cowling et al., 2012). Third, we demonstrate that our approach outperforms a state-of-the-art "thinking" LLM across various two-player games, including novel (or "OOD") ones which we create, to avoid contamination issues with the training set of the LLM.

2 Background

Interactions in multiplayer games can be described using the formalism of extensive-form games (Kuhn, 1953; Shoham & Leyton-Brown, 2009; Albrecht et al., 2024; Murphy, 2025): there is a set $\mathcal{N} = \{1, 2, \dots, n\}$ of n players that take discrete actions $a \in \mathcal{A}$. Sequences of actions are called histories $h \in \mathcal{H}$; all games start at the initial empty history, and end at terminal histories $\mathcal{Z} \subseteq \mathcal{H}$. There is a special player called *chance* (also sometimes called nature), c, which plays with a known, fixed (stochastic) policy—the chance outcome distribution—e.q., representing dice rolls and card draws. Due to chance events being explicitly represented by the game environment, each history h can be thought of as a unique transcription of a game (either finished or in progress) and as a "ground truth" state known only to the environment. At every history h, there is a player to act $\tau(h) \in \mathcal{N} \cup \{c\}$, and a set of legal actions $\mathcal{A}(h) \subseteq \mathcal{A}$. Formally defining states in partially-observable (imperfect information) games can be tricky, and we defer this to Appendix B to couple it with the description of the search method (policy generation). Agents encode policies to take actions $\pi(h) \in \Delta(\mathcal{A})$, where $\Delta(\cdot)$ represents a discrete probability distribution. For each agent i, the goal is to find a policy that maximize its own cumulative reward $\sum_{t=1}^{T} r^{i}(h_{t})$. However, in the multiagent setting each individual objective jointly depends on choices of other agents.

Our game environments are based on OpenSpiel (Lanctot et al., 2019): each implementation provides logic to determine legal actions, transitions from one ground truth state to the next, rewards, and player observations in a general way. However, the agent does not know the true environment model. Instead, it must learn the code world model by using an LLM applied to a text description of the game, together with example game play data, as described in detail in Sec. 4. Given the learned CWM, we pick the best move by using existing game solvers: for perfect information games, we use MCTS, and for imperfect information games, we use Information Set MCTS (see Appendix B). In both cases, we optionally augment the search algorithm with a learned value function, and in the case of ISMCTS, we augment the search algorithm with a hidden state estimator. We also tried learning a policy using PPO applied to the (partially observed) CWM: see Appendix D for details.

3 Related Work

There is a growing interest in evaluating the abilities of LLMs to play games, as exemplified by the recent release of Kaggle Game Arena¹, as well as other recent work (Costarelli et al., 2024; Duan et al., 2024; Verma et al., 2025; Hu et al., 2025a; Sun et al., 2025;

¹See https://www.kaggle.com/game-arena.

Cipolina-Kun et al., 2025; Hu et al., 2025b; Guertler et al., 2025). Similar to these papers, our aim is to design LLM-based agents that play text-based games. Furthermore, like ggbench (Verma et al., 2025), we assess the generality of our agents using novel games, that are (by construction) out-of-distribution (OOD) for the LLM. However, rather than using the LLM directly as a policy, we focus on using the LLM to generate a CWM, to which we then apply standard solvers, such as (IS)MCTS or PPO.

There are a few other papers that also use a model-based approach, similar to ours. "World-Coder" generates a set of CWM hypotheses from trajectory data using LLM-powered code synthesis, stores each hypothesis (candidate model) in a tree, and uses Thompson sampling to decide which hypothesis to ask the LLM to improve, see (Tang et al., 2024a). Given the learned CWM, WorldCoder uses ReAct-style methods (Yao et al., 2022) for decision-making. GIF-MCTS (Dainese et al., 2024) developed a similar method, but uses MCTS for agent decision-making. Our work extends this past work by considering strategic multiagent environments, synthesizing value functions (to speed up (IS)MCTS), and synthesizing and refining inference functions (to handle imperfect information games).

Imperfect information games can be considered a special kind of (multi-agent) partially observable Markov decision process (POMDP). Learning such models from observational data is notoriously difficult. In very recent work, Curtis et al. (2025) introduce "POMDP Coder", which learns a partially observed CWM. However, unlike us, they assume the hidden states are observed in hindsight (at the end of the trajectory). By contrast, we also consider a "closed deck" scenario, in which the hidden states are never observed. In addition, Curtis et al. (2025) use a determinized belief space planner (related to the POMCP method of Silver & Veness (2010)), whereas we use ISMCTS (see Appendix B) or PPO (Appendix D).

There are other many other ways to use LLMs for reasoning in games and multiagent systems. A recent line of work focuses on using LLMs to construct game-theoretic models of arbitrary scenarios in order to derive and deploy intelligent, strategic policies. Gemp et al. (2024) treats an LLM as an environment transition operator, controllable via instruction sets. An extensive-form game tree is explicitly constructed in OpenSpiel and an equilibrium over instruction sets is computed. Daskalakis et al. (2024) demonstrates how to design a game tree for Romeo and Juliet with the assistance of an LLM, subsequently modifying the tree so that the classic story lies in the support of its Nash equilibrium. Xu et al. (2025) embeds several observed Werewolf dialogues in a latent space, clusters the messages to form a finite action space and resulting game tree, and then runs counterfactual regret minimization on this discrete latent representation to derive a policy. Mensfelt et al. (2024a) proposed an approach to automatically translate natural language descriptions of small bimatrix games to logic representations (similarly in Mensfelt et al. (2024b)). Most closely related to this work, Deng et al. (2025) automated the construction of explicit (imperfect-information) extensiveform game trees from natural language descriptions of games, including a debugging module to ensure the resulting Gambit (Savani & Turocy, 2024) representation was valid. In contrast to this work, they only conditioned on game descriptions (rules) not observed trajectories and applied their pipelines to games with game trees containing at most 25 decision nodes (Kuhn Poker); code-world models offer the potential to scale to much larger game instances in some cases due to their more efficient encoding of repeat transitions.

4 Methods

At a high level, when confronted with a new game, our general game playing agent follows these steps: First, it plays a few games to completion using a random policy. The data collected during each game forms a *trajectory*, which consists of observations, rewards, legal actions, and states at each timestep. Second, it uses a textual description of the rules of the game, plus the generated trajectories, to learn a CWM². Finally, the agent plays the game in an arena against other opponents, using an MCTS policy built on top of the synthetic CWM. For imperfect information games (IIGs) we use ISMCTS instead of MCTS. If all the

²Note: We could potentially update the CWM after each step of game play, as we acquire new data, but in this paper, we learn the model up-front, given the initial offline trajectories and game description, for reasons of efficiency.

synthesized elements are correct, as the amount of play-time compute increases, the playing behavior of our agent gets closer to optimal. Thus, in contrast with LLM-as-a-policy agents, we shift the burden on the LLM from producing a good policy to producing a good world model, which in turn enables planning methods to turn compute into playing performance.

4.1 Synthesizing the Code World Model

A CWM is a playable, approximate copy of a target game. It contains functions providing logic to update the game state when an action is taken (transition function, which includes a termination), the legal actions given a state, the observation given a state (observations and state differ in the case of IIGs), the distribution for chance nodes, and the reward function for a state. All these functions are deterministic, with randomness entering the game only through the actions of the chance player. To synthesize a new CWM, we provide the LLM with the game's rules and offline trajectories, and demand that it creates a CWM following the OpenSpiel API (Lanctot et al., 2019) format. See Appendix G for prompt details.

A single-shot generation of the CWM will often be insufficient to produce a correct implementation of the game unless we add some kind of corrective feedback. Thus we subject the initial CWM to iterative refinement (Dainese et al., 2024; Tang et al., 2024b) to improve its quality. For refinement, a series of unit tests are automatically generated from the offline trajectories. For each transition in an offline trajectory, unit tests are generated in order to check the correctness of the CWM predictions as compared with the original trajectory (states, observations, rewards, legality of actions), and the absence of execution errors.

In the case of IIGs, this process requires that the offline trajectories contain not only the observations of the game and the actions of the players, but also the hidden states and the actions of all other players (including chance). The post-hoc availability of hidden states, an assumption also used in concurrent work (Curtis et al., 2025), can sometimes be unrealistic. Sec 4.4 introduces a novel approach to handle CWM learning from partially observed trajectories.

Unit tests are binary, so we can measure the *transition accuracy* as the rate of correctness of such tests. We refine the CWM until perfect transition accuracy (1.0) is achieved or our refinement budget runs out. We feed back the stack traces from failed unit tests to the LLM to help the refinement. We consider two separate approaches to refinement:

Conversation (sequential refinement). This is a serial "chat mode" approach, in which the stack trace of a newly failed unit test is appended to our previous interactions with the LLM to create the new prompt, and a new CWM addressing the unit test failure is requested. Failed unit tests derived from the offline trajectories are submitted to the LLM until all pass.

Tree search. Just like in the REx approach (Tang et al., 2024b;a), we maintain multiple CWMs in a refinement tree structure, and use Thompson sampling to choose which CWM to refine next, favoring those that either have high transition accuracy or have been refined few times. Each LLM call consists of a fresh prompt that contains the CWM chosen to be refined, the refinement instructions, and the stack trace of a failed unit test for that CWM. The prompts and hyperparameters used during synthesis are presented in Appendix G.

4.2 Synthesizing inference functions for IIGs

One of the novelties of our work is the synthesis of inference functions to enable the use of ISMCTS planning with the learned CWM at play time in imperfect information games (IIGs). To see why this is necessary, note that ISMCTS requires that at each game step t the agent can estimate the hidden state of the game s_t , as explained in Appendix B. More precisely, at play time, agent i must be able to sample from its belief state $p_M(s_t|o_{1:t}^i, a_{1:t}^i)$, where M is the estimated CWM³. Since exact inference incurs an exponential cost in the worst case, we ask the LLM to synthesize code to approximately sample from the posterior,

 $^{^{3}}$ For players other than i, we assume a uniform prior on the legal actions defined by the CWM. Only the support of this prior affects our approach, as we will focus on posterior support, see below.

utilizing only agent *i*'s actions $a_{1:t}^i$ and observations $o_{1:t}^i$ so far from the offline trajectory. We consider two alternative approaches to achieve this goal: hidden history inference and hidden state inference. We describe these below.

Hidden history inference. Since all the functions in the CWM are deterministic, the posterior over the hidden state s_t can be obtained from the posterior over the action history h_t , which includes the actions of the chance player. In this approach, the agent controlling player i asks the LLM to create a function that samples $\tilde{h}_t \sim p_M(h_t|o_{1:t}^i, a_{1:t}^i)$. The CWM can then be used to execute \tilde{h}_t and recreate a history of hidden states $\tilde{s}_{1:t}$ and observations $\tilde{o}_{1:t}^i$. A unit test is created for each time step t in which player i acts, verifying that the sampled values match the run time evidence (i.e., $\tilde{o}_t^i = o_t^i$ and $\tilde{a}_t^i = a_t^i$). This allows refinement (on the offline trajectories) to be applied to the inference function.

Once the refined inference function passes all unit tests⁴ (i.e., inference accuracy is 1.0), we can claim that the sampled \tilde{h}_t belongs to the support of $p_M(h_t|o^i_{1:t},a^i_{1:t})$, and therefore, the \tilde{s}_t generated by this process belongs to the support of $p_M(s_t|o^i_{1:t},a^i_{1:t})$. Although this does not guarantee that \tilde{s}_t is correctly distributed, the correct support is already very informative, given the extremely sparse support of state posteriors in games. Furthermore, this approach guarantees that the sampled posterior state \tilde{s}_t is a valid CWM state. Note that at play time the (test) inference accuracy can drop below 1.0 (depending on how well the synthesized inference code generalizes to novel observations), meaning that the approximate posterior samples might not always belong to the support of the actual posterior. However, \tilde{s}_t is still guaranteed by construction to be a valid hidden state in the CWM.

Hidden state inference. Rather than obtaining a state posterior sample indirectly through the action history, it is also possible to ask the LLM to create code that directly samples $\tilde{s}_t \sim p_M(s_t|o^i_{1:t},a^i)$. Then, the CWM can be used to obtain \tilde{o}_t from \tilde{s}_t . Correctness of the inference function can be partially validated by a unit test at each time step that verifies that the sampled values match the actual observations, $\tilde{o}_t = o_t$. CWM refinement can then be used to improve the synthesized inference function. State inference is potentially much simpler than full history inference, but it cannot guarantee that the produced sample \tilde{s}_t belongs to the support of the posterior, nor that it constitutes a valid CWM hidden state, because it ignores the dependency between consecutive states.

4.3 Synthesizing value functions

Another novelty of our work is the synthesis of value functions to speed up and improve value estimation in MCTS and ISMCTS. This can be faster (and potentially more accurate) than estimating the value of a new leaf node through random rollouts. To synthesize a deterministic value function V(s) to estimate the value of the (potentially hidden) state at leaf nodes, we can prompt the LLM to generate code, just as we did for learning the CWM. However, value functions are not refined, since there is no ground truth to compare to. Instead, multiple functions are generated and the best one is selected through a tournament.

4.4 Open deck vs closed deck during training

So far we assumed that the offline trajectories (used to train the CWM) contained hidden state information even for IIGs. Concurrent work Curtis et al. (2025) also assumes the ability to peek at hidden states. We refer to this setup as *open deck* synthesis⁵. This setup is justified in several practical scenarios, such as in a cooperative training environment where players share information to learn the mechanics of the game, during the design phase of a new game where developers have full access to the state, or when a human expert provides fully annotated "open-book" demonstrations to bootstrap an agent's understanding.

⁴Unlike the CWM functions, inference functions are stochastic (samplers). Thus, their unit tests are potentially stochastic, but for correct inference functions they will deterministically pass.

⁵We want to emphasize that in our open deck setting, hidden state information is only available in the offline trajectories to aid CWM synthesis, and not during actual game play. Thus the players only ever see observations, but the CWM learner may see hidden states (in the open deck setting).

However, there are scenarios in which the agent can only ever access its own observations and actions, so that the open deck assumption is violated. This would be the case, e.g., if the agent plays a novel game online. We refer to this scenario as *closed deck* synthesis; to the best of our knowledge, this scenario has not been addressed in prior CWM work.

To handle this scenario, we propose to combine the pieces introduced in the previous sections to build a regularized CWM "autoencoder". The idea is as follows: we ask the LLM to generate a CWM and a hidden history inference function, just like above, but we drop all the unit tests that are not verifiable without access to the hidden information (i.e., those checking the transition accuracy between consecutive hidden states), and we just keep the ones that we can verify (i.e., checking the result of mapping observations to hidden states and back to observations). We additionally add unit tests to a few iterations of random play ensuring that there are no execution errors. In other words, we refine based on the inference accuracy and lack of execution errors. This generates a kind of autoencoder, where the inference function acts as an encoder, producing a hidden sequence of actions \tilde{h}_t from $o_{1:t}^i, a_{1:t}^i$ and the CWM acts as a decoder, recreating the observations and actions from the latent h_t . Instead of a bottleneck, or a regularization term, the game rules and the required OpenSpiel API (used in the unit tests) introduced in the context of the LLM act as regularizers to prevent trivial latent spaces from being discovered. Valid posterior histories \tilde{h}_t (i.e., those that pass all unit tests) can be used to obtain a lower bound on the likelihood of the CWM, as follows: $p_M(o_{1:t}^i) = \sum_{h_t} p_M(o_{1:t}^i|h_t) p_M(h_t) \le p_M(o_{1:t}^i|\tilde{h}_t) p_M(\tilde{h}_t) = p_M(\tilde{h}_t).$ (The last equality follows because $p_M(o_{1:T}^i|\tilde{h}_t)=1$ when all unit test pass.) This lower bound is tightest when h_t is the maximum a posteriori, but is valid for any sample.

5 Experiments

Following the approach described in Sec. 4, we build an agent, which we call CWM-(IS)MCTS, which performs CWM synthesis (using either open or closed deck trajectories), and then plays using MCTS or ISMCTS. (We also tried learning a policy using PPO; see Appendix D for details.) We measure the playing abilities of our agent on multiple games against three other agents: A random legal action executor called Random; an (IS)MCTS agent that has access to the game's ground truth (GT) code, including inference functions but not value functions, which we call GT-(IS)MCTS; and an LLM as a policy, which we call Gemini 2.5Pro (we use "dynamic thinking", rather than specifying a thinking budget). All methods have access to the same data: the rules of the game as text and 5 offline trajectories. (IS)MCTS approaches always run 1,000 simulations before taking an action, using either the value function or 10 rollouts (in which all players act randomly) to determine the initial value of a new leaf node. A sketch of the information flow for each agent is given in Appendix F.

To validate the generality of our approach we use both perfect and imperfect information games, as well as well-known and OOD games. The perfect information games are: Tictac-toe, Connect four, Backgammon, Generalized tic-tac-toe (OOD), and Generalized chess (OOD). The imperfect information games are: Leduc poker, Bargaining, Gin rummy, Quadranto (OOD), and Hand of war (OOD). The out-of-distribution (OOD) games are not part of the LLMs training set, and have been created by us for these experiments. See Appendix H for the rules of each game.

5.1 Synthesis accuracy

The CWM agent operates by synthesizing a CWM of the game (and potentially other auxiliary functions) prior to playing the game, see Sec. 4 for details. We use Gemini 2.5 Pro for synthesis. For the concrete prompts used during synthesis, see Appendix G. For examples of synthesized code, see Appendix I.

Refinement attempts to increase the fraction of units tests that pass, iterating until all pass or the budget for LLM calls is exhausted. The fraction of unit tests that the CWM passes is the *training transition accuracy*, and the fraction of tests that the inference function passes is the *training inference accuracy*. To check for overfitting to the offline trajectories, after

synthesis, we measure the accuracy on a separate test set of 10,000 transitions, randomly sampled from 100 games where each player is randomly assigned a random policy or MCTS on the ground truth game code. This yields the test transition accuracy and test inference accuracy. The test set is never used to train on; instead it is used to estimate the accuracy of the learned CWMs. Finally, at play time against the LLM as a policy, online transitions are observed, and again used to assess the accuracy of the CWM and inference functions.

5.1.1 Perfect information games

For perfect information games, we find that we can learn a correct CWM for all the games, and that the resulting learned models have high test (generalization) accuracy. Both conversation and tree search work very well in this setting. Appendix C contains precise numbers (Tables 4 and 5, respectively), and shows the quick convergence of the CWM with the number of LLM calls (Fig. 6). We will stick with tree search for the remainder of this paper, since its ability to backtrack confers it additional resilience in harder settings.

5.1.2 Imperfect information games, open deck

In the case of imperfect information games (open deck learning), we find that the transition accuracy of the learned CWMs is very high, except for Gin rummy, where the training accuracy is just 84% and the test accuracy is 79%. See Table 1 for details. We hypothesize this is due to its high degree of logical and procedural complexity. Unlike games with more uniform rules, Gin rummy involves a multi-stage scoring phase (knocking, laying off melds, calculating deadwood, and checking for undercuts) that is difficult for the LLM to capture perfectly in code from a small number of trajectories. This highlights a key frontier for CWM synthesis: mastering games with intricate, multi-step procedural subroutines.

We also measure the inference accuracy obtained by the synthetic inference functions. We tried both hidden history and hidden state inference (see Sec.4.2). Results with hidden history inference (shown in Table 1 and Fig.1) are slightly better, so this will be the method of choice for the CWM-ISMCTS agent. (The results with hidden state inference are provided in Appendix C, Table 6 and Fig. 7.) Results for 3 of the 5 games are good, but once again we see that results for Gin rummy are quite poor (inference accuracy is only about 52%), and to a lesser extent Hand of war (inference accuracy is about 94%), even though CWM accuracy for Hand of war is good (about 98%). This suggests that hidden history inference is harder than learning the transition dynamics from a fully observed sequence of trajectories.

Table 1: Imperfect info. games, CWM refinement via tree search, hidden history inference.

Game	OOD	transi	transition accuracy			ence accu	racy	# LLM calls
Game	OOD	train	test	online	train	test	online	# LLW Cans
Bargaining	Х	1.0000	0.9827	1.0000	1.0000	1.0000	1.0000	23.0
Leduc poker	X	1.0000	0.9977	0.9942	1.0000	1.0000	1.0000	4.4
Gin rummy	X	0.8360	0.7881	0.9044	0.6550	0.5189	0.9678	701.2
Quadranto	✓	1.0000	1.0000	1.0000	1.0000	0.9864	0.9916	6.0
Hand of war	~	1.0000	0.9814	0.9868	1.0000	0.9357	1.0000	144.0

5.1.3 Imperfect information games, closed deck

Finally, we consider CWM synthesis with refinement in the novel closed deck setup in which no hidden information is available, not even post-hoc. The results in Table 2 show degradation on the synthesis quality with respect to the open deck setting of Table 1. Despite this, game play performance does not degrade significantly, as we show in the next section.

5.2 Arena: Game play performance

In this section, we test how the previous synthesis results translate into playing performance against other opponents in our game arena. Since the CWM synthesis process is stochastic,

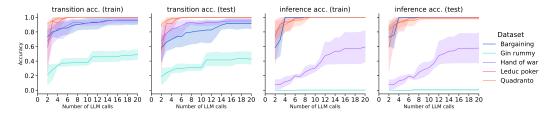


Figure 1: Evolution of the transition and inference accuracy with the number of LLM calls for imperfect games with refinement via tree search and hidden history inference.

Table 2: Imperfect information games, hidden history inference, closed deck.

Game	OOD	infer	ence accur	acy	# LLM calls
Game	OOD	train	test	online	# LLW Cans
Bargaining	Х	1.00000	0.67359	0.76000	88.2
Leduc poker	X	1.00000	0.97080	0.96585	9.0
Gin rummy	X	0.06215	0.10397	0.53953	467.2
Quadranto	~	1.00000	0.95183	0.96085	99.0
Hand of war	~	0.98125	0.92846	0.94835	475.4

we repeat it 5 times, automatically rejecting bad samples (see Appendix E), and pick a random CWM for each match. Results correspond to the average of 100 matches.

5.2.1 Perfect information games

All of our perfect information games are ternary-outcome games, so we are limited to win, lose, or draw (W/L/D). Fig. 2 shows the performance of our CWM-MCTS agent, when acting as Player 0 or Player 1, against three different competitors. A player forfeits when it fails to provide a valid action in the allotted time. The middle pair of bars of each panel show CWM-MCTS playing against GT-MCTS, an upper bound for performance that uses the ground truth (GT) code of the game for planning. Both agents are similarly good, without either of them clearly winning in any of the games. This highlights the quality of our code synthesis. CWM-MCTS is able to beat Gemini 2.5Pro (which is used as a policy) in all the considered games. For detailed numerical results, see Table 7 in Appendix C. We used a synthetic value function for Gen. tic-tac-toe, see Fig. 9 for the ablation without value function.

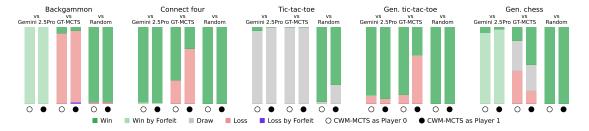


Figure 2: W/L/D rates for game play between CWM-MCTS and three opponents. CWMs are refined via tree search and hidden history inference.

5.2.2 Imperfect information games, open deck

Our imperfect information games contain a mixture of ternary-outcome games, zero-sum games and general-sum games (see Table 3 for a summary of all the games' characteristics). Win/loss/draw rates and payoff distributions are shown in Fig. 3. Except for Hand of war, CWM-ISMCTS beats or matches Gemini 2.5Pro in all imperfect information games. In the case of Gin rummy, this should be interpreted as Gemini 2.5Pro being a very weak player, you

can check its forfeit rate in Table 14. For Leduc poker, although our average performance is superior, we also observe high variance. For Bargaining we used a synthetic value function, which results in a significant improvement when CWM-ISMCTS acts as player 1 (see Fig. 9 in Appendix C for the corresponding ablation). We did not observe an improvement or degradation in performance when value functions were applied to the other games.

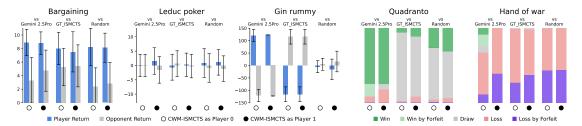


Figure 3: W/L/D rates and payoff distributions for game play between CWM-ISMCTS and three opponents. CWMs are refined via tree search and hidden history inference, *open deck*.

5.2.3 Imperfect information games, closed deck

Finally, we consider the closed deck setting, in which games are strictly partially observable, and no hidden state information or actions from other players are available in the offline trajectories. Results degrade w.r.t. the open deck setting, but CWM-ISMCTS-Closed continues to beat or match Gemini 2.5Pro (with high variance in the case of Leduc poker). We hypothesize that the non-intuitive improvement of CWM-ISMCTS-Closed at Hand of war w.r.t. the open deck setting could be due to the freedom to synthesize simpler state spaces when playing closed deck. Refer to Tables 12 and 13 in Appendix C for detailed results.

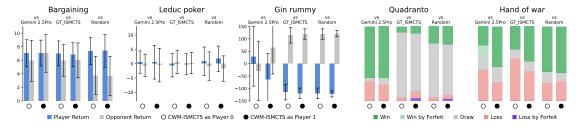


Figure 4: W/L/D rates and payoff distributions for game play between CWM-ISMCTS and three opponents. CWMs refined via tree search and hidden history inference, closed deck.

6 Discussion

In this work we extend the existing CWM framework by considering two-player games, performing value function code synthesis to improve player performance, introducing the concept of "inference as code" to enable state estimation in imperfect information games, and providing a learning algorithm (based on code-based autoencoders) to enable learning in the novel closed deck (strict partial observability) setting. Our results show the superiority of this approach with respect to LLMs as policies on multiple perfect and imperfect information games, including newly created ones.

However, we also notice that our method struggles to learn the rules of Gin rummy, an imperfect information game with intricate logic, especially in the very challenging closed deck setting. In future work, we hope to extend our method to enable active and online learning of the world model, so the agent can more effectively discover the true hidden causal mechanisms underlying each game (c.f., (Geng et al., 2025)). In addition, we would like to extend the technique to handle open-world games with free-form text and/or visual interfaces, so as to evaluate it on larger sets of novel games, see (Ying et al., 2025).

References

- Stefano V. Albrecht, Filippos Christianos, and Lukas Schäfer. *Multi-Agent Reinforcement Learning: Foundations and Modern Approaches*. MIT Press, 2024. URL https://www.marl-book.com.
- Lucia Cipolina-Kun, Marianna Nezhurina, and Jenia Jitsev. Game reasoning arena: A framework and benchmark for assessing reasoning capabilities of large language models via game play, 2025. URL https://arxiv.org/abs/2508.03368.
- Anthony Costarelli, Mat Allen, Roman Hauksson, Grace Sodunke, Suhas Hariharan, Carlson Cheng, Wenjie Li, Joshua Clymer, and Arjun Yadav. Gamebench: Evaluating strategic reasoning abilities of llm agents, 2024. URL https://arxiv.org/abs/2406.06613.
- Rémi Coulom. Efficient selectivity and backup operators in monte-carlo tree search. In H. Jaap van den Herik, Paolo Ciancarini, and H. H. L. M. (Jeroen) Donkers (eds.), *Computers and Games*, pp. 72–83, Berlin, Heidelberg, 2007. Springer Berlin Heidelberg. ISBN 978-3-540-75538-8.
- Peter I. Cowling, Edward J. Powley, and Daniel Whitehouse. Information set Monte Carlo tree search. *IEEE Transactions on Computational Intelligence and AI in Games*, 4: 120–143, 2012.
- Aidan Curtis, Hao Tang, Thiago Veloso, Kevin Ellis, Joshua Tenenbaum, Tomás Lozano-Pérez, and Leslie Pack Kaelbling. LLM-guided probabilistic program induction for POMDP model estimation. arXiv preprint arXiv:2505.02216, 2025.
- Nicola Dainese, Matteo Merler, Minttu Alakuijala, and Pekka Marttinen. Generating code world models with large language models guided by monte carlo tree search. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), Advances in Neural Information Processing Systems, volume 37, pp. 60429–60474. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/6f479ea488e0908ac8b1b37b27fd134c-Paper-Conference.pdf.
- Constantinos Daskalakis, Ian Gemp, Yanchen Jiang, Renato Paes Leme, Christos Papadimitriou, and Georgios Piliouras. Charting the shapes of stories with game theory. In *NeurIPS Creative AI Track*, 2024.
- Shilong Deng, Yongzhao Wang, and Rahul Savani. From natural language to extensive-form game representations. In *Proceedings of the 24th International Conference on Autonomous Agents and Multiagent Systems*, pp. 593–601, 2025.
- Jinhao Duan, Renming Zhang, James Diffenderfer, Bhavya Kailkhura, Lichao Sun, Elias Stengel-Eskin, Mohit Bansal, Tianlong Chen, and Kaidi Xu. Gtbench: Uncovering the strategic reasoning capabilities of llms via game-theoretic evaluations. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), Advances in Neural Information Processing Systems, volume 37, pp. 28219–28253. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/2024/file/3191170938b6102e5c203b036b7c16dd-Paper-Conference.pdf.
- Ian Gemp, Roma Patel, Yoram Bachrach, Marc Lanctot, Vibhavari Dasagi, Luke Marris, Georgios Piliouras, Siqi Liu, and Karl Tuyls. Steering language models with gametheoretic solvers. In *Agentic Markets Workshop at ICML 2024*, 2024.
- Jiayi Geng, Howard Chen, Dilip Arumugam, and Thomas L Griffiths. Are large language models reliable AI scientists? assessing reverse-engineering of black-box systems. arXiv [cs.LG], May 2025. URL http://arxiv.org/abs/2505.17968.
- Leon Guertler, Bobby Cheng, Simon Yu, Bo Liu, Leshem Choshen, and Cheston Tan. Textarena, 2025. URL https://arxiv.org/abs/2504.11442.
- Lanxiang Hu, Mingjia Huo, Yuxuan Zhang, Haoyang Yu, Eric P. Xing, Ion Stoica, Tajana Rosing, Haojian Jin, and Hao Zhang. lmgame-bench: How good are llms at playing games?, 2025a. URL https://arxiv.org/abs/2505.15146.

544

545 546

547

548

549

550

551

552

554

555

556

558

559

561

562

563

565 566

567

568

569

570

571 572

573 574

575

576

577

578

579

580

582

583 584

585

586

588

589

590

- Sihao Hu, Tiansheng Huang, Gaowen Liu, Ramana Rao Kompella, Fatih Ilhan, Selim Furkan Tekin, Yichang Xu, Zachary Yahn, and Ling Liu. A survey on large language model-based 542 game agents, 2025b. URL https://arxiv.org/abs/2404.02039. 543
 - Daniel Kahneman. A perspective on judgement and choice. American Psychologist, 58: 697-720, 2003.
 - Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980, 2014.
 - Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *Proceedings* of the 17th European Conference on Machine Learning, ECML'06, pp. 282293, Berlin, Heidelberg, 2006. Springer-Verlag. ISBN 354045375X. doi: 10.1007/11871842_29. URL https://doi.org/10.1007/11871842_29.
 - H. W. Kuhn. Extensive games and the problem of information. Annals of Mathematics Studies, 28:193–216, 1953.
 - Marc Lanctot, Edward Lockhart, Jean-Baptiste Lespiau, Vinicius Zambaldi, Satyaki Upadhvay, Julien Pérolat, Sriram Srinivasan, Finbarr Timbers, Karl Tuyls, Shayegan Omidshafiei, Daniel Hennes, Dustin Morrill, Paul Muller, Timo Ewalds, Ryan Faulkner, János Kramár, Bart De Vylder, Brennan Saeta, James Bradbury, David Ding, Sebastian Borgeaud, Matthew Lai, Julian Schrittwieser, Thomas Anthony, Edward Hughes, Ivo Danihelka, and Jonah Ryan-Davis. OpenSpiel: A framework for reinforcement learning in games. CoRR, abs/1908.09453, 2019. URL http://arxiv.org/abs/1908.09453.
 - Yi Liao, Yu Gu, Yuan Sui, Zining Zhu, Yifan Lu, Guohua Tang, Zhongqian Sun, and Wei Yang. Think in games: Learning to reason in games via reinforcement learning with large language models, 2025. URL https://arxiv.org/abs/2508.21365.
 - Agnieszka Mensfelt, Kostas Stathis, and Vince Trencsenyi. Autoformalization of game descriptions using large language models. arXiv preprint arXiv:2409.12300, 2024a.
 - Agnieszka Mensfelt, Kostas Stathis, and Vince Trencsenyi. Autoformalizing and simulating game-theoretic scenarios using llm-augmented agents. arXiv preprint arXiv:2412.08805, 2024b.
 - Kevin Murphy. Reinforcement learning: An overview, 2025. URL https://arxiv.org/abs/ 2412.05265.
 - Anian Ruoss, Grégoire Delétang, Sourabh Medapati, Jordi Grau-Moya, Li Kevin Wenliang, Elliot Catt, John Reid, Cannada A. Lewis, Joel Veness, and Tim Genewein. Amortized planning with large-scale transformers: A case study on chess. In A. Globerson, L. Mackey, D. Belgrave, A. Fan, U. Paquet, J. Tomczak, and C. Zhang (eds.), Advances in Neural Information Processing Systems, volume 37, pp. 65765–65790. Curran Associates, Inc., 2024. URL https://proceedings.neurips.cc/paper_files/paper/ 2024/file/78f0db30c39c850de728c769f42fc903-Paper-Conference.pdf.
 - Rahul Savani and Theodore L. Turocy. Gambit: The package for computation in game theory, version 16.2.0 edition, 2024. https://www.gambit-project.org.
 - John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347, 2017.
 - John Schultz, Jakub Adamek, Matej Jusup, Marc Lanctot, Michael Kaisers, Sarah Perrin, Daniel Hennes, Jeremy Shar, Cannada Lewis, Anian Ruoss, Tom Zahavy, Petar Velikovi, Laurel Prince, Satinder Singh, Eric Malmi, and Nenad Tomaev. Mastering board games by external and internal planning with language models. In Proceedings of the Forty-Second International Conference on Machine Learning (ICML), 2025. URL https:// arxiv.org/abs/2412.12119.
 - Y. Shoham and K. Leyton-Brown. Multiagent Systems: Algorithmic, Game-Theoretic, and Logical Foundations. Cambridge University Press, 2009.

- David Silver and Joel Veness. Monte-carlo planning in large pomdps. Advances in neural information processing systems, 23, 2010.
- Haoran Sun, Yusen Wu, Peng Wang, Wei Chen, Yukun Cheng, Xiaotie Deng, and Xu Chu. Game theory meets large language models: A systematic survey with taxonomy and new frontiers, 2025. URL https://arxiv.org/abs/2502.09053.
- Hao Tang, Keya Hu, Jin Peng Zhou, Si Cheng Zhong, Wei-Long Zheng, Xujie Si, and Kevin Ellis. Code repair with LLMs gives an exploration-exploitation tradeoff. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*, November 2024a. URL https://openreview.net/pdf?id=o863gX6DxA.
- Hao Tang, Darren Key, and Kevin Ellis. Worldcoder, a model-based llm agent: Building world models by writing code and interacting with the environment, 2024b. URL https://arxiv.org/abs/2402.12275.
- Vivek Verma, David Huang, William Chen, Dan Klein, and Nicholas Tomlin. Measuring general intelligence with generated games, 2025. URL https://arxiv.org/abs/2505.07215.
- Zelai Xu, Wanjun Gu, Chao Yu, Yi Wu, and Yu Wang. Learning strategic language agents in the werewolf game with iterative latent space policy optimization. In Forty-second International Conference on Machine Learning, 2025.
- Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models., 2022. URL https://arxiv.org/abs/2210.03629.
- Lance Ying, Katherine M Collins, Prafull Sharma, Cedric Colas, Kaiya Ivy Zhao, Adrian Weller, Zenna Tavares, Phillip Isola, Samuel J Gershman, Jacob D Andreas, Thomas L Griffiths, Francois Chollet, Kelsey R Allen, and Joshua B Tenenbaum. Assessing adaptive world models in machines with novel games. <code>arXiv</code> [cs.AI], July 2025. URL http://arxiv.org/abs/2507.12821.

A Information on the games

A summary of the games that we use in our experiments is given in Table 3.

Table 3: Details of the games that we use. The columns have the following meaning: OOD: whether the game is novel (no source code on the internet); Observability: Full means perfect information game, partial means imperfect information game; Payoff: W/L/D means Win/Lose/Draw, General means general sum; # Actions: number of possible actions; Obs. dim.: dimensionality of the observation tensor IS dim.: dimensionality of the information set (i.e., game's hidden state).

Name	OOD	Observability	Payoff	# Actions	Obs. dim.	IS dim.
Backgammon	X	Full	W/L/D	1352	200	
Connect four	X	Full	W/L/D	7	126	
Tic-tac-toe	X	Full	W/L/D	9	27	
Gen. tic-tac-toe	✓	Full	W/L/D	36	108	
Gen. chess	~	Full	W/L/D	5555	250	
Bargaining	Х	Partial	General	121	93	309
Leduc poker	X	Partial	Zero-sum	3	16	30
Gin rummy	X	Partial	Zero-sum	241	644	655
Quadranto	✓	Partial	W/L/D	5	9	7
Hand of war	✓	Partial	W/L/D	16	27	73

B Information Set Monte Carlo Tree Search

Recall that a history encodes the sequence of actions taken by all players, including chance. But in an imperfect information game, not all aspects of the history are observable. For instance, in a game of poker, h contains information about the cards held by all players (as chosen by the dealers actions), but some of this information is private and hence not known by some players. After an action is executed and added to the history $(h_{t-1}, a_t) \equiv h_t$, each player $i \in \mathcal{N}$ perceives individual observations $o_t^i(h_t)$. The state (from the perspective of an agent i) is then a function of $o_{1:t}^i$, e.g., just the last observation.

To choose actions in an IIG, we can use the Information Set MCTS method of (Cowling et al., 2012), which we now describe. First, recall that in classical MCTS, there is a root node corresponding to the current state of the game which all simulations start from, and non-root nodes which correspond to states that occur after the root state. At each node, statistics such as average values state-action values, $\hat{Q}(s,a)$, and simulation counts are maintained. The main differences in ISMCTS are: (i) the simulations start at a distribution of possible ground truth states and (ii) statistics are maintained and aggregated across information states with respect to the current player.

Figure 5 contains an example with a simplified poker game with a deck of three cards (Jack, Queen, King). In this example, the current player has received the King as a private card and no actions have yet been taken, so there are only two ground truth states: the opponent could have either the Queen or the Jack. An iteration first samples the Queen and continues with this ground truth state h_0 , sampling actions, and generating histories h_1, h_2, h_3 , and so on until the first node not in the tree is encountered. It is then added to the tree, and a random rollout policy takes over until a terminal state. The dotted boxes are the analogs of nodes stored in a tree (or lookup table) and correspond to information states. Return estimates (i.e., Q-value statistics) and visit counts are maintained in these nodes as in classical MCTS (Coulom, 2007) (aggregated over different samplings of ground truth states), and UCB is used to select actions in the standard way (Kocsis & Szepesvári, 2006).

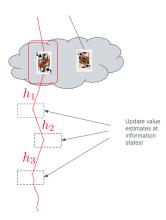


Figure 5: ISMCTS. A search tree is built over possible ground truth histories (e.g. h_1, h_2, \ldots). Because the player cannot distinguish between certain histories, statistics are aggregated at the level of information sets (dotted boxes), which group all histories that appear identical to the player.

C ADDITIONAL EXPERIMENTAL RESULTS

In the interest of space, some additional experimental results are included in this section.

C.1 Synthesis

C.1.1 Accuracy of Learned Transition and Inference functions

Comparing Table 4 and 5, it is apparent that even though both options work reasonably well, tree search has the edge, both in terms of accuracy (higher) and number of LLM calls (lower).

Table 4: Perfect information games, refinement via tree search.

Game	OOD	trans	ition accur	racy	# LLM calls
Game	OOD	train	test	online	# LLW Cans
Backgammon	Х	1.00000	0.99932	1.00000	16.8
Connect four	X	1.00000	1.00000	1.00000	2.0
Tic-tac-toe	X	1.00000	1.00000	1.00000	2.0
Gen. tic-tac-toe	✓	1.00000	1.00000	1.00000	2.4
Gen. chess	~	1.00000	1.00000	1.00000	5.2

Table 5: Perfect information games, refinement via conversation.

Game	OOD	trans	transition accuracy					
Game	OOD	train	test	online	# LLM calls			
Backgammon	Х	1.00000	0.99944	1.00000	13.2			
Connect four	X	1.00000	1.00000	1.00000	3.2			
Tic-tac-toe	X	1.00000	1.00000	1.00000	2.0			
Gen. tic-tac-toe	~	1.00000	1.00000	1.00000	2.4			
Gen. chess	•	1.00000	1.00000	1.00000	4.2			

Table 6: Imperfect info. games, refinement via tree search, hidden state inference.

Game	OOD	transi	transition accuracy			ence accu	racy	# LLM calls
Game	OOD	train	test	online	train	test	online	# LLW Cans
Bargaining	Х	1.0000	0.9482	0.8712	1.0000	1.0000	1.0000	32.8
Leduc poker	X	1.0000	0.9854	0.9942	1.0000	1.0000	1.0000	4.2
Gin rummy	X	0.8999	0.8293	0.8909	1.0000	0.9513	0.9738	813.4
Quadranto	✓	1.0000	1.0000	0.9991	1.0000	0.9911	0.9876	7.4
Hand of war	~	1.0000	0.9782	0.9806	1.0000	1.0000	1.0000	28.0

${ m C.1.2}$ Accuracy of learned transition and inference functions vs number of LLM calls

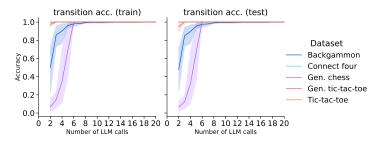


Figure 6: Evolution of the transition accuracy of the best generated CWM with the number of LLM calls for perfect games (with CWM refinement via tree search).

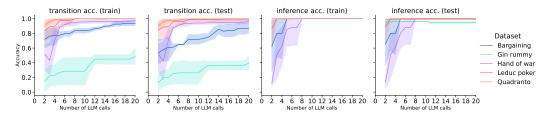


Figure 7: Evolution of the transition and inference accuracy with the number of LLM calls for imperfect games with refinement via tree search and hidden state inference.

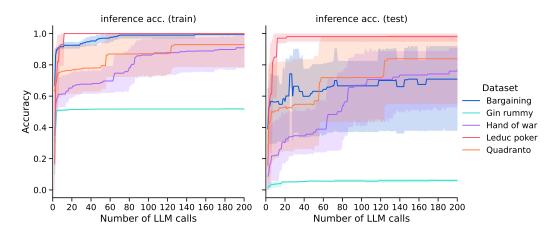


Figure 8: Evolution of the inference accuracy with the number of LLM calls for imperfect games with refinement via tree search with closed deck.

C.1.3 Tree search Settings

We use the following settings for treesearch throughout our experiments.

- heuristic_weight=5.0: Weight on the heuristic value (higher means more exploitation). The heuristic weight C adjusts the parameters α and β of the **Beta** prior on each arm Tang et al. (2024b). In particular, we set $\alpha = 1 + C \times h$ and $\beta = 1 + (1 C) \times h$, where the heuristic value h is the average pass rate of the unit tests.
- num_retries=500: Number of retries for tree search.
- num_tests_on_init=5: Number of tests of each type to include on the first synthesis.
- num_tests_on_error=1: Number of failed tests of each type to include during code refinement.
- min_heuristic_value_on_init=0.01: Minimum heuristic value to consider a node for expansion on initialization.
- min_heuristic_value_gain=0.01: Minimum heuristic value gain to consider a node for expansion.

C.2 Detailed Per-game arena results

For games in which the outcomes are win, lose or draw, we show the frequency of these 3 outcomes in 3 different columns, for each agent. For games with arbitrary payoff (Bargaining, LeDuc-Poker, Gin Rummy), we show the payoff to each player in 2 different columns. We consider the case when our agent acts as Player 0 or Player 1, and show these in different rows, to account for first-mover advantage.

For imperfect information games, we show results for hidden history inference (open deck learning), hidden state inference (open deck learning), and hidden history inference (closed deck learning).

For games in which the outcomes are win, lose or draw, we also report (in small font) the number of games with outcome that were forfeited vs the total number of games with that outcome. (A forfeit means the agent has either thrown an exception or tried to execute an illegal action, since our game arena API does not allow the agent to see which actions are legal at a given point in the game.)

C.2.1 Perfect information games

Table 7: Win rates using CWM refinement via tree search against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Р	Ge	emini 2.5 Pro		G	Т-МСТ	rs	I	Randon	1
Game	1	Win (forfeit/n)	Loss (forfeit/n)	Draw (n)	Win	Loss	Draw	Win	Loss	Draw
Backgammon		1.00 (100/100) 1.00 (100/100)	0.00 (0/0) 0.00 (0/0)	0.00 (0) 0.00 (0)	0.08 0.07	0.92 0.93	0.00	0.98 0.98	0.02 0.02	0.00
Connect four		1.00 (2/100) 1.00 (1/100)	0.00 (0/0) 0.00 (0/0)	0.00 (0)	$0.69 \\ 0.28$	$0.31 \\ 0.72$	0.00	1.00 1.00	0.00	0.00
Tic-tac-toe	X	0.05 (0/5) 0.00 (0/0)	0.00 (0/0) 0.00 (0/0)	0.95 (95) 1.00 (100)	0.00	0.00	1.00 1.00	$0.97 \\ 0.75$	0.00	$0.03 \\ 0.25$
Gen. tic-tac-toe	X	0.89 (0/89) 0.93 (0/93)	0.10 (0/10) 0.07 (0/7)	0.01 (1) 0.00 (0)	$0.88 \\ 0.37$	0.12 0.63	0.00 0.00	1.00 1.00	0.00	0.00
Gen. chess	 Å	1.00 (92/100) 1.00 (97/100)	0.00 (0/0) 0.00 (0/0)	0.00 (0) 0.00 (0)	0.18 0.49	$0.43 \\ 0.17$	0.39 0.34	1.00 1.00	0.00	0.00

C.2.2 HIDDEN HISTORY INFERENCE

Table 8: Payoffs using CWM refinement via tree search and hidden history inference against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Р	Gemini	2.5 Pro	GT-ISN	ICTS	Rand	lom
Game	1	Us	Them	Us	Them	Us	Them
Dangaining	(\$	8.90	3.31	8.01	5.26	8.21	2.41
Bargaining	€)	8.80	4.73	7.51	5.44	8.12	2.81
T - d1		-0.03	0.03	-0.65	0.65	0.86	-0.86
Leduc poker		1.55	-1.55	0.24	0.24	1.09_	-1.09
C :	٢	120.54	-120.54	-115.62	115.62	-4.92	4.92
Gin rummy	<i>e</i>	123.00	-123.00	-115.62	115.62	-15.99	15.99

Table 9: Win rates using CWM refinement via tree search and hidden history inference against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Р	Ge	Gemini 2.5 Pro			GT-ISMCTS			Random		
Game	1	Win (forfeit/n)	Loss (forfeit/n)	Draw (n)	Win	Loss	Draw	Win	Loss	Draw	
Quadranto		0.91 (16/91) 0.75 (0/75)	0.08 (0/8) 0.18 (0/18)	0.01 (1) 0.07 (7)	0.06 0.11	0.02 0.08	0.92 0.81	0.27 0.31	0.03 0.06	0.70 0.63	
Hand of war	•	0.35 (16/35) 0.33 (0/33)	0.56 (11/56) 0.62 (39/62)	0.09 (9) 0.05 (5)	0.20 0.31	$0.72 \\ 0.61$	0.08	0.33 0.33	0.57 0.62	0.10 0.05	

C.2.3 HIDDEN STATE INFERENCE

Table 10: Payoffs using CWM refinement via tree search and hidden state inference against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Р	Gemini	2.5 Pro	GT-ISN	1CTS	Rand	lom
Game	•	Us	Them	Us	Them	Us	Them
Dangaining		8.48	4.32	7.46	4.17	7.70	2.72
Bargaining	€ □	7.98	6.42	7.25	6.04	7.78	3.47
T - d1		1.75	-1.75	0.16	-0.16	1.19	-1.19
Leduc poker		0.37	-0.37	0.41	0.41	1.12	-1.12
C:	٢	66.42	-66.42	-114.39	114.39	-28.29	28.29
Gin rummy	r	121.77	-121.77	-121.77	121.77	-4.92	4.92

Table 11: Win rates using CWM refinement via tree search and hidden state inference against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game P	D	Ge	Gemini 2.5 Pro			GT-ISMCTS			Random		
Game		Win (forfeit/n)	$Loss\ ({\rm forfeit/n})$	Draw (n)	Win	Loss	Draw	Win	Loss	Draw	
Quadranto		0.58 (16/58) 0.37 (0/37)	0.40 (0/40) 0.54 (0/54)	$0.02_{(2)}$ $0.09_{(9)}$	$0.13 \\ 0.14$	0.02 0.01	$0.85 \\ 0.85$	0.19 0.26	$0.04 \\ 0.07$	0.77 0.67	
Hand of war	•	0.41 (16/41) 0.44 (0/44)	0.49 (0/49) 0.40 (0/40)	0.10 (10) 0.16 (16)	$0.25 \\ 0.42$	0.61 0.41	0.14 0.17	$0.59 \\ 0.63$	0.28 0.24	0.13 0.13	

C.2.4 HIDDEN HISTORY INFERENCE WITH CLOSED DECK LEARNING

Table 12: Payoffs using CWM refinement via tree search with closed deck against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Р	Gemini 2.5 Pro		GT-ISN	ICTS	Rand	om
Game		Us	Them	Us	Them	Us	Them
Dangaining	(\$	7.03	5.91	7.01	5.98	7.37	3.76
Bargaining	€	7.07	7.02	6.83	6.01	7.39 _	3.68
T - 1 1		0.57	-0.57	-0.56	0.56	0.86	-0.86
Leduc poker		0.49	-0.49	-0.21	0.21	1.71	-1.71
C:	r	29.52	-29.52	-114.39	114.39	-119.31	119.31
Gin rummy		-63.96	63.96	-119.31	119.31	-121.77	121.77

Table 13: Win rates using CWM refinement via tree search with closed deck against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game	D	Ge	emini 2.5 Pro		GT	-ISMC	TS	I	Randon	1
Game	1	Win (forfeit/n)	Loss (forfeit/n)	Draw (n)	Win	Loss	Draw	Win	Loss	Draw
Quadranto		$0.69 (0/69) \\ 0.71 (2/71)$	$\begin{array}{c} 0.26 \ \tiny{(0/26)} \\ 0.22 \ \tiny{(0/22)} \end{array}$	$0.05_{(5)}$ $0.07_{(7)}$	$0.08 \\ 0.13$	$0.05 \\ 0.09$	$0.87 \\ 0.78$	$0.23 \\ 0.27$	$0.05 \\ 0.05$	$0.72 \\ 0.68$
Hand of war	•	0.41 (16/41) 0.54 (0/54)	0.42 (0/42) 0.25 (0/25)	$0.17_{(17)}$ $0.21_{(21)}$	$0.31 \\ 0.47$	0.57 0.40	0.12 0.13	$0.61 \\ 0.62$	0.24 0.26	0.15 0.12

C.3 Forfeit rates for non-ternary-outcome games

Table 14: Forfeit rates for non-ternary-outcome games using CWM refinement via tree search and hidden history inference against multiple opponents. This is the rate at which each agent forfeits the game by failing to execute a legal action. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Р	Gemini	2.5 Pro	GT-IS	MCTS	Random		
Game	1	Us	Them	Us	Them	Us	Them	
Di	\$	0.00	0.00	0.00	0.00	0.00	0.00	
Bargaining	(ED)	_ 0.00 _	0.01	0.00	0.00	$_{0.00}_{-}$	_ 0.00	
Leduc poker		0.00	0.00	0.00	0.00	0.00	0.00	
Leduc poker		_ 0.00 _	0.00	0.00	0.00	$_{0.00}_{-}$	_ 0.00	
C:		0.01	0.99	0.94	0.00	0.04	0.00	
Gin rummy	<i>r</i>	0.00	1.00	0.94	0.00	0.13	0.00	

Table 15: Forfeit rates for non-ternary-outcome games using CWM refinement via tree search and hidden state inference against multiple opponents. This is the rate at which each agent forfeits the game by failing to execute a legal action. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Р	Gemini	2.5 Pro	GT-IS	MCTS	Random		
Game	1	Us	Them	Us	Them	Us	Them	
Dii	\$	0.00	0.00	0.00	0.00	0.00	0.00	
Bargaining	€)	0.00	0.00	0.00	0.00	_0.00_	0.00	
Leduc poker	•	0.00	0.16	0.00	0.00	0.00	0.00	
Leduc poker		_ 0.00 _	0.00	0.00	0.00	$_{0.00}_{-}$	0.00	
C:		0.23	0.77	0.93	0.00	0.23	0.00	
Gin rummy		0.00	0.99	0.99	0.00	0.04	0.00	

Table 16: Forfeit rates for non-ternary-outcome games using CWM refinement via tree search with closed deck against multiple opponents. This is the rate at which each agent forfeits the game by failing to execute a legal action. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Р	Gemini	2.5 Pro	GT-IS	MCTS	Random		
Game	1	Us	Them	Us	Them	Us	Them	
Dangaining	(\$)	0.00	0.00	0.00	0.00	0.00	0.00	
Bargaining	€)	_ 0.00 _	_0.00	0.00	0.00	_0.00_	0.00	
Ladua nalan		0.00	0.00	0.00	0.00	0.00	0.00	
Leduc poker		_ 0.00 _	_0.00	0.00	0.00	$_{0.00}_{-}$	0.00	
Cira muna		0.38	0.62	0.93	0.00	0.97	0.00	
Gin rummy		0.76	0.24	0.97	0.00	0.99	0.00	

C.4 VALUE FUNCTION ABLATIONS

As explained in the main text, the purpose of value functions is to speed up (IS)MCTS by providing a better value initialization for new leaf nodes. This can also result in higher quality selections for a fixed budget. Synthetic value functions are generated by the LLM in one-shot, and its usefulness assessed via a tournament ran on top of the synthesized CWM. Agents using different value functions (or potentially no value function) compete against each other the synthesized CWM to evaluate performance.

The use of value function only delivered improvements in the case of Gen. tic-tac-toe and Bargaining, so our agent only used value functions when playing those games. Note that the choice to use value functions or not can be assessed before actual online game play, by having the agent play locally (with and without using a value function) on its own synthetic CWM as a proxy, and assessing which option is most beneficial.

Fig. 9 shows the ablation corresponding to not using a value function in Gen. tic-tac-toe and Bargaining.

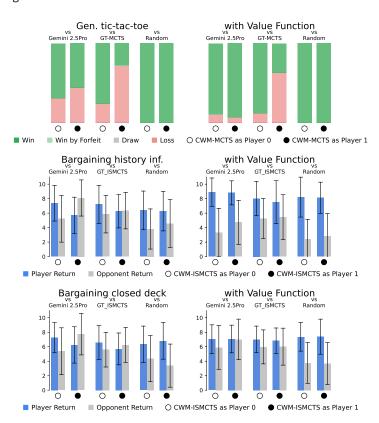


Figure 9: Ablation for Gen. tic-tac-toe and Bargaining. Effect of using synthesized value functions (right column) vs not (left column) to improve planning in CWMs.

D PLANNING WITH PPO INSTEAD OF (IS)MCTS

D.1 TRAINING A PPO AGENT ON TOP OF A CWM

The CWM agent discussed in the main paper relies on (IS)MCTS to take actions within its learned CWM. While effective, this online planning process can be slow. We investigate an alternative approach: amortizing the planning computation into a reactive policy, trained with the PPO algorithm (Schulman et al., 2017).

We entirely learn this PPO policy within the learned CWM environment. For each game, we train a PPO-CWM agent (acting either as Player 0 or Player 1) to maximize its rewards against an opponent that uniformly picks a legal action.

Mapping JSON observations to 1D tensors: The CWM represents observations in the JSON format provided by OpenSpiel, whereas the actor-critic networks we use (which are based on MLPs and RNNs) requires fixed-size 1D arrays as input.⁶ Consequently, we need a procedure to map each JSON into a flat tensor representation. We generate this mapping programmatically by prompting a LLM as shown below, providing the CWM training sequences as examples.

You are an expert reinforcement learning researcher and Python programmer.

Your task is to implement the following two functions which form a bijective pair:

```
def observation_to_tensor(obs) -> np.ndarray: # 1D
   ...

def tensor_to_observation(tensor) -> np.ndarray: # 1D
   ...

An example input dataset is as follows:
{example}
```

First reason about the problem and possible corner cases. Finally output only the resulting two functions without any placeholders.

Architecture. Our PPO agent uses an actor-critic architecture. For perfect information games, the actor and critic networks share a common feature extractor consisting of two 256-unit fully-connected layers with tanh activations. The actor head is a final linear layer that outputs logits for each action, which are then masked to ensure only legal moves are considered. The critic head is a separate linear layer that outputs a single scalar value.

For imperfect information games, we augment this architecture with a recurrent neural network to process historical information. An input observation x_t is first passed through a 256-unit linear layer (with tanh activation). The result is fed into an RNN along with the previous hidden state h_{t-1} to produce an output vector. This output is concatenated with the original input x_t and passed through a final 256-unit hidden layer (with tanh activation) before being fed to the actor and critic heads as described above.

PPO training. The PPO-CWM agent is trained for a total of 10M agent steps inside the CWM, using the hyperparameters in Table 17. From the two-player trajectories collected, we extract the single-agent sequence of observations, actions, and rewards corresponding to the PPO-CWM agent. This filtered data is used to compute the advantages and the final PPO loss objective.

For each game, and each player, we train 5 PPO-CWM agents with different seeds, and select the one with the highest win rate against the random opponent for final evaluation. This agent is then benchmarked in the Arena, as described in Sec. 5.2. We include matches against our CWM MCTS agent to compare both approaches for leveraging the CWM.

⁶We could use transformers, which can handle JSON strings, for the actor and critic, but such models would be much slower to train.

Table 17: PPO hyperpameters.

Module	Hyperparameter	Value
Environment	Number of environments Rollout horizon in environment	50 100
Advantage	$\gamma top \lambda$	0.99 0.95
Loss	ϵ clipping Value loss coefficient Entropy loss coefficient start Entropy loss coefficient end Entropy loss coefficient schedule	0.2 0.5 0.1 0.01 Linear
Learning	Optimizer Learning rate Max. gradient norm Learning rate annealing Number of minibatches (MFRL) Number of epochs (MFRL)	Adam (Kingma & Ba, 2014) 0.0003 0.5 False 10

D.2 Results

Arena results are presented in Tables 18 to 22. Note that PPO-CWM was not trained on Gin Rummy due to the poor performance of the CWM on that game. All the CWMs in this section have been trained on 100 offline trajectories.

PPO-CWM vs. Random. PPO-CWM outperforms the random agent for all the games.

PPO-CWM vs. Gemini 2.5 Pro. Our PPO-CWM agent outperforms or matches Gemini in all the games. For perfect information games, PPO-CWM wins in Backgammon, Generalized Chess and Tic-Tac-Toe; and exhibits mixed results (winning as one player and losing as the other) in Connect Four and Generalized Tic-Tac-Toe. For imperfect information games, for both open deck and closed deck, PPO-CWM wins in Bargaining and Quadranto, and ties in Hand of War and Leduc Poker.

PPO-CWM vs. CWM MCTS. For perfect information games, where the learned CWM is a near-perfect replica of the environment, PPO-CWM is outperformed by our CWM-MCTS agent. The only exception is Generalized Tic-Tac-Toe when PPO-CWM acts as Player 0. For imperfect information games, PPO-CWM wins in two games (Hand of War and Bargaining) and loses in the other two games (Leduc poker and Quadranto).

D.2.1 Games with Perfect Information

Table 18: PPO-CWM win rates using CWM refinement via tree search against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game	D	C	WM MCTS		Ge	emini 2.5 Pro		G	Т-МСТ	S	Random		
Game	1	Win (forfeit/n)	Loss (forfeit/n)	Draw (n)	Win (forfeit/n)	Loss (forfeit/n)	Draw (n)	Win	Loss	Draw	Win	Loss	Draw
D 1		0.01 (0/1)	0.99 (0/99)	0.00 (0)	1.00 (100/100)	0.00 (0/0)	0.00 (0)	0.02	0.98	0.00	0.92	0.08	0.00
Backgammon		0.03 (0/3)	0.97 (0/97)	0.00 (0)	1.00 (100/100)	0.00 (0/0)	0.00 (0)	0.01	0.99	0.00	0.94	0.06	0.00
C		0.00 (0/0)	1.00 (0/100)	0.00 (0)	0.92 (0/92)	0.08 (0/8)	0.00 (0)	0.00	1.00	0.00	1.00	0.00	0.00
Connect four	<u> </u>	0.00 (0/0)	1.00 (0/100)	0.00 (0)	0.02 (0/2)	0.98 (0/98)	0.00 (0)	0.00	1.00	0.00	0.99	0.01	0.00
m: ,	X	0.00 (0/0)	0.00 (0/0)	1.00 (100)	0.00 (0/0)	0.00 (0/0)	1.00 (100)	0.00	0.00	1.00	1.00	0.00	0.00
Tic-tac-toe	0	0.00 (0/0)	1.00 (0/100)	0.00(0)	0.87 (0/87)	0.12 (0/12)	0.01 (1)	0.00	1.00	0.00	0.91	0.01	0.08
C	X	0.45 (0/45)	0.55 (0/55)	0.00 (0)	0.91 (0/91)	0.09 (0/9)	0.00 (0)	0.54	0.46	0.00	1.00	0.00	0.00
Gen. tic-tac-toe	0	0.04 (0/4)	0.96 (0/96)	0.00 (0)	0.38 (0/38)	0.62 (0/62)	0.00 (0)	0.05	0.95	0.00	0.99	0.01	0.00
C1	8	0.00 (0/0)	1.00 (0/100)	0.00 (0)	0.94 (90/94)	0.06 (0/6)	0.00 (0)	0.00	1.00	0.00	1.00	0.00	0.00
Gen. chess	*	0.08 (0/8)	0.92 (0/92)	0.00 (0)	0.95 (5/95)	0.05 (0/5)	0.00 (0)	0.08	0.92	0.00	1.00	0.00	0.00

D.2.2 HIDDEN HISTORY INFERENCE

Table 19: PPO-CWM win rates using CWM refinement via tree search and hidden history inference against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Ъ	С	WM MCTS		Ge	mini 2.5 Pro		GT-ISMCTS			Random		
Game	1	Win (forfeit/n)	Loss (forfeit/n)	Draw (n)	Win (forfeit/n)	Loss (forfeit/n)	Draw (n)	Win	Loss	Draw	Win	Loss	Draw
01		0.04 (0/4)	0.68 (0/68)	0.28 (28)	0.75 (0/75)	0.25 (0/25)	0.00 (0)	0.00	0.64	0.36	0.55	0.16	0.29
Quadranto	Quadranto	0.08 (0/8)	0.34 (1/34)	0.58 (58)	0.75 (0/75)	0.15 (0/15)	0.10 (10)	0.04	0.58	0.38	0.56	0.11	0.33
II 1 . C	•	0.63 (36/63)	0.26 (0/26)	0.11 (11)	$0.34\ (0/34)$	0.45 (0/45)	0.21 (21)	0.30	0.60	0.10	0.63	0.23	0.14
Hand of war	4	0.69 (36/69)	0.27 (0/27)	0.04 (4)	$0.54\ (0/54)$	0.31 (0/31)	0.15 (15)	0.52	0.39	0.09	0.69	0.24	0.07

Table 20: PPO-CWM payoffs using CWM refinement via tree search and hidden history inference against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Р	CWM	MCTS	Gemini 2.5 Pro		GT-IS	MCTS	Random		
Game	1	Us	Them	Us	Them	Us	Them	Us	Them	
Dii		7.94	4.29	8.41	3.82	8.16	4.76	7.67	2.91	
Bargaining	€)	7.72	4.16	8.56	5.11	8.18	4.25_{-}	_7.36_	2.81	
T		-1.89	1.89 1.53	0.16	-0.16	-1.39	1.39	0.85	-0.85	
Leduc poker		-1.53	1.53	-0.60	0.60	-2.57	2.57	2.16	-2.16	

D.2.3 HIDDEN STATE INFERENCE

Table 21: PPO-CWM win rates using CWM refinement via tree search and hidden state inference against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Þ	C	WM MCTS		Ge	mini 2.5 Pro		GT-ISMCTS			Random		
Game	1	Win (forfeit/n)	Loss (forfeit/n)	Draw (n)	Win (forfeit/n)	Loss (forfeit/n)	Draw (n)	Win	Loss	Draw	Win	Loss	Draw
01		0.00 (0/0)	0.63 (0/63)	0.37 (37)	0.76 (0/76)	0.17 (0/17)	0.07 (7)	0.00	0.60	0.40	0.56	0.10	0.34
Quadranto		0.00 (0/0)	0.88 (0/88)	0.12 (12)	0.62 (0/62)	0.34 (0/34)	0.04 (4)	0.00	0.88	0.12	0.66	0.19	0.15
II 1	•	0.63 (0/63)	0.26 (0/26)	0.11 (11)	0.34 (2/34)	0.45 (0/45)	0.21 (21)	0.19	0.69	0.12	0.67	0.22	0.11
Hand of war	0.54 (0/54)	0.38 (0/38)	0.08 (8)	0.50 (0/50)	0.31 (0/31)	0.19 (19)	0.49	0.45	0.06	0.58	0.25	0.17	

Table 22: PPO-CWM payoffs using CWM refinement via tree search and hidden state inference against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Þ	CWM	MCTS	Gemini 2.5 Pro		GT-IS	MCTS	Random		
Game	1	Us	Them		Them	Us	Them	Us	Them	
Bargaining	\$	8.17 7.62	4.61 4.51	8.51 8.46	3.83 4.78	8.10 7.76	4.48 4.45	7.87 8.05	3.31 2.51	
Leduc poker	•	-1.31 -2.01	1.31 2.01	-0.02 1.67	0.02 -1.67	-2.17 -2.02	2.17 2.02	0.58 2.11	-0.58 -2.11	

D.2.4 HIDDEN HISTORY INFERENCE WITH CLOSED DECK LEARNING

Table 23: PPO-CWM win rates using CWM refinement via tree search with closed deck against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Þ	С	WM MCTS		Ge	mini 2.5 Pro		GT	-ISMC	TS	Random		
Game	1	Win (forfeit/n)	Loss (forfeit/n)	Draw (n)	Win (forfeit/n)	Loss (forfeit/n)	Draw (n)	Win	Loss	Draw	Win	Loss	Draw
0 1		0.04 (0/4)	0.72 (0/72)	0.24 (24)	0.87 (0/87)	0.11 (0/11)	0.02 (2)	0.00	0.58	0.42	0.41	0.19	0.40
Quadranto		0.05 (0/5)	0.37 (0/37)	0.58 (58)	0.66 (0/66)	0.11 (0/11)	0.23 (23)	0.05	0.46	0.49	0.41	0.17	0.42
II 1 C		1.00 (100/100)	0.00 (0/0)	0.00 (0)	0.31 (1/31)	0.57 (0/57)	0.12 (12)	0.38	0.50	0.12	0.63	0.19	0.18
Hand of war	*	1.00 (100/100)	0.00 (0/0)	0.00 (0)	0.54 (0/54)	0.31 (0/31)	0.15 (15)	0.43	0.47	0.10	0.65	0.16	0.19

Table 24: PPO-CWM payoffs using CWM refinement via tree search with closed deck against multiple opponents. For each game, results in the first (second) row correspond to our agent going first (second).

Game	Р	CWM	MCTS	Gemini 2.5 Pro		GT-IS	MCTS	Random		
Game	1	Us	Them	Us	Them	Us	Them	Us	Them	
Danmaining	(\$	5.88	6.04	7.20	4.91	5.94	6.04	7.23	3.81	
Bargaining	€ D	6.77	5.61	_ 7.76_	5.66	_ 6.97	5.25	_7.82_	3.48	
Leduc poker	•	-1.17	$1.17 \\ 0.25$	-0.47	0.47	-1.35	1.35	0.70	-0.70	
Leduc poker		-0.25	0.25	-0.48	0.48	-1.54	1.54	2.17	-2.17	

1459 1460

1461

1462

1463

1464

1465

1466

1467 1468

1469

1470

1471

1472

1473 1474 1475

1476 1477

1478

1479 1480

1481

1482

1483 1484

1485

1486

1487

1488 1489

1490

1491

1492 1493

1494

1495

1496 1497 1498

1499

1501

1502 1503 1504

1505

1506

1507

1508

1509 1510 1511

AUTOMATIC REJECTION OF BAD CWM SAMPLES

The CWM refinement process can occasionally produce a low-quality CWM. This is rarely the case for perfect information games, where more information is available for refinement and unit tests are more strict, but we have observed this happening in the case of imperfect information games. To reduce this effect, in the case of imperfect information games, we sample 5 CWMs, create a CWM-ISMCTS agent from each one, and make those agents compete against each other. Agents are then ranked according to the average payoff obtained in those competitions. Agents that are worse than the best scoring agent by more than 10% of the observed utility range are rejected.

Since we do not have access to the ground truth game for these competitions, the agents use the CWM of one of them as a stand-in for the actual game. We call the CWM used to play the game the host. This means that we have 2 possible hosts \times 5 agents acting as Player 0×5 agents acting as Player 1. This results in a total of 50 possible matches. Since the outcome of a match is stochastic, we repeat each match 50 times. Execution failures or the execution of illegal actions during these games result in both players losing the game.

\mathbf{F} SKETCH OF INFORMATION FLOW OF EACH AGENT

Here we provide a sketch of the information flow for each the agents. Of course, many details are omitted, and the prompts are highly simplified, see Appendix G for the actual prompts.

```
def llm_agent_generator(LLM, rules, traj):
         prompt = (f"You are playing a game with these rules: {rules}.\n"
                   f"Example trajectories: {traj}.\n")
        def policy(action_obs_history):
           return LLM(prompt + f"Action-observation history: {action_obs_history}. "
                      "Pick the next best action.")
         return policy
       def cwm_agent_perfect_info_generator(LLM, rules, traj, GT=False):
        M = induce_cwm(LLM, rules, traj) if not GT else ground_trutn_M
        V = induce_value_fn(LLM, rules, traj, M)
        def policy(action_obs_history):
           return MCTS(action_obs_history[-1], M, V)
        return policy
       def cwm_agent_imperfect_info_generator(LLM, rules, traj, GT=False):
         (M, I) = induce_cwm_pomdp(LLM, rules, traj) if not GT else ground_trutn_MI
1500
        V = induce_value_fn(LLM, rules, traj, M)
        def policy(action_obs_history): return ISMCTS(action_obs_history, M, V, I)
        return policy
      def induce_cwm_zero_shot(LLM, rules, traj):
         prompt = (f"You are playing a game with these rules: {rules}.\n"
                   f"Generate python code that matches this API: {fn_signature}\n"
                   f"The code should pass these unit tests: {make_tests(traj)}\n")
         return LLM(prompt)
```

```
1512
             System and agent prompts
1513
1514
       G.1
             Tree Search
1515
1516
       Our tree search prompt is:
1517
        You are an expert python programmer who is building the game of {game_name}.
1518
       Here is a description of the game:
1519
       {game_desc}
1520
1521
       The goal is to implement a python function with the following signature.
1522
       # START FUNCTION SIGNATURE
1523
       {function_signature}
1524
       # END FUNCTION SIGNATURE
1525
1526
       The original implementation is as follow. Please try to refine the original code.
1527
       # START CODE BLOCK
1528
       {orig_code}
1529
       # END CODE BLOCK
1530
1531
       Your code should satisfy the following unit tests.
1532
       Your code should fix the TODO errors in the comments of the unit tests, if any.
1533
       # START UNIT TESTS
1534
       {test_code}
       # END UNIT TESTS
1535
1536
1537
       Do not repeat the unit tests, only return the functions.
1538
       Do not leave placeholders.
1539
       Do not repeat the function signature.
1540
       Do not copy the unit tests.
1541
1542
       Only produce code that is compact.
1543
       Do write comments explaining what the code does.
1544
       Do use helper functions to reduce code duplication.
1545
       Start by reasoning about the game and the unit tests.
1546
       Also reason about the errors and possible fixes.
1547
1548
       Finally, try to write {num_targets} versions of the code.
       Make sure each code is in a different code blocks starting with ```python.
1549
1550
        function_signature contains the function definition for the LLM to fill out, while test_code
1551
       defines the properties (expressed as unit tests) that the resulting code needs to satisfy.
1552
1553
        function_signature and test_code both depend on if the game is a perfect or imperfect
       information game, whether it is being learnt in an open or closed deck fashion, and if the
       inference is perform via hidden history or hidden state inference. These variations are
1555
       defined in the following sections.
1556
1557
        Finally, orig_code is the code being refined at each iteration. On the first iteration, this
1558
       paragraph is not present.
1559
1560
        G.2 Perfect information games
1561
1562
       function_signature is defined as follows:
1563
       Action: str
1564
        State: dict[str, Any]
1565
       PlayerObservation: dict[str, Any]
```

1619

player_id = {player_id}

```
def apply_action(state: State, action: Action) -> State:
1567
                """Returns the new state after an action has been taken."""
1568
             def get_current_player(state: State) -> int:
1569
                """Returns current player, with -1 for chance and -4 for terminal."""
1570
1571
             def get_player_name(player_id: int) -> str:
1572
                """Returns the name of the player, with 'chance' for -1, and 'terminal' for -4."""
1574
             def get_rewards(state: State) -> list[float]:
                """Returns the rewards per player from their last action."""
1575
1576
             def get_legal_actions(state: State) -> list[Action]:
1577
                """Returns legal actions that can be taken in current state."""
1578
             def get_observations(state: State) -> list[PlayerObservation]:
1579
                """Returns the observation for player.""
1580
1581
             test_code tests the transition between two states, testing each of the API calls defined in
1582
             function_signature. Here is an example transition unit test for tic tac toe, where the board
1583
             is provided as a flat 1D array:
1584
             class TestTransition2(unittest.TestCase):
1585
                def test_transition_2(self):
1586
                   state = {'board': [None, None, None, None, 'x', None, 'o', None, None], '
1587
                          current_player_mark': 'x'}
1588
1589
                   self.assertEqual(0, get_current_player(state))
                   {\tt self.assertEqual('0', get\_player\_name(0))}
1590
                   self.assertEqual([0.0, 0.0], get_rewards(state))
1591
                   self.assertEqual([{'board': [None, None, None, None, 'x', None, 'o', None, None], '
1592
                          current_player_mark': 'x'}, {'board': [None, None, None, 'x', None, 'o',
1593
                          None, None], 'current_player_mark': 'x'}], get_observations(state))
1594
                   self.assertSetEqual(set(['x(0,0)', 'x(0,1)', 'x(0,2)', 'x(1,0)', 'x(1,2)', 'x(2,1)',
                             'x(2,2)']), set(get_legal_actions(state)))
1595
                   {\tt self.assertEqual(\{'board'\colon [None,\ None,\ 'x',\ 'x',\ None,\ 'o',\ None,\ None],\ 'x',\ None,\ 'o',\ None,\ None],\ 'x',\ None,\ 'y',\ None,\ None,\ None,\ None],\ 'y',\ None,\ No
1596
                          current_player_mark': 'o'}, apply_action(state, 'x(1,0)'))
1597
1598
             If this test has failed, the LLM is provided with the python error message in the form of
1599
             a comment before the test. The use of self.assertEqual style functions ensures that the
             LLM is provided with a rich description of how the expected and actual data structures
             vary.
1602
1603
             G.3 HIDDEN HISTORY INFERENCE FUNCTION SYNTHESIS, OPEN DECK
1604
1605
             function_signature starts with the version from Section G.1, then adds the inference def-
             inition:
1606
             def resample_history(obs_action_history: list[tuple[PlayerObservation, Action | None]],
                    player_id: int) -> list[Action]:
                """Stochastically sample one of many potential history of actions for all players(
1609
                       including 'chance' and 'terminal')
1610
1611
               This is given only a single player's observations and actions, and needs to recreate
1612
                       the player_id's observations
1613
1614
             unit_text again starts with the definition from Section G.1 then adds the following test for
1615
            added inference function:
1616
             state = INITIAL_STATE
1617
             obs_action_history = {obs_action_history}
1618
             obs_and_action_iter = iter(obs_action_history)
```

current_player_obs, current_player_action = next(obs_and_action_iter)

```
for action in resample_history(obs_action_history, player_id):
1621
         print(f"In state {{state}}}")
1622
         if get_current_player(state) == player_id:
           self.assertEqual(current_player_obs, get_observations(state)[player_id])
1623
           print(f"Recreated observation {{current_player_obs}}")
1624
           self.assertEqual(current_player_action, action)
1625
           current_player_obs, current_player_action = next(obs_and_action_iter)
1626
         print(f"Taking action {{action}}")
         state = apply_action(state, action)
1628
1629
         next(obs_and_action_iter)
1630
         raise ValueError('Failed to iterate through all observations.')
        except StopIteration:
        self.assertEqual(player_id, get_current_player(state))
1633
1634
```

 where INITIAL_STATE is provided at the beginning of the unit tests and is the static first state of the game. obs_action_history is the history of observations and actions for player player_id for which we want to resample the history of actions that lead to the current observations.

Note the presence of print statements inside the unit test. The last ten lines of standard output are provided to the LLM in addition to the error message.

G.4 HIDDEN STATE INFERENCE FUNCTION SYNTHESIS

function_signature again starts with the version from Section G.1, then adds the inference function definition:

```
def resample_state(obs_action_history: list[tuple[PlayerObservation, Action | None]],
    player_id: int) -> list[int]:
"""Stochastically sample one of the reachable statess for player given the observation
    and action history that recreates the player's observation."""
```

unit_test again starts with the definition from Section G.1 then adds the following test for added inference function above:

```
obs_action_history = {obs_action_history}
player_id = {player_id}
resampled_state = resample_state(obs_action_history, player_id)
self.assertEqual(obs_action_history[-1][0], get_observations(resampled_state)[player_id])
```

G.5 HIDDEN HISTORY INFERENCE FUNCTION SYNTHESIS, CLOSED DECK

```
function_signature is similar to that in Section G.2:
```

```
def resample_history(obs_action_history: list[tuple[PlayerObservation, Action | None]],
    player_id: int, last_is_terminal: bool) -> list[Action]:
    """Stochastically sample one of many potential histories of actions for all players(
        including 'chance' and 'terminal')
    given only a single player's observations and actions.

It needs to recreate the player_id's observations.
last_is_terminal indicates if the last player observation is from end of game when
        player_id is -4."""
```

Note the extra argument last_is_terminal. This indicates that the final observation in obs_action_history is of the terminal state. This allows adding tests that resample the entire game from the beginning to the terminal state, testing the ability of the LLM to

1675

1676

1677

1697

1699

1700

1701

1702

1724

17251726

1727

predict the final reward of the player. In open deck, the transition tests cover this. For simplicity, we assumed that the rewards are terminal but this is easy to relax.

The corresponding unit_tests for the inference function is:

```
1678
        state = INITIAL_STATE
1679
        obs_action_history = {obs_action_history}
1680
       player_id = {player_id}
        last_is_terminal = {ends_in_terminal}
       obs_and_action_iter = iter(obs_action_history)
1682
       current_player_obs, current_player_action = next(obs_and_action_iter)
1683
        for action in resample_history(obs_action_history, player_id, last_is_terminal):
1684
         print(f"In state {{state}}")
1685
         if get_current_player(state) == player_id:
1686
           self.assertEqual(current_player_obs, get_observations(state)[player_id])
           print(f"Recreated observation {{current_player_obs}}")
1687
           self.assertEqual(current_player_action, action)
1688
           current_player_obs, current_player_action = next(obs_and_action_iter)
1689
         print(f"Taking action {{action}}")
         state = apply_action(state, action)
         next(obs_and_action_iter)
1693
         raise ValueError('Failed to iterate through all observations.')
       except StopIteration:
1695
         pass
1696
```

Again, this is very similar to unit_test in Section G.2, but also covers the terminal state of the game and it's associated reward.

Note that no transition unit tests are added as we do not have access to the state. However, just testing the inference function is not enough to ensure that the resulting closed deck game is playable. Instead, a random play test is added to unit_test:

```
1703
       state = {initial_state}
1704
       rg = np.random.RandomState({seed})
1705
        for it in range(1000): # upper bound on game length
1706
         current_player = get_current_player(state)
         rewards = get_rewards(state)
1707
         assert len(rewards) == 2
1708
         print (f"State is {{state}}, current player is {{current_player}}, rewards are {{
1709
              rewards}}")
1710
         if current_player == -4: # Game over
1711
           break
1712
         if current_player in [0,1]: # Real players
1713
           print(f"Observation for current player is {{get_observations(state)[current_player
1714
                ]}}")
1715
         else:
1716
           assert current_player == -1
1717
         legal_actions = get_legal_actions(state)
1718
         chosen_action = rg.choice(legal_actions)
1719
         print(f"Taking action {{repr(chosen_action)}} from {{len(legal_actions)}} options,
1720
              first 10 are {{[*legal_actions][:10]}}")
1721
         state = apply_action(state, chosen_action)
       else:
1722
         raise ValueError(f"Game did not end after 1000 steps.")
1723
```

This tests that if every player randomly picks a valid move, the game will correctly play and terminate. Note that we assume access to the static and deterministic initial state of the game, before any chance nodes have taken place. This could also be synthesized by the LLM instead.

```
1728
        G.6
             RESAMPLING THE STATE AT GAME PLAYING TIME FOR IMPERFECT INFORMATION
1729
              GAMES
1730
        When playing the game, we allow the system to up to 10 tries to get a valid state that
1731
        produces the current observations:
1732
1733
1734
       for retry in range(10):
         json_state = {start_state}
1735
         try:
1736
           actions = resample_history(obs_action_history, player_id)
1737
           for action in actions:
1738
             json_state = apply_action(json_state, action)
1739
             state_log.append(json_state)
         except Exception as e: # Running generated code, could raise anything.
1740
           continue
1741
         recreated_obs = get_observations(json_state)[player_id]
1742
         if recreated_obs == obs_action_history[-1][0]:
1743
           return json_state
1744
1745
       Additionally, if the ISMCTS process fails due to, e.g., poor understanding of the game
1746
        termination criteria in the CWM, we fall back to resampling the state and then return a
       uniformly sampled legal action from that state.
1747
1748
        G.7 Value function synthesis
1749
1750
       Our value function synthesis function prompt is
1751
1752
       You are an expert python programmer. You are playing the game {game}, and need
1753
        to synthesize a value function for monte carlo tree search.
1754
1755
       {game_description}
1756
       For reference, the game is implemented as follow
1757
1758
       {code}
1759
1760
       The function you need to write is:
1761
       {value_function}
1762
       It should return the reward at terminal states, and otherwise an estimate of the
1763
       value for each non-terminal states.
1764
1765
       It should always be a float:
1766
       {player_tests}
1767
       Terminal states should match rewards:
1768
       {terminal_tests}
1769
1770
       To write a good value function first reason about the game and produce a heuristic value
1771
             that is informative, and do not just output zeros everywhere other than terminal
1772
            states.
       Finally ONLY output the new value_function, do not output any other text, code,
1773
       explanations or placeholders.
1774
       The response code must be a single CODE BLOCK that uses this format:
1775
       The opening fence: ```python The closing fence: ```
1776
1777
1778
       Where value_function is
1779
1780
       def value_function(state: dict[str, Any], player_id: int) -> float:
1781
```

"""Returns the value estimate for player_id in state.

```
1782
1783
         For terminal states the function returns the true return. For ongoing play
         the function should return a value estimate that reflect the winning potential
1784
1785
         of the player with given player_id.
1786
        111
1787
1788
        player_tests and terminal_tests is the list of example
1789
        111
1790
       {current_player}
1791
        self.assertIsInstance(value_function(state, {current_player}), float)
1792
        if {current_player} == pyspiel.PlayerId.TERMINAL:
1793
           rewards = get_rewards(state)
           for player in range(len(rewards)):
1794
               self.assertEqual(rewards[player], value_function(state, player))
1795
1796
1797
1798
1799
1800
1801
1802
1803
1804
1805
1806
1807
1808
1809
1810
1811
1812
1813
1814
1815
1816
1817
1818
1819
1820
1821
1822
1823
1824
1825
1826
1827
1828
1829
1830
1831
1832
1833
1834
1835
```

1836 Η GAME RULES 1837 1838 H.1 Backgammon 1839 1840 Backgammon is a two-player board game that combines strategy and luck. The object is to 1841 move all of your checkers off the board before your opponent does. Here's a 1842 breakdown of the rules: 1843 1844 **The Board and Setup: ** 1845 * **The Board:** The board consists of 24 narrow triangles called **points**. These points are grouped into four quadrants of six points each: 1847 * **Inner/Home Board:** The quadrant closest to each player's starting position. 1848 * **Outer Board:** The quadrant further from each player's starting position. * **The Bar: ** The area in the middle of the board, separating the two sides. 1849 * **The Bear-Off Area: ** The area off the board where checkers are moved once they reach 1850 the player's home board. 1851 * **Checkers: ** Each player has 15 checkers of one color (typically black and white). 1852 * **Dice:** Two dice are used to determine movement. 1853 * **Doubling Cube (Optional but common):** A cube with the numbers 2, 4, 8, 16, 32, and 1854 64, used to increase the stakes of the game. 1855 **Initial Setup:** 1856 1857 Each player's 15 checkers are set up in a specific configuration on the points: 1858 1859 * 2 checkers on the opponent's 24-point. * 5 checkers on the opponent's 13-point. 1860 * 3 checkers on their own 8-point. 1861 * 5 checkers on their own 6-point. 1862 1863 **Gameplay:** 1864 1. **Starting the Game:** Each player rolls one die. The player with the higher roll 1865 goes first. If the rolls are the same, they roll again until one player rolls higher. The player who goes first uses the numbers rolled on *both* dice to make 1867 their first move. 1868 2. **Rolling the Dice:** On subsequent turns, each player rolls two dice. 1869 1870 **Moving Checkers:** After rolling the dice, the player must move their checkers 1871 according to the numbers rolled. 1872 * **Separated Moves:** Each die represents a separate move. You can move one checker 1873 the distance of one die's roll and another checker the distance of the other die's roll. 1874 * **Combined Move:** You can move one checker the combined distance of both dice 1875 rolls, but *only if* the point you would land on for the first die's roll is not 1876 blocked (see "Blocked Points" below). 1877 * **Mandatory Moves:** You must move your checkers if possible. If you can only make one of the two moves indicated by the dice, you must make that move. If you can 1879 make both, you must make both. * **No Legal Moves: ** If you cannot make any legal moves based on the dice roll, 1880 your turn ends. 1881 1882 4. **Point Direction:** You always move your checkers from your opponent's inner board 1883 towards your own home board. The points are numbered 1 to 24, where 24 is the latest point in your opponent's inner board. Each move makes the checker move to 1884 smaller numbered points. 1885 1886 5. **Landing on a Point:** 1887 * **Empty Point:** You can land on an empty point. 1888 * **Point Occupied by Your Own Checkers:** You can land on a point occupied by any 1889 number of your own checkers.

* **Point Occupied by Opponent's Checkers:**

```
1890
               * **Blots:** If a point is occupied by *only one* of your opponent's checkers, it
1891
                    's called a "blot." If you land on a blot, you "hit" the opponent's checker.
1892
                    The hit checker is placed on the **bar**.
               * **Blocked Points:** If a point is occupied by *two or more* of your opponent's
1893
                   checkers, it is "blocked." You *cannot* land on a blocked point.
1894
1895
        6. **Entering from the Bar:** If a player has checkers on the bar, they must re-enter
1896
            them onto the board before making any other moves.
           * **Re-entry Points:** You can re-enter a checker from the bar onto a point in your
1898
                opponent's home board that corresponds to the number rolled on a die. For
                example, if you roll a 3, you can re-enter a checker onto your opponent's 3-
1899
                point.
1900
```

- * **Blocked Re-entry:** If the corresponding point in your opponent's home board is blocked by two or more of your opponent's checkers, you cannot re-enter using that die roll.
- * **Priority:** You must use any available die rolls to re-enter checkers from the bar. If you can re-enter one checker but not the other based on your dice roll, you must still re-enter the one you can. If you cannot re-enter any checkers, your turn ends.
- 7. **Doubles:** If you roll doubles (e.g., two 4s), you can use each number *four times *. So, two 4s means you have four moves of 4. You can use these moves in any combination, as long as they are legal. Each turn allows you to make two moves only . So if a player rolls a double, they take an extra turn, making at most two moves in each of the turns.
- 8. **Bearing Off:** If you don't have any checkers outside of your home board or at the bar, you can "bear off" chekers (moving them off the board) that are at your home board.
 - * **Bearing Off Rolls:** To bear off a checker, you must roll the exact number that the checker is on to move it off the board from its current point. For example, if a checker is on your 4-point, you need to roll a 4 to bear it off.
 - * **Higher Rolls:** If you roll a number higher than the highest point occupied by your pieces, you can still bear off a piece. However, you must bear off a piece from the highest occupied point. For example, if your highest occupied point is the 4-point, and you roll a 6, you can bear off a piece only from the 4-point, you cannot bear off a piece from lower points.
 - * **Lower Rolls:** If you roll a number lower than the point your checker is on, you
 can still move a checker from a higher point the distance of the roll (if legal
), or you must move a checker from a lower point the distance of the roll if
 possible. You cannot bear off a checker if you have checkers on higher points in
 your home board that can be moved by the dice roll.
 - * **Blocked Bear Off:** You cannot bear off a checker if any of your checkers are still on the bar or outside of your home board. You must bring all your checkers into your home board before bearing off.
 - * **Moving Pieces Within Home Board:** Instead of bearing off, you can also move your checkers within your home board using the dice rolls.
- 9. **Action Notation:** Player moves are typically represented using a specific notation . Each turn consists of at most two moves. Each move isrepresented by a string of the form "Move checker at X using Y roll", where "X" is the position of the checker being moved, and "Y" indicates if the move is done based on the dice roll with the higher or lower number.
 - The first position is always a number between 1 and 24 or "Bar", and it is presented in each player's perspective, where 24 is the latest point in the opponent's inner board.
 - For a double roll (e.g., a 3-3, granting four moves of 3 as per rule 7 "Doubles")
 - It is assumed these four moves are made by the player in two stages: First, providing two moves, then the player gets a second turn then providing the subsequent two moves. For instance, a player might move two checkers, each making two 3-point moves. Next, the player would have a second turn, providing two more moves of 3.
- 10. **Board Notation:** The board is represented with 2 arrays of 24 numbers,

```
1944
               where each number is either 0 (empty) or a number between 1 and 15
1945
               (indicating the number of checkers of that color on that point).
1946
               The first array is for the first player, and the second is for the second.
               The board ordering is from second player's persective. Starts from first
1947
               player's home base and ends at the second player's home base. The first
1948
               index is the latest point in the second player's inner board, meaning
1949
               position 24 for the first player and 1 for the second player.
1950
1951
       **Winning the Game: **
1952
       The first player to bear off all 15 of their checkers wins the game.
1953
1954
        **Optional Rules (Commonly Used):**
1955
1956
        * **The Doubling Cube:**
           * **Offering a Double:** At the start of their turn, *before* rolling the dice, a
1957
                player can offer to "double" the stakes of the game.
1958
           * **Accepting a Double:** The opponent can either accept or decline the double. If
1959
                they decline, they lose the game immediately and the current stake is paid. If
1960
                they accept, the stakes are doubled, and the opponent now "owns" the doubling
1961
                cube, meaning they are the only one who can offer the next double.
           * **Subsequent Doubles:** The owner of the cube can offer to redouble at the start
1962
                of their turn. The stakes continue to double with each accepted redouble (2, 4,
1963
                8, 16, etc.).
1964
        * **Gammon and Backgammon:** These are ways to win with higher stakes.
1965
           * **Gammon:** If a player bears off all their checkers before the opponent has borne
1966
                 off *any* checkers, the winner wins a "gammon," which is typically worth double
1967
                 the value of the doubling cube.
           * **Backgammon: ** If a player bears off all their checkers before the opponent has
1968
                borne off *any* checkers and the opponent still has one or more checkers on the
1969
                bar or in the winner's home board, the winner wins a "backgammon," which is
1970
                typically worth triple the value of the doubling cube.
1971
       **Key Concepts and Strategy:**
1972
1973
        * **Hitting Blots:** Hitting your opponent's checkers puts them on the bar and disrupts
1974
            their progress.
1975
        * **Making Points:** Occupying points with two or more of your checkers creates "blocks"
1976
             that prevent your opponent from moving past. Strategic point-making is crucial.
        * **Prime: ** Creating a "prime" (six consecutive blocked points) can severely hinder
1977
            your opponent's movement.
1978
        * **Running: ** Moving your checkers quickly towards your home board.
1979
        * **Positioning:** Carefully considering where to move your checkers to maximize your
1980
            options and limit your opponent's.
1981
        * **Risk vs. Reward:** Balancing the risk of leaving blots with the potential for making
1982
             good moves.
1983
       Backgammon is a game with layers of strategy that unfold as you play. While the dice
1984
            introduce an element of chance, skillful play, understanding probability, and
1985
            strategic decision-making significantly influence the outcome. Enjoy the game!
1987
       H.2 Connect four
1988
1989
        ### Rules of Connect Four
1990
1991
           **Setup:** Connect Four is played on a 6-row by 7-column vertical grid, which starts
1992
             completely empty.
           **Players and Marks:** There are two players: Player 0 uses the 'x' mark and Player
1993
            1 uses the 'o' mark.
1994
           **Turns:** Player 0 ('x') always goes first, and turns alternate between players.
1995
           **Making a Move:** On your turn, you choose a column to drop your mark into. The
1996
            mark will fall to the lowest unoccupied square within that chosen column.
1997
            Attempting to drop a mark into a column that is already full is an invalid move;
            you must choose a column with at least one empty square to complete your turn.
```

- 1998 **Winning the Game:** The winner is the first player to get four of their marks in a 1999 row (horizontally, vertically, or diagonally). The game ends immediately as soon 2000 as a winning line is formed. 2001 **Drawing the Game: ** If all 42 squares on the grid are filled and neither player has won, the game ends in a draw. 2002 **End the Game: ** The game only concludes upon a win or a draw. A player must make a 2003 move on their turn as long as there is at least one valid move available on the 2004 board. 2005 **Move Notation:** Use the move notation '[mark][col]', where col is the 0-indexed column you are dropping your mark into. For example, 'x3' means Player 0 ('x') 2006 drops their mark into the fourth column from the left (column index 3). 2007 2008 2009 H.3 TIC-TAC-TOE 2010 ### Rules of Tic-Tac-Toe 2011 2012 **Setup: ** Tic-Tac-Toe is played on a 3x3 grid, which starts completely empty. 2013 **Players and Marks:** There are two players: Player 0 uses the 'x' mark and Player 2014 1 uses the 'o' mark. 2015 **Turns:** Player 0 ('x') always goes first, and turns alternate between players. 2016 **Making a Move:** On your turn, you must place your mark in a single, unoccupied square. Attempting to place a mark in an already occupied square is an invalid move 2017 ; you must choose an empty square to complete your turn. 2018 **Winning the Game:** The winner is the first player to get three of their marks in 2019 a row (horizontally, vertically, or diagonally). The game ends immediately as soon 2020 as a winning line is formed. 2021 **Drawing the Game:** If all nine squares on the grid are filled and neither player has won, the game ends in a draw. 2022 **End the Game:** The game only concludes upon a win or a draw. A player must make a 2023 move on their turn as long as there is at least one valid move available on the 2024 2025 **Move Notation:** Use the move notation 'mark(row,col)', where row and col are 0indexed. For example, 'x(0,0)' means Player 0 ('x') places their mark in the top-2026 left square. 2027 2028 2029 H.4 GEN. TIC-TAC-TOE 2030 Generalized Tic-Tac-Toe (6x6, Win Length 4, 2 Players) 2031 2032 1. Overview: This is a two-player strategy game played on a 6x6 grid. The goal 2033 is to be the first player to form a continuous line of four of your own marks. 2034 This game is a specific configuration of a generalized Tic-Tac-Toe framework. 2035 2036 2. Game Setup: 2037 Board: A 6x6 grid of cells (36 cells in total), with rows and columns numbered 2038 2039 Players: Two players. Conventionally, one player uses 'x' and the other uses 'o'. 2040 Starting State: The board is initially empty. 2041 3. Gameplay: 2042 2043 Players take turns placing their mark on an unoccupied cell on the board. 2044 A designated player (e.g., Player 'x') makes the first move. 2045 The game continues with players alternating turns.
 - 4. Winning Condition:

A player wins if they are the first to place four of their marks in an unbroken straight line.

2051 This line can be:

2046

2047 2048

2049

```
2052
       * Horizontal: Four marks in the same row.
2053
        * Vertical: Four marks in the same column.
2054
        * Diagonal: Four marks along any of the board's diagonal lines (both directions).
2055
        5. Draw Condition:
2056
2057
       If all cells on the 6x6 board are filled with marks, and neither player has
2058
       achieved a line of four of their marks, the game is a draw.
2059
2060
       6. End of Game:
2061
       The game concludes immediately when either:
2062
       One player achieves a winning line of four marks (that player is the winner).
2063
       All cells are filled, and no winning line exists (the game is a draw).
2064
       7. Key Parameters for this Specific Variant:
2065
2066
       Number of Rows: 6
2067
       Number of Columns: 6
2068
       Winning Line Length: 4
2069
       Number of Players: 2
2070
2071
       H.5 Gen. Chess
2072
2073
2074
       The game of generalized chess is a two player game where each player controls a
2075
       collection of pieces and wins by capturing the target piece from the other
2076
       player. Each kind of game piece has a specific pattern of movements that it can
2077
       execute. A piece can execute any one of its available moves as long as that move
        stays on the board and doesn't land on another of that player's pieces. If the
2078
       piece lands on an opponent piece, it captures the opponent piece and removes it
2079
       from the board. Allowed piece movements are not the same as in standard chess.
2080
2081
       Actions are described using board coordinates. For a 5x5 board, rows are labeled
2082
       A-E from top to bottom, and columns are labeled 1-5 from left to right. A move
        from a starting square to a destination square is written as 'start_to_end',
2083
       for example, 'A2_to_C2' means move the piece from square A2 to square C2.
2084
2085
       Passing a turn is specified as 'PASS'.
2086
2087
       This 'army5x5a' variant of generalized chess is played on a 5x5 board.
2088
       It includes the following pieces, with their corresponding set of allowed moves:
2089
         - general: [(1, 0), (-1, 0), (0, 1), (0, -1), (0, -2), (0, 2)]
2090
         - infantry: [(1, 0), (2, 0), (1, -1), (1, 1), (-1, 0)]
2091
        - cavalry: [(0, 3), (1, 2), (2, 1), (3, 0)]
2092
       Game pieces are depicted with the following symbols: 'general': 'X', 'infantry': 'I', '
2093
            cavalry': 'V'. Player 0 pieces are upper-case while Player 1 pieces are lower-case.
2094
2095
       The 'general' is the target piece. Capturing this piece wins the game.
2096
2097
2098
       H.6 Bargaining
       The rules of "bargaining" aren't fixed and formal like a board game with a rulebook.
2100
            Instead, it's a dynamic social process of negotiation where two or more parties
2101
            attempt to reach a mutually agreeable outcome on a price or terms for a product,
2102
            service, or agreement. Here's a breakdown of the core principles and common "rules"
2103
             of bargaining, understood more as strategies and expectations:
2104
2105
       **Core Principles of Bargaining:**
```

- * **Mutual Desire for an Agreement:** Both parties generally want to reach a deal, even
 if their initial positions are far apart.
 - * **Information Asymmetry:** One party often has more information than the other, which can influence the negotiation.
 - * **Iterative Process:** Bargaining usually involves a series of offers and counteroffers.
 - * **Focus on Value:** Bargaining is about perceived value what each party believes the item or service is worth.
 - * **Potential for Compromise:** Both parties are usually expected to give a little to reach an agreement.
 - **Implicit "Rules" or Common Strategies:**

2108

2109

2110

2111

2112

2113

2114

2115

2116

2119

2131

2132

2133

2134

2138

2139

2140

2141

- These are not hard-and-fast rules, but rather common practices and expectations that guide the negotiation:
- **Know Your Limits (Walk-Away Point):** Before starting, each party should have a clear idea of the maximum (for a buyer) or minimum (for a seller) price they are willing to accept. This is your "reservation point."
- * **Buyer's Perspective:** Start lower than what you're willing to pay.
- 3. **Justify Your Offers:** Simply stating a price is less effective than explaining *
 why* you're offering that price. Reference market value, condition of the item,
 your budget, etc.
 - 4. **Make Concessions Incrementally:** Don't jump straight to your walk-away point. Make small concessions with each counter-offer. This signals a willingness to negotiate while still trying to get the best possible deal.
- 5. **Signal Willingness to Walk Away (But Don't Bluff Too Much):** Letting the other party know you're willing to walk away if you don't get a satisfactory price can be a powerful tactic. However, repeated or unbelievable threats can undermine your credibility.
 - 6. **Listen Actively and Ask Questions:** Pay attention to the other party's offers, reasoning, and potential underlying needs. Asking questions can reveal information and build rapport.
- 7. **Be Patient:** Bargaining takes time. Don't rush the process.
- 2144 8. **Maintain a Respectful Tone:** Even if the negotiation becomes difficult, try to
 2145 maintain a polite and respectful demeanor. Aggression can shut down the
 2146 conversation.
- 9. **Consider Non-Price Factors:** While price is central, bargaining can also involve other terms like delivery time, payment method, warranties, or additional items included.
- 2150
 2151
 10. **Know When to Stop:** If it's clear you won't reach an agreement that meets your needs, it's okay to respectfully end the negotiation.
- 2153 11. **Be Prepared to Walk Away:** If you can't reach an agreement within your limits, you must be prepared to walk away. This is crucial for maintaining your boundaries.
- 2155
 2156
 2157
 12. **The Final Offer:** Often, one party will indicate their "final offer." This suggests they are unwilling to make further concessions. However, this isn't always truly final and can be tested with a counter-offer.
- 2159 13. **The Art of the Counter-Offer:** Respond to offers with a counter-offer that is a concession from your previous position, but still moves you closer to your goal.

```
2160
2161
        **Situational Differences:**
2162
2163
        The "rules" of bargaining can vary depending on the context:
2164
        * **Cultural Norms:** Bargaining is much more common and expected in some cultures (e.g.
2165
            ., bazaars in many parts of the world) than others (e.g., retail stores in most
2166
            Western countries).
2167
        * **Type of Item/Service:** Bargaining for a car is different than bargaining for a
2168
            small trinket at a market.
2169
        * **Power Dynamics:** Who has more leverage in the negotiation can significantly impact
            the process.
2170
2171
        **In summary, the "rules" of bargaining are less about strict regulations and more about
2172
             strategic communication, understanding the other party's perspective, and being
            prepared to make concessions to reach a mutually acceptable agreement. It's a
2173
            negotiation dance where both parties are trying to get the best possible outcome
2174
            within their own limits.**
2175
2176
2177
        H.7 Leduc Poker
2178
2179
        Leduc Poker is a simplified two-player poker game, ideal for AI research, that uses a
2180
            small deck to focus on core poker concepts like betting strategy and imperfect
2181
            information.
2182
       Here is a detailed breakdown of the rules to clarify legal moves. Note that in this
2183
            implementation, the "Check" action is not available; players must use "Call"
2184
            instead. A call may be zero-cost if there is no outstanding bet to match.
2185
2186
        **1. Setup & Preliminaries**
2187
        * **Players:** 2.
2188
          **Deck:** 6 cards (two Jacks, two Queens, two Kings).
          **Blinds:** Before cards are dealt, mandatory bets are posted:
2189
               Player 1 (P1) posts a **Small Blind** of 1 unit.
2190
               Player 2 (P2) posts a **Big Blind** of 2 units.
2191
        * **The Deal:** Each player receives one private card, face down.
2192
        **2. Core Betting Rules**
2193
        * **Raise Sizing:** The amount to raise is fixed.
2194
           * **Round 1:** The raise amount is **2 units**.
2195
              **Round 2:** The raise amount is **4 units**.
2196
          **Total Betting Cap:** The total betting cap for each round is a maximum of **two
2197
            raises**.
2198
           **Acting First: ** Player 1 (the small blind) acts first in both betting rounds (pre-
            flop and post-flop).
2199
2200
        **3. Round 1: Pre-Flop Betting**
2201
        This round occurs before the public card is revealed.
2202
          **P1's First Action:** P1 must act on P2's 2-unit Big Blind.
2203
              **Fold: ** Forfeit the 1-unit blind. P2 wins the pot.
2204
               **Call:** Match the 2 units by putting in 1 more unit.
2205
               **Raise:** Make a 2-unit raise, for a total of 4 units (P1 puts in 3 units). The
2206
                total betting cap has been reached.
2207
           **P2's Action:**
            * If P1 **called**, P2 can **Call** (a zero-cost action, as bets are equal) to end
2208
                the round, or **Raise** (by putting in 2 more units to make it 4 total).
2209
               If P1 **raised**, P2 can only **Call** (by putting in 2 more units) or **Fold**.
2210
                The betting cap has been reached.
2211
        * **P1's Second Action (if necessary):** If P1 called and P2 then raised, the action
2212
            returns to P1. P1 can only **Call** (by putting in 2 more units) or **Fold**.
2213
```

4. The Flop: Public Card

```
2214
       After Round 1 betting concludes, one public card is dealt face-up. This card is shared
2215
            by both players.
2216
2217
       **5. Round 2: Post-Flop Betting**
       This round occurs after the flop. There are no blinds.
2218
2219
           **P1's First Action:**
2220
           * **Call:** Make a zero-cost call to pass the turn (as there is no outstanding bet)
2221
2222
           * **Raise:** Make a 4-unit raise.
           **P2's Action:**
2223
           * If P1 **called** (at zero-cost), P2 can also **Call** (at zero-cost, ending the
2224
               round) or **Raise** 4 units.
2225
               If P1 **raised**, P2 can **Call** (matching the 4 units), **Raise** (by putting
2226
                in another 4 units, for a total bet of 8), or **Fold**. The total betting cap
                has been reached.
           **Subsequent Actions:**
2228
              If P2 **raised** (after P1's initial zero-cost call), the action returns to P1,
2229
                who can **Call** (the 4 unit bet), **Raise** (to 8 total), or **Fold**. The
2230
                total betting cap has been reached.
2231
           * If a player **raises**, the other player can only **Call** or **Fold**, as the
2232
               betting cap has been reached.
2233
        **6. Showdown & Hand Ranking**
2234
       If neither player folds, a showdown occurs after Round 2 betting.
2235
2236
           **Hand:** A player's hand is their private card combined with the public card.
2237
           **Hand Ranks (best to worst):**

    **Pair:** Two cards of the same rank (e.g., J-J). Higher pairs beat lower pairs.

2238
           2. **High Card:** If no one has a pair, the player with the highest card wins (K > Q
2239
                > J).
2240
           **Ties:** If both players have the same hand rank (e.g., both have a King-high), the
2241
             pot is split.
2242
       **7. Winning**
2243
       A player wins the pot either by being the only one left after the other folds, or by
2244
            having the best hand at showdown.
2245
2246
       H.8 GIN RUMMY
2247
2248
        # The Game of Gin Rummy
2249
2250
       Gin Rummy is a two-player card game played with a standard 52-card deck. The
2251
        primary objective is to form "melds" in your hand, which are either sets of
2252
        three or four cards of the same rank (e.g., 7h 7c 7d) or runs of three or more
       cards of the same suit in sequence (e.g., 4h 5h 6h). Cards not part of any meld
2253
       are referred to as "deadwood." The value of deadwood cards corresponds to their
2254
       rank (Aces are 1 point, face cards are 10, and number cards are their face
2255
        value). The ultimate goal is to minimize the point value of your deadwood.
2256
2257
       A round of Gin Rummy concludes when a player "knocks." A player can choose to
       knock on their turn if the total point value of their deadwood is less than or
2258
       equal to a predetermined "knock card" value. Announcing "gin" is a special type
2259
       of knock where a player has no deadwood at all.
2260
2261
        # Player Hand Information: This section provides details about your own hand.
2262
2263
       Deadwood: This calculates the current point total of the cards in
2264
       your hand that are not part of a valid meld (a set or a run). Minimizing this
2265
       value is the primary goal.
2266
2267
        The Card Grid: This is a visual representation of the cards you currently hold.
```

It is organized logically for easy parsing:

```
2268
2269
        Rows: Each of the four rows corresponds to a suit, in the order of Spades (top
2270
        row), Clubs, Diamonds, and Hearts (bottom row).
2271
        Columns: The columns represent the rank of the cards, ordered from Ace on the
2272
        far left to King on the far right.
2273
2274
       Here are also some example moves:
2275
2276
       Player: 0 Action: Pass
        Player: 1 Action: Draw upcard
2277
        Player: 1 Action: Jc
2278
        Player: 0 Action: 3d
2279
        Player: 1 Action: Draw stock
2280
        # Action Legality is Dictated by Game Phase: Before selecting a move, you must
2282
        first check the phase.
2283
2284
        If the phase is Draw, the only valid actions are Draw upcard or Draw stock.
2285
2286
        If the phase is Discard, the only valid actions are to discard a specific card
        from your hand (e.g., Action: 4c) or to Knock.
2287
2288
        A player cannot discard a card until after they have successfully drawn one.
2289
2290
        # Special Case: The First Turn of the Round
2291
2292
        The very first turn of a round has a unique rule. The non-dealer has the first
2293
        option on the initial upcard.
2294
2295
        The non-dealer can either take the upcard (Draw upcard) or Pass.
2296
        If the non-dealer passes, the dealer then has the same choice: take the upcard
2297
        or pass.
2298
2299
        If both players pass on the initial upcard, the non-dealer must then start their
2300
        turn by drawing from the stock pile. After this initial sequence, play continues
        with the standard draw/discard phases.
2301
2302
2303
        # Knocking:
2304
        When a player knocks in Gin Rummy, the round immediately ends and a specific
        sequence of scoring, known as the "layoff," begins. Here is a detailed
2305
        breakdown of what happens.
2306
2307
        1. The Knock and Laying Down Hands
2308
        First, the player who is knocking (the "knocker") lays their hand face up on the
2309
        table, organising their cards into melds (sets and runs) and separating their
2310
        unmelded cards, known as "deadwood."
2311
        What the Player Needs to Do After Knocking
2312
        After sending Action: Knock, the player must follow a strict, multi-step process
2313
        to lay down their hand for scoring.
2314
2315
        Step 1: Declare Your Melds
        The player must now explicitly declare their melds to the game, one by one.
2316
        For exampple, if the agent's hand contains two valid runs:
2317
2318
       A run of clubs: 7c8c9cTc
2319
2320
       A run of diamonds: 9dTdJdQd
2321
        Correct First Move in the Knock Phase:
```

```
2322
       Player: 1 Action: 7c8c9cTc
2323
2324
       Player: 1 Action: 9dTdJdQd
2325
        Step 2: Declare Subsequent Melds
2326
        After the agent declares its first meld, it will receive a new observation. The
2327
        game will still be in Phase: Knock. The Valid actions will now include any
2328
        remaining melds that can be made from the cards left in the hand.
2329
2330
        Note: The value of the knocker's deadwood must be 10 points or less (or the value of
2331
        the designated knock card for that round). Face cards are worth 10 points, aces
2332
        are 1 point, and all other cards are their numerical value.
2333
2334
        2. The Opponent's Turn: Laying Off
        Next, the defending opponent lays down their own hand, also separating their
        melds from their deadwood. Crucially, the opponent then gets the opportunity to
2336
        "lay off" any of their own deadwood cards by adding them to the knocker's melds.
2337
2338
        For example:
2339
        If the knocker has a meld of three Kings (Ks Ks Ks), and the opponent has the
2340
        fourth King (Ks) as deadwood, they can add it to the knocker's set, thus
2341
        eliminating those 10 points from their deadwood count.
2342
2343
        If the knocker has a run of 5h 6h 7h, the opponent can lay off a 4h or an 8h
2344
        to extend the run.
2345
        The knocker is not allowed to lay off any of their deadwood on the opponent's melds.
2346
2347
        3. Scoring the Hand
2348
        After the opponent has finished laying off their cards, both players calculate
2349
        the final value of their remaining deadwood. The scoring for the hand is then
        determined in one of three ways:
2350
2351
        a) A Successful Knock
2352
        If the knocker's deadwood count is lower than the opponent's deadwood count,
2353
        the knocker scores the difference between the two counts.
2354
        Example: The knocker has 7 points of deadwood. The opponent initially has 35
2355
        points, but after laying off a 10-point card, their deadwood is reduced to 25.
2356
        The knocker scores 18 points (25 - 7).
2357
2358
        b) An Undercut
2359
        If the opponent, after laying off their cards, has a deadwood count that is
        equal to or less than the knocker's count, they have "undercut" the knocker. In
2360
        this scenario, the opponent scores the difference in points (if any) plus a
2361
        bonus, which is typically 25 points.
2362
2363
        Example: The knocker has 8 points. The opponent has 6 points after layoffs. The
2364
        opponent scores 2 points (8 - 6) plus a 25-point bonus, for a total of 27
2365
        points.
2366
        c) Going Gin
2367
        If the knocker has a deadwood count of zero, this is called "going gin." The
2368
        knocker receives a bonus (typically 25 points) in addition to the full value of
        the opponent's entire deadwood count. When a player goes gin, the opponent is
2369
        not allowed to lay off any of their cards.
2370
2371
        Example: A player goes gin. Their opponent has 42 points of deadwood. The
2372
        ginning player scores 42 points plus a 25-point gin bonus, for a total of 67
2373
        points.
2374
        What if the Stock Pile Runs Out?
2375
```

If the stock pile is reduced to its last two cards and the player who drew the

third to last card discards without knocking, the hand is declared a draw. No points are awarded to either player, and the deal passes to the next player for a new round.

If you see Phase: Wall in an observation, it means:

The Stock Pile is Exhausted: The round has concluded because there are no more cards to be drawn from the stock.

No Player Has Knocked: Neither you nor your opponent were able to knock by the time the last card was drawn.

The Hand is a Draw: No points are awarded to either player for this round. The hand is over.

No Action is Required: The game is in a terminal state for the current round. The only thing to do is to acknowledge the result and wait for the next hand to be dealt. The deal will typically pass to the player who didn't deal the drawn hand. Thus `Player: X Action: Pass` must be provided as action.

H.9 Quadranto

Quadranto is a partially observable game in which two players try to catch each other in a $4\ \text{by }4\ \text{matrix}.$

The 4 by 4 matrix is divided in 4 quadrants. At the beginning, player 0 is randomly placed in the top left quadrant and player 1 is randomly placed in the bottom right quadrant.

During their turn, each player can choose to move in each of the four cardinal directions, "Left", "Right", "Up", "Down". Or they can choose to "Stay", which means they remain where they are. When a player moves, if it lands on the same location where the other player is, it wins and the game ends.

The observation tells the player where it is located and in which *quadrant* the opponent player is located. Therefore, neither player knows exactly where the other player is located until the very moment in which one player catches the other.

If the players perform a total of 20 moves without catching each other, the game ends in a draw, both players get 0 points. If one catches the other, the winning player gets +1 points and the losing player gets -1 points.

H.10 HAND OF WAR

Hand of War is a strategic card game where choosing your cards wisely is key to victory. You'll manage a hand of cards, adding a layer of tactical decision-making to every round as you aim to capture all of your opponent's cards.

Objective:

* The goal of Hand of War is to capture as many of your opponent's cards.

Setup:

- * **Shuffle and Deal:** Thoroughly shuffle the deck. Deal the entire deck evenly between two players, face down.
- * **Form Hands:** Each player draws the top three cards from their draw pile.
- 2429 **Gameplay (The "Battle"):**

```
2430
        * **Choose a Card:** Simultaneously, both players select one card from their
2431
       hand and place it face down.
2432
       * **Reveal and Compare:** Both players flip their chosen cards.
        * **Higher Card Wins:** The player with the higher-ranking card wins the
2433
       battle and takes both cards, placing them at the bottom of their win pile.
2434
        * **Card Ranking:** Ace (High), K, Q, J.
2435
        * **Draw New Cards:** After the battle, players draw from their draw pile to
2436
       replenish their hand to three cards.
2437
2438
       **"Showdown" (When Cards Tie):**
2439
           **Declaration:** If cards are of the same rank, a "Showdown" occurs.
2440
        * **Face-Down Cards:** Each player places 1 card from their draw pile face
2441
       down.
2442
       * **Choose Battle Card:** Players choose one card from their hand and place it
       face up.
2443
          **Determine Showdown Winner:** Higher battle card wins all cards in the
2444
        Showdown.
2445
           **Another Tie: ** Repeat Showdown process (burn 1, choose card).
2446
          **Draw After Showdown:** Players replenish their hand to three cards.
2447
       **Game End Conditions:**
2448
2449
       The game can end in one of two ways:
2450
2451
       1. **Winning by Capturing All Cards:** The game ends immediately if one player
2452
       possesses all 16 cards. That player is the winner.
2453
        2. **Winning by Win Pile Count (Draw Pile Depleted):** If a player's draw pile
2454
       becomes empty and they are required to perform an action they cannot complete
2455
        (such as drawing a card to replenish their hand after a battle, or burning a
2456
        card during a Showdown), the game ends immediately. In this scenario:
2457
           * The number of cards in each player's win pile is counted.
              The player with more cards in their win pile wins the game.
2458
           * If both players have an equal number of cards in their win piles, the
2459
           game results in a draw.
2460
              Note: The win pile is never reshuffled into the draw pile during play.
2461
2462
       **Notation**
       Player 0 always moves first, followed by player 1. Thus if the previous move was
2463
       a player 0 move, then the next move is a player 1 move.
2464
2465
       For convenience, a list of played cards is provided in PubliclyRevealedCards,
2466
       but not playable anymore.
2467
2468
2469
2470
2471
2472
2473
2474
2475
2476
2477
2478
2479
2480
2481
2482
```

I Sample synthesized CWMs and inference functions

We show the resulting CWM after tree search refinement on perfect and imperfect information games, and for the latter we include both open and closed deck inference. This includes inference and value functions when they exist.

I.1 TIC-TAC-TOE (PERFECT INFORMATION)

2484

2485 2486

2487

2488

2489 2490

```
import collections
2492
             from typing import Any, cast
2493
             # Type aliases for clarity as defined in the problem
2494
             Action = str
State = dict[str, Any]
2495
             PlayerObservation = dict[str, Any]
2496
             # Game constants for mapping players, marks, and defining winning conditions.
2497
             _MARK_TO_PLAYER = {'x': 0, 'o': 1}
             WINNING LINES = [
2498
                 # Rows
2499
                 (0, 1, 2), (3, 4, 5), (6, 7, 8),
                 # Columns
2500
                 (0, 3, 6), (1, 4, 7), (2, 5, 8),
                 # Diagonals
2501
                 (0, 4, 8), (2, 4, 6)
2502
            # --- Helper Functions ---
2503
2504
            \label{liststr} \mbox{def \_check\_winner(board: list[str \mid None]) -> str \mid None:}
                   Checks if there is a winner on the board, returning the winner's mark."""
2505
               for line in _WINNING_LINES:
                 p1, p2, p3 = line
# Check if all three cells in a line are the same and not empty.
2506
2507
                  \textbf{if} \ board[p1] \ \textbf{and} \ board[p1] \ == \ board[p2] \ == \ board[p3] : \\
                   return board[p1]
2508
             \begin{array}{lll} \textbf{def \_is\_game\_over(board: list[str \mid None])} \ -> \ bool: \\ """ \textbf{Checks if the game has ended either by a win or a draw."""} \end{array} 
2510
              return _check_winner(board) is not None or all(cell is not None for cell in board)
2511
2512
             # --- Core Game Functions ---
2513
             def apply_action(state: State, action: Action) -> State:
2514
                  "Returns the new state after an action has been taken."""
               mark = action[0]
2515
               row = int(action[2])
               col = int(action[4])
2516
               # Create a copy of the board to modify.
new_board = state['board'][:]
index = row * 3 + col
2517
2518
               new_board[index] = mark
2519
               # A game is over if there is a winner or the board is full.
2520
               if is game over(new board):
2521
                 next_player_mark = None
2522
                 # Alternate turns between 'x' and 'o'
                 next_player_mark = 'o' if mark == 'x' else 'x'
2523
2524
               return {'board': new board, 'current player mark': next player mark}
2525
             def get_current_player(state: State) -> int:
               """Returns current player, with -1 for chance and -4 for terminal."""
mark = state['current_player_mark']
2526
               if mark is None:
2527
                 return -4 # Terminal state
2528
               return _MARK_TO_PLAYER[mark]
2529
             def get_player_name(player_id: int) -> str:
                 ""Returns the name of the player, with 'chance' for -1, and 'terminal' for -4."""
2530
               if player_id == -4:
2531
                 return 'terminal
               if player_id == -1:
    return 'chance'
2532
               return str(player_id)
             def get_rewards(state: State) -> list[float]:
2534
                  'Returns the rewards per player from their last action."""
2535
               # Rewards are only non-zero in a terminal state.
               if get_current_player(state) != -4:
2536
2537
               winner_mark = _check_winner(state['board'])
               if winner_mark == ':
```

```
2538
                return [1.0, -1.0] # Player 0 wins
if winner_mark == 'o':
    return [-1.0, 1.0] # Player 1 wins
2539
2540
                return [0.0, 0.0] # Draw
2541
2542
              def get_legal_actions(state: State) -> list[Action]:
                 """Returns legal actions that can be taken in current state."""
# No legal actions if the game is over.
2543
                if get_current_player(state) == -4:
2544
                   return []
2545
                actions = []
2546
                mark = state['current_player_mark']
# An action is legal for any empty (None) square.
for i, cell in enumerate(state['board']):
    if cell is None:
2547
2548
                     row, col = divmod(i, 3)
2549
                actions.append(f"{mark}({row},{col})")
return actions
2550
              def get_observations(state: State) -> list[PlayerObservation]:
    """Returns the observation for player."""
# Fully_observable game, so both players see the complete state.
2551
2552
                return [state, state]
2553
2554
2555
2556
2557
2558
2559
2560
2561
2562
2563
2564
2565
2566
2567
2568
2569
2570
2571
2572
2573
2574
2575
2576
2577
2578
2579
2580
2581
2582
2583
2584
2585
2586
2587
2588
2589
2590
2591
```

```
2592
          I.2 Bargaining (imperfect information, open deck)
2593
2594
2595
2596
2597
2598
2599
2600
2601
2602
2603
2604
2605
2606
2607
2608
2609
2610
2611
2612
2613
2614
2615
2616
2617
2618
2619
2620
2621
2622
2623
2624
2625
2626
2627
2628
2629
2630
2631
2632
2633
2634
2635
2636
2637
          import copy
          import itertools
2638
          import random
from typing import Any, Dict, List, Optional, Tuple
2639
2640
          # Type Aliases for clarity
          State = Dict[str, Any]
2641
          Action = str
PlayerObservation = Dict[str, Any]
2642
          # --- Game Constants ---
ITEMS = sorted(['X', 'Y', 'Z'])
2643
2644
          NUM_PLAYERS = 2
MAX_TURNS = 10
2645
          MAX_ITEM_VALUE = 10
          \# This is a fixed set of chance outcomes, likely for deterministic testing or analysis.
```

```
2646
                       _CHANCE_ACTIONS = sorted(['1,2,3;8,1,0;4,0,2', '1,4,1;4,1,2;2,2,0', '2,2,1;1,1,6;0,4,2', '1,4,1;9,0,1;2,2,0', '1,4,1;5,1,1;0,1,6', '4,1,1;2,1,1;1,0,6', '3,1,1;1,4,3;0,2,8', '1,1,3;0,1,3;1,3;2', '1,3,1;2,2,2;10,0,0', '1,2,2;2,3,1;4,0,3', '1,4,1;6,1,0;8,0,2', '1,1,3;7,3,0;0,4,2', '1,5,1;4,0,6;3,1,2', '3,3,1;3,0,1;0,2,4', '1,2,3;8,1,0;7,0,1', '4,1,2;0,6,2;2,2,0', '2,1,2;3,2,1;4,2,0', '1,3,1;4,2,0;8,0,2', '2,1,3;3,1,1;0,10,0', '1,3,1;6,1,1;4,1,3', '2,2,1;3,0,4;2,1,4', '3,3,1;1,1,4;3,0,1', '1,2,3;0,5,0;3,2,1',
2647
2648
                                                                  '1.3,1;1,2,3;3,1,4',
2649
                         1,3,1;0,1,7;6,0,4',
2650
                        '1,3,1;2,2,2;1,2,3',
                                                                                                                                                                                                          '2,2,1;1,3,2;5,0,0', '1,3,1;4,2,0;1,1,6'
                       '3,2,1;2,1,2;1,3,1',
2651
                                                                                                                                                                                                                                                        '3,1,1;1,0,7;0,8,2'
                         1,1,3;6,1,1;0,1,3',
                                                                  '4,1,1;1,4,2;2,1,1',
                                                                                                                                                                                                           '2,1,2;3,4,0,1,2,1
'1,1,5;3,2,1;8,2,0', '3,3,1;2,1,1,2,1,
'1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2;10,0,0;2,1,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' '1,4,2' 
2652
                        '1,2,2;4,0,3;2,1,3',
                       '2,1,4;1,8,0;3,0,1',
                                                                                                                                                                                                          '2,1,2;2,2,2;1,8,0',
2653
                                                                                                                                                                                                         '2,1,2,1,8,9,1,6,1', '1,3,1;3,1,4;10,0,0'
'3,1,1;0,5,5;3,0,1', '3,1,1;1,0,7;2,4,0'
'3,1,1;0,3,7;1,1,6', '2,1,4;0,2,2;2,2,1'
                       '1,4,1;5,1,1;2,0,8',
2654
                        '1,3,1;1,3,0;7,0,3',
                                                                                                                                                             '3,1,1;2,1,3;0,7,3',
'2,3,1;1,2,2;2,1,3',
                       '2,2,1;2,1,4;2,3,0',
2655
                                                                                                               '2,3,1;0,3,1;2,2,0', '1,1,4;0,6,1;1,5,1',
'2,1,2;1,4,2;3,2,1', '2,1,4;3,0,1;2,6,0',
'2,1,2;3,2,1;3,4,0', '1,1,3;3,4,1;1,9,0',
                                                                                                                                                                                                          '1,1,5;10,0,0;3,2,1',
'1,1,5;1,4,1;4,1,1',
                       '1,3,1;2,0,8;0,3,1',
'4,1,1;0,0,10;1,2,4',
                                                                   '4,2,1;1,0,6;0,2,6',
'1,1,3;1,9,0;7,0,1',
                                                                                                                                                                                                                                                        '3,1,1;1,5,2;1,5,2
                                                                                                                                                                                                                                                        '2,2,1;1,3,2;3,2,0'
2656
                                                                    '3,1,1;2,2,2;0,8,2',
                                                                                                                                                                                                           '2,4,1;2,1,2;2,0,6',
                        '2,2,1;3,0,4;0,2,6',
                                                                                                                                                             '1,4,2;4,1,1;4,1,1',
'2,2,1;1,1,6;3,1,2',
2657
                                                                                                                                                                                                           '1,2,2;6,1,1;0,1,4',
'1,3,1;4,2,0;3,1,4',
                       '3,1,1;0,1,9;2,4,0',
'3,1,1;3,1,0;0,3,7',
                                                                   '1,1,4;1,1,2;5,5,0',
'2,1,4;5,0,0;1,4,1',
                                                                                                                 '3,1,1;3,1,0;2,0,4',
'4,1,1;1,5,1;0,4,6',
                                                                                                                                                                                                                                                        12,3,1;0,2,4;4,0,2
                                                                                                                                                                                                                                                        '3,1,1;0,2,8;1,1,6'
2658
                                                                                                                 '5,1,1;1,1,4;1,1,4',
'1,3,1;0,2,4;2,1,5',
                       '3,1,1;1,3,4;2,4,0',
'1,5,1;8,0,2;2,1,3',
'3,3,1;0,1,7;2,0,4',
                                                                   '4,1,1;1,3,3;0,6,4',
'1,3,1;4,2,0;5,0,5',
                                                                                                                                                             '1,1,3;3,1,2;2,2,2',
'1,3,1;4,1,3;3,2,1',
'4,1,1;1,6,0;1,4,2',
                                                                                                                                                                                                          '1,3,2,8,0,1;1,3,0', '1,1,5;0,5,1,8,2,0'

'2,3,1;1,1,5;1,2,2', '2,2,1;2,0,6;0,1,8'

'2,2,2;1,2,2;2,1,2', '5,1,1;1,5,0;1,1,4'
2659
                                                                   '1,3,3;4,0,2;1,1,2',
                                                                                                                 '1,4,1;1,2,1;2,2,0',
                                                                                                                                                                                                                                                        '5,1,1;1,5,0;1,1,4',
                                                                                                                '2,1,3;0,1,3;3,1,1',
'1,4,2;6,1,0;0,2,1',
2660
                       '3,3,1;2,1,1;0,1,7',
'4,2,1;0,2,6;2,1,0',
                                                                                                                '1,2,3;0,2,2;2,1,2', '2,2,1;3,1,2;3,2,0',
1,5,1;3,1,2;4,1,1', '1,2,2;0,4,1;4,1,2',
'3,1,2;1,1,3;2,2,1', '2,2,2;0,2,3;1,4,0',
'2,2,2;0,2,3;3,0,2', '2,1,2;5,0,0;2,4,1',
'1,1,3;7,0,1;1,6,1', '3,2,1;1,2,3;2,2,0',
'3,1,1;0,5,5;1,2,5', '1,2,3;10,0,0;5,1,1',
'2,2,2;2,0,3;0,3,2', '2,4,1;3,0,4;3,1,0',
'3,3,1;2,1,1;3,0,1', '1,1,3;6,1,1;1,3,2',
'2,4,1,4,0,2,1,1,4',
'1,3,1,1,1,3,1,4',1,1,3,1,1,4',1,3,1,0',
2661
                                                                    '2,1,2;1,8,0;2,4,1',
                                                                                                                                                                                                          '3,1,1;1,1,6;1,0,7, '1,1,1,3;5,5,0;1,0,0,1,1,1,4;5,1,1;2,4,1', '1,1,3;5,5,0;1,0,1,0,1,0,1,0,1,1,3;5,0,1,0,1,0,1,0,1,1,3;5,0,0,1,0,4,1,3'
                         1,3,1;4,0,6;0,3,1',
                       '3,1,1;2,0,4;1,7,0',
                                                                    '2,1,2;5,0,0;1,2,3',
2662
                                                                                                                                                                                                          '1,1,3;9,1,0;6,1,1',
                        '3,3,1;2,0,4;0,3,1',
                                                                    '1,1,3;6,1,1;0,4,2',
                                                                    '2,2,2;5,0,0;1,1,3',
                                                                                                                                                                                                           '3,1,1;0,4,6;2,1,3',
                         1,1,3;1,3,2;4,6,0',
                                                                                                                                                                                                                                                        1,3,1;3,0,7;2,1,
2663
                       '2,1,2;0,2,4;4,2,0',
                                                                    '1,1,5;5,0,1;5,5,0',
                                                                                                                                                                                                            '1,4,1;0,1,6;9,0,1',
                                                                                                                                                                                                                                                         1,1,5;2,3,1;7,3,0
                                                                                                                                                                                                           '5,1,1;0,2,8;1,3,2',
'2,1,3;4,2,0;2,0,2',
2664
                         1,5,1;2,1,3;0,1,5',
                                                                     '1,3,1;2,1,5;0,3,1',
                                                                                                                                                                                                                                                        '3,2,1;3,0,1;0,1,8',
                                                                    '1,3,1;1,3,0;3,1,4',
                                                                                                                '2,1,2;0,4,3;2,0,3',
                                                                    '2,1,2;0,8,1;4,2,0',
                                                                   '2,2,1;1,4,0;2,0,6',
'3,1,3;1,4,1;1,1,2',
2666
                        '4,2,1;0,1,8;1,2,2',
                       '1,1,3;3,4,1;1,3,2',
                                                                    '1,2,2;4,0,3;2,3,1',
                                                                                                                                                                                                          '2,1,2;2,4,1;3,0,2', '2,1,3;2,6,0;0,1,3'
'2,2,2;0,4,1;2,0,3', '2,2,2;0,1,4;2,3,0'
'1,2,2;6,1,1;0,4,1', '1,5,1;5,0,5;3,1,2'
2667
                       '2,2,2;1,0,4;3,1,1',
'3,1,1;0,5,5;1,6,1',
                                                                                                                 '1,2,3;7,0,1;3,2,1', '1,4,1;3,0,7;0,1,6', '4,2,1;0,1,8;2,0,2', '2,2,1;0,3,4;4,0,2',
                                                                   '1,5,1;5,1,0;2,0,8',
2668
                       '3,1,1;1,0,7;1,5,2',
                                                                   '2,1,2;4,2,0;1,0,4',
                                                                                                                 '4,1,2;1,2,2;1,6,0',
                                                                                                                                                             '2,3,2;4,0,1;1,2,1',
                                                                                                                                                                                                           '1,2,2;0,1,4;8,0,1', '2,1,4;3,4,0;2,2,
'3,1,3;2,4,0;0,1,3', '3,2,1;1,2,3;2,1,
                                                                                                                                                              '2,1,2;4,2,0;0,2,4',
                                                                                                                                                                                                                                                        '2.1.4:3.4.0:2.2.1
                       '2,1,2;0,8,1;3,0,2',
'4,1,2;1,6,0;2,0,1',
                                                                   '4,1,1;1,2,4;1,0,6',
'2,1,3;3,4,0;1,5,1',
                                                                                                                 '5,1,1;0,7,3;1,2,3',
'4,1,2;0,6,2;1,6,0',
2669
                                                                                                                                                             '1,2,2;2,2,2;2,2,2',
                                                                                                                 '4,1,1;1,0,6;0,1,9',
'1,2,4;0,5,0;4,1,1',
                                                                                                                                                             '2,2,1;3,1,2;1,1,6',
'2,1,2;1,0,4;1,6,1',
                                                                                                                                                                                                           '2,2,1;2,2,2;2,1,4',
'1,1,4;1,5,1;4,6,0',
                        '1,4,1;9,0,1;0,2,2',
                                                                    '2,2,1;0,3,4;1,0,8',
2670
                                                                                                                                                                                                                                                        1,1,4;1,5,1;0,6,1
                                                                    '1,3,1;4,2,0;5,1,2',
                        '4,1,1;2,2,0;1,0,6',
2671
                                                                   '1,5,1;2,1,3;5,0,5',
                                                                                                                 '1,4,1;1,1,5;5,1,1',
                                                                                                                                                              '1,3,1;0,1,7;5,1,2',
                                                                                                                                                                                                          '1,2,2;8,0,1;4,1,2',
                        '3,1,1;1,3,4;1,5,2',
                                                                                                                                                                                                         '2,1,3;1,5,1;0,10,0',
'2,1,4;3,4,0;1,4,1',
                                                                   '1,4,1;6,1,0;2,1,4',
'3,1,1;0,3,7;1,3,4',
                                                                                                                '1,2,4;4,1,1;0,1,2', '3,2,1;1,0,7;2,2,0', '3,1,2;0,10,0;1,3,2', '1,2,4;0,1,2;2,0,2',
                        '3,3,1;0,2,4;1,2,1',
                                                                                                                                                                                                                                                         1,2,2;0,1,4;6,1
2672
                                                                                                                                                                                                                                                        12,2,2;1,3,1;0,2,3
                         1,4,1;8,0,2;2,2,0
                                                                      '3,1,3;1,7,0;1,4,1',
                                                                                                                  '2,4,1;1,0,8;0,2,2',
                                                                                                                                                              '1,1,4;4,2,1;1,1,2',
                                                                                                                                                                                                           '2,1,2;3,2,1;5,0,0',
2673
                                                                   '2,3,2;2,2,0;0,2,2',
'1,1,3;0,10,0;2,2,2',
                                                                                                                 '4,1,1;2,0,2;1,4,2',
'3,1,1;3,1,0;0,1,9',
                                                                                                                                                               '1,4,1;7,0,3;4,1,2',
                                                                                                                                                                                                           '3,1,1;1,7,0;0,4,6',
                        '1,3,1;9,0,1;0,1,7',
                                                                                                                                                                                                                                                        '3.2.2:2.1.1:2.0.2
                                                                                                                                                              '1,1,3;3,7,0;3,4,1',
2674
                        '2,2,1;1,3,2;3,0,4',
                                                                                                                                                                                                           '2,2,2;1,0,4;1,1,3',
                                                                                                                                                                                                                                                        1,3,1;7,1,0;9,0,1
                                                                                                                 '1,3,1;3,2,1;0,1,7',
'1,2,4;2,0,2;2,4,0',
                                                                    '3,1,2;2,0,2;2,2,1',
                                                                                                                                                              '1,1,3;2,8,0;4,0,2',
                                                                                                                                                                                                           '2,3,1;0,1,7;2,0,6',
2675
                                                                                                                                                              '3,1,1;1,0,7;1,4,3',
                                                                                                                                                                                                           '1,4,1;1,2,1;1,1,5',
                                                                                                                                                                                                                                                         11,1,3;9,1,0;3,4,1
                        '1,4,1;0,1,6;6,0,4',
                                                                     '1,1,4;0,2,2;2,8,0',
                        '2,2,1;1,4,0;2,2,2',
                                                                    '3,1,1;0,1,9;1,5,2',
                                                                                                                 '3,1,1;0,1,9;2,2,2',
                                                                                                                                                              '1,3,3;4,2,0;1,1,2',
                                                                                                                                                                                                           '1,1,3;1,0,3;5,5,0',
                                                                                                                                                                                                                                                        '4,2,1;1,2,2;0,1,8'
2676
                                                                                                                                                              '2,2,2;1,2,2;2,3,0',
                       '1,4,1;4,1,2;0,1,6',
                                                                   '1,3,1;1,1,6;2,2,2',
                                                                                                                 '2,2,2;2,2,1;0,2,3',
                                                                                                                                                                                                           '1,1,4;4,2,1;9,1,0',
                                                                                                                                                             '3,1,3;2,1,1;0,1,3',
'1,3,1;0,3,1;2,2,2',
                                                                                                                  '4,1,1;1,3,3;0,7,3',
                                                                                                                                                                                                            '2,1,2;0,4,3;3,4,0',
                                                                                                                                                                                                                                                         '1,4,1;1,0,9;4,1,2
2677
                        '4,1,2;1,2,2;1,2,2',
                                                                     '1,4,2;2,1,2;2,0,4',
                                                                   '1,1,4;5,1,1;4,6,0',
                                                                                                                                                                                                           '3,1,2;1,1,3;0,2,4',
                         5,1,1;0,1,9;1,2,3',
                                                                                                                 '1,4,2;0,0,5;4,1,1',
                                                                                                                                                                                                                                                        '2,2,3;0,2,2;2,3,0',
                                                                                                                                                                                                          '1,5,1;0,0,10;3,1,2',
2678
                                                                   '3,1,2;3,1,0;0,8,1',
                                                                                                                 '5,1,1;1,2,3;0,1,9', '4,2,1;1,1,4;0,4,2',
                                                                                                                                                                                                                                                        1,2,2;2,0,4;6,1
                       '2,4,1;0,2,2;1,1,4',
                                                                                                                                                                                                            '1,4,1;1,1,5;1,1,5',
                                                                                                                 '4,2,1;0,4,2;1,3,0',
'3,1,1;1,3,4;0,8,2',
                                                                                                                                                                                                                                                         '1,4,1;0,1,6;8,0,2
                                                                     '1.2.2:6.0.2:8.1.0'.
                                                                                                                                                             '2,1,2;0,4,3;2,4,1',
'3,1,2;2,0,2;1,7,0',
                         1,1,4;3,3,1;8,2,0',
2679
                        '2,2,2;4,1,0;2,0,3',
                                                                   '2,4,1;1,2,0;3,0,4',
                                                                                                                                                                                                           '1,4,1;1,2,1;3,1,3',
                                                                                                                                                                                                                                                        1,1,3;4,3,1;2,8,0
                                                                                                                '3,1,1;1,3,4;0,8,2', '3,1,2;2,0,2;1,7,0', '1,4,1;1,2,1;3,1,3', '1,1,3;4,3,1;2,8,0'
'3,1,1;1,6,1;1,5,2', '2,1,4;3,0,1;1,8,0', '1,1,3;4,0,2;6,4,0', '2,2,1;0,3,4;1,3,2'
'3,1,2;0,10,0;2,2,1', '3,2,1;2,2,0;1,2,3', '1,3,1;1,2,3;4,2,0', '2,4,1;1,2,0;0,2,2'
'1,1,3;0,7,1;3,1,2', '2,1,2;2,4,1;2,6,0', '1,1,3;2,5,1;7,0,1', '1,3,1;0,0,10;2,2,2'
'1,1,3;3,4,1;3,7,0', '1,4,1;5,1,1;1,2,1', '1,4,1;6,1,0;1,2,1', '1,3,2;3,1,2;6,0,2'
'1,1,4;6,0,1;2,8,0', '2,2,1;1,3,2;2,2,2', '1,1,3;3,1,2;9,1,0', '2,1,4;2,2,1;3,0,1'
'1,1,3;1,9,0;4,3,1', '4,1,1;2,4;0,2,8', '1,1,3;6,1,1,0,10,0', '2,2,1;2,4;2,3,0'
'3,1,1;1,1,6;0,6,4', '1,3,1;1,3,0;1,0,9', '2,2,2;2,2,1;3,0,2', '3,1,2;0,0,5;1,5,1'
                        '4,1,2;0,8,1;2,2,0',
                                                                   '4,2,1;0,3,4;2,0,2',
'4,2,1;1,2,2;1,0,6',
2680
                        '4,1,1;1,4,2;0,3,7',
2681
                       '3,1,1;2,4,0;2,3,1',
                                                                   '2,1,2;2,4,1;0,0,5',
                       '2,2,1;2,1,4;5,0,0',
'1,5,1;3,0,7;2,1,3',
                                                                    '2,3,1;3,1,1;1,0,8',
'4,1,2;1,2,2;0,0,5',
2682
                       '2,4,1;2,0,6;3,1,0',
                                                                   '2,2,2;0,2,3;1,0,4',
2683
                        '4.1.2:1.6.0:1.4.1'.
                                                                    '1,2,3;5,1,1;1,3,1',
                                                                                                                '3,1,1;1,1,6;0,6,4', '1,3,1;1,3,0;1,0,9', '2,2,2;2,1;3,0,2', '3,1,2;0,0,5;1,5,1' '2,1,2;3,4,0;0,4,3', '3,2,2;0,5,0;2,1,1', '1,5,1;5,1,0;0,1,5', '1,2,2;8,0,1;6,1,1' '1,1,3;6,1,1;5,2,1', '1,1,4;2,8,0;0,6,1', '2,1,2;2,2,2;4,0,1', '3,1,3;0,10,0;1,4,1' '1,3,2;10,0,0;5,1,1', '2,1,2;3,4,0;0,8,1', '1,4,2;4,1,1;4,0,3', '3,1,2;1,3,2;2,4,0' '1,4,1;3,0,7;3,1,3', '2,2,2;3,1,1;2,1,2', '2,1,2;3,2,1;1,6,1', '1,3,2;3,1,2;2,4,1' '2,2,1;1,6;0,3,4', '1,1,3;1,0,3;1,3,2', '1,2,2,6;1,1;2,0,4', '1,3,2;3,1,2;2,2,1' '2,1,3;2,0,2;3,4,0', '2,1,4;1,0,2;0,2,2', '3,1,1;0,9,1;3,1,0', '1,5,1;3,0,7;1,1,4'
                        '1,3,3;4,0,2;4,2,0',
                                                                    '1,2,2;4,2,1;6,1,1',
2684
                       '2,1,2;1,2,3;2,6,0',
                                                                   '2,1,4;1,4,1;2,2,1',
'1,3,1;4,2,0;0,1,7',
2685
                        '1,2,4;2,2,1;10,0,0',
                       '2,2,2;1,4,0;0,4,1',
'1,5,1;6,0,4;3,1,2',
                                                                    '1,1,3;1,0,3;1,9,0',
2686
                                                                   '1,3,1;0,1,7;7,1,0',
                                                                   '1,4,1;2,2,0;5,1,1',
                        '2.2.1:1.2.4:2.0.6'.
                                                                   '1,4,2;6,1,0;6,0,2',
'1,1,5;0,5,1;2,3,1',
'1,1,3;0,1,3;2,5,1',
                                                                                                                1,3,2;4,2,0;2,0,4', '3,1,1;0,10,6;1,2,5', '1,3,2;3,1,2;7,1,0',
'3,1,1;1,2,5;0,1,9', '1,1,3;3,1,2;10,0,0', '1,1,3;6,4,0;0,4,2',
'1,4,1;8,0,2;2,1,4', '1,1,4;7,3,0;1,1,2', '1,3,1;2,2,2;7,1,0',
                                                                                                                                                                                                                                                        '1,1,4;0,2,2;3,7,0',
                        '1,4,1;1,2,1;9,0,1',
                                                                                                                                                                                                                                                        12,2,1;1,0,8;1,3,2
2688
                       '2,2,2;4,0,1;2,3,0',
                                                                                                                                                                                                           '1,3,1;2,2,2;7,1,0',
                                                                                                                                                                                                                                                        '3,1,1;1,0,7;3,1,0'
                         4.1.1:1.0.6:1.1.5
2689
                                                                   '1,3,1;1,1,6;6,1,1', '1,3,3;1,2,1;4,0,2', '3,1,1;0,10,0;1,3,4', '3,1,1;1,7,0;2,2,2', '1,4,1;0,2,2;1,0,9', '5,1,1;0,4,6;1,5,0', '1,1,5;8,2,0;1,4,1', '1,2,4;4,1,1;8,1,0', '3,1,1;0,0,10;1,1,6', '1,3,1;4,1,3;7,0,3', '1,2,4;2,0,2;8,1,0', '1,1,3;2,2,2;6,1,1',
                       '2,2,1;3,2,0;1,0,8',
                                                                                                                                                                                                                                                        1,4,1;1,1,5;3,0,7
                       '2,1,4;2,2,1;1,0,2',
2690
                                                                  '3,1,1,0,0,10;1,1,0,,
'3,3,1;0,0,10;1,2,1',
'1,4,2;0,2,1;4,0,3',
'1,4,1;1,2,1;6,1,0',
'1,2,2;8,1,0;0,3,2',
'1,2,2;8,1,0;0,3,2',
                                                                                                                                                                                                                                                        1,1,3;6,1,1;2,2,2
                         5,1,1;0,6,4;1,0,5',
                                                                                                                 '13,21,21,12,13,3', '13,14,8,0,27,1,0', '1,2,31,0,3;4,3,0', '1,2,22,0,3,2;8,1,0', '14,1;1,2,1;6,1,0', '12,4;4,1,1;6,2,0', '3,2,1;0,0,10;1,3,1', '3,1,1;1,4,3;0,0,10',1,2,2;8,1,0;0,3,2', '13,3,1;2,1,5;3,2,1', '1,1,4;5,5,0;3,3,1', '2,1,2;3,0,2;3,4,0',1,2,0,6;0,2,2', '1,1,3;2,8,0;3,1,2', '1,1,3;7,0,1;0,7,1', '2,3,1;2,1,3;3,1,1',1,4;1,9,0;4,2,1', '3,2,1;0,1,8;1,1,5', '4,1,1;0,4,6;1,3,3', '1,4,1;4,1,2;6,0,4',2,2;1,2,0,6;0,2,6', '2,2,2;0,4,1;1,2,2', '1,4,1;6,0,4;0,2,2', '1,2,2;4,2,1;6,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2,0',1,1;5,2
                        '1,2,2;6,0,2;2,3,1',
2691
                       '2,2,2;1,4,0;1,2,2',
'2,1,2;3,2,1;3,0,2',
2692
                       '1,3,1;7,1,0;6,0,4',
2693
                        '1,4,1;0,2,2;4,1,2',
                                                                   '1,1,5;9,1,0;1,4,1',
'3,1,2;1,5,1;3,1,0',
                                                                                                                '1,1,4;1,9,0;4,2,1',
'2,2,1;2,0,6;0,2,6',
                       '3,1,3;0,7,1;1,7,0',
2694
                                                                                                                '1,1,5;2,3,1;6,4,0',
'4,1,1;1,1,5;0,10,0',
                                                                                                                                                             '2,1,2;1,4,2;3,4,0', '1,1,4;4,2,1;2,8,0', '1,1,4;1,5,1;1,1,2', '1,1,5;7,3,0;1,4,1',
                                                                   '1,2,3;1,3,1;4,3,0',
'1,3,1;3,0,7;7,1,0',
                       '3,1,3;1,4,1;2,4,0',
                                                                                                                                                                                                                                                        '1,3,1;6,1,1;4,2,0'
                                                                                                                                                                                                                                                        '4,2,1;2,1,0;0,1,8'
                        11.2.2:4.0.3:0.3.2
                                                                    '1,2,2;6,1,1;2,2,2',
                                                                                                                 '2,2,2;0,4,1;2,3,0',
                                                                                                                                                              '1,4,1;3,1,3;5,0,5',
                                                                                                                                                                                                           '3,2,1;0,4,2;3,0,1',
                        '1,2,3;2,1,2;2,4,0',
                                                                                                                                                                                                                                                         '2,4,1;2,1,2;3,0,4',
2696
                                                                                                                '1,1,5;0,10,0;1,4,1', '1,1,3;3,7,0;6,1,1', '2,3,1;1,2,2;0,3,1', '3,1,1;0,7,3;1,0,7'
'2,2,2;3,1,1;2,2,1', '2,4,1;1,1,4;3,0,4', '2,1,3;4,2,0;1,5,1', '1,2,2;6,1,1;10,0,0'
                                                                   '2,3,1;4,0,2;1,2,2',
'1,4,1;0,1,6;5,0,5',
                       '2,3,1;2,1,3;3,0,4',
                         1,2,2;0,3,2;4,0,3',
2697
                                                                                                                                                            '2,1,3;1,8,0;1,2,2',
'2,2,2;0,1,4;3,1,1',
                                                                                                                '2,2,2;1,1,3;0,4,1',
'1,1,4;8,2,0;4,2,1',
                       '4,1,1;0,7,3;1,0,6'
2698
                        '1,1,5;1,4,1;10,0,0'.
                                                                                                                                                                                                                                                         12.2.1:3.1.2:2.2.2
                                                                                                                '2,1,3;2,0,2;1,8,0', '1,4,1;3,1,3;1,1,5', '1,1,3;6,1,1;5,5,0', '2,1,2;1,0,4;4,2,0',
                                                                   '1,1,3;3,1,2;1,3,2',
                                                                                                                                                                                                                                                        1,3,1;5,1,2;0,3,1
                       '3,1,1;1,4,3;1,5,2',
2699
                                                                   '4,2,1;0,2,6;1,0,6',
                                                                                                                '1,5,1;6,1,5;3,0,7', '2,2,1;6,2,6;1,4,0,'
'2,1,3;3,4,0;2,3,1', '3,1,2;6,2,4;1,5,1',
'1,1,3;3,1,2;6,1,1', '2,1,2;5,0,0;3,2,1',
'2,2,1;0,1,8;2,1,4', '1,4,1;1,2,1;0,1,6',
                        11.3.1:7.1.0:4.1.31.
                                                                     '4,2,1;1,2,2;1,1,4',
                                                                                                                                                                                                                                                          12.1.2:2.4.1:2.4.1
                        '2,3,1;0,2,4;2,1,3',
                                                                    '1,2,4;6,2,0;0,1,2',
                                                                                                                                                                                                                                                        '2,1,2;1,6,1;2,4,1
                        '2,1,3;1,5,1;2,3,1',
                                                                   '1,3,3;1,1,2;1,0,3',
'1,3,1;0,0,10;5,1,2',
                                                                                                                                                                                                           '1,1,3;1,9,0;4,0,2',
                                                                                                                                                                                                                                                         1,1,3;3,1,2;1,6,1
                        '1,2,2;8,1,0;4,0,3',
                                                                                                                                                                                                                                                         1.3,1;4,2,0;1,0,9
                                                                     '2,2,2;4,1,0;2,1,2',
                                                                                                                 '2,3,1;1,0,8;1,1,5'50'3,3,1;1,1,4;1,2,1',
'1,1,4;3,7,0;4,2,1',50'3,1,1;1,3,4;1,6,1',
                                                                                                                                                                                                           '3,1,2;1,7,0;1,1,3',
                                                                                                                                                                                                                                                        '1,3,1;6,1,1;6,0,4
                                                                                                                                                                                                           '3,1,1;0,1,9;1,0,7',
                                                                                                                                                                                                                                                         '2.2.2:3.0.2:1
                                                                                                                '5,1,1,0,5,5;1,0,5', '2,2,2;0,2,3;2,0,3', '2,1,2;4,2,0;1,2,3', '2,2,2;1,3,1;2,0,3', '3,1,2;2,4,0;0,2,4', '2,2,2;4,1,0', '4,2,1;1,1,4;0,2,6', '4,1,2;2,2,0;0,8,1', '3,1,1;0,2,8;2,1,3', '1,1,3;4,0,2;5,5,0', '2,1,4;4,2,0;2,2,1', '2,2,2;3,2,0;0,2,3',
                         2,4,1;0,1,6;1,2,0',
                                                                   '1,1,4;5,1,1;6,0,1',
                                                                                                                                                                                                                                                          11,4,1;4,1,2;5,1
                                                                   '3,1,1;0,4,6;1,1,6',
                         1,4,1;6,1,0;4,1,2',
```

'4.1.1:1.2.4:0.5.5

'1,1,3;0,1,3;7,0,

'3,2,2;2,1,1;0,1,4',

'1,3,1;7,0,3;1,2,3',

'1,1,3;4,0,2;5,5,0',

```
2700
2701
                       # --- Helper Functions ---
2702
                       def _parse_quantities(q_str: str) -> Dict[str, int]:
                              ""Parses a quantity string like '1,2,0' into a dictionary."""
return (item: int(q) for item, q in zip(ITEMS, q_str.split(',')))
2703
2704
                       def _format_quantities(quantities: Dict[str, int]) -> str:
    """Formats a quantity dictionary into a string like '1,2,0'."""
2705
                               return ",".join(str(quantities.get(item, 0)) for item in ITEMS)
2706
2707
                       def _create_agreement(state: State, offering_player: int, offered_quantities: Dict[str, int]) -> List[Dict[str, int]]:
                                 ""Creates the final agreement structure based on an accepted offer.'
2708
                               shares = [{}, {}]
shares[offering_player] = offered_quantities
2709
                              other_player = 1 - offering_player
# The other player gets the remainder of the item pool.
2710
                               shares[other_player] = {
2711
                                      item: state['pool'][item] - offered_quantities.get(item, 0) for item in ITEMS
2712
2713
                      # --- Core API Functions ---
2714
                       def apply_action(state: State, action: Action) -> State:
    """Returns the new state after an action has been taken."""
2715
                               new_state = copy.deepcopy(state)
2716
                              player_id = get_current_player(new_state)
2717
                              if player_id == -1: # Chance player sets up the game.
   pool_str, v0_str, v1_str = action.split(';')
   new_state['pool'] = _parse_quantities(pool_str)
2718
2719
                                      new_state['player_0_values'] = _parse_quantities(v0_str)
new_state['player_1_values'] = _parse_quantities(v1_str)
new_state['current_player'] = '0'
2720
2721
                                      return new state
2722
                               if "agrees" in action:
                                      # A player agrees to the last offer, ending the game.
last_offer = new_state['offer_history'][-1]
2723
                                      new_state['agreement'] = _create_agreement(new_state, last_offer['player'], last_offer['quantities'])
new_state['current_player'] = None # Mark as a terminal state.
2724
2725
                               elif "offers" in action:
2726
                                      # A player makes a new offer.
new_state['num_turns'] += 1
2727
                                       quantities = _parse_quantities(action.split(' offers ')[1])
2728
                                      new_offer = {
                                                'num_turn': new_state['num_turns'],
2729
                                                'player': player_id
                                               'quantities': quantities
2730
                                      new_state['offer_history'].append(new_offer)
2731
2732
                                      # If turn limit is reached, this offer becomes a forced, zero-reward agreement.
if new_state['num_turns'] >= MAX_TURNS:
2733
                                              new_state['current_player'] = None
                                              new_state['agreement'] = _create_agreement(new_state, player_id, quantities)
2734
                                       else:
2735
                                              new_state['current_player'] = str(1 - player_id) # Switch to other player.
2736
                              return new_state
2737
                      def get_current_player(state: State) -> int:
    """Returns current player, with -1 for chance and -4 for terminal."""
2738
                               player = state.get('current_player')
2739
                               if player == 'chance':
    return -1
2740
                              if player is None:
                                      return -4
2741
                              return int(player)
2742
                       \begin{tabular}{ll} \be
2743
                               if player_id == -1:
2744
                                      return 'chance
                              if player_id == -4:
2745
                                      return 'terminal
2746
                              return str(player_id)
2747
                       def get_rewards(state: State) -> list[float]:
    """Returns the rewards per player from their last action."""
2748
                               # Rewards are only given for a voluntary agreement. A forced agreement
# at the turn limit (MAX_TURNS) results in zero reward for both.
2749
                               if not state.get('agreement') or state['num_turns'] >= MAX_TURNS:
2750
                                      return [0.0] * NUM_PLAYERS
2751
                                rewards = []
                               for i in range(NUM_PLAYERS):
    player_values = state[f'player_{i}_values']
2752
2753
                                      player_share = state['agreement'][i]
# Reward is the total value of items received by the player.
```

```
2754
                     reward = sum(player_share.get(item, 0) * player_values.get(item, 0) for item in ITEMS)
2755
                     rewards.append(float(reward))
                return rewards
2756
            def get_legal_actions(state: State) -> list[Action]:
    """Returns all legal actions for the current player."""
2757
2758
                 player_id = get_current_player(state)
                 if player_id == -1: # Chance player
2759
                     return _CHANCE_ACTIONS
                 if player_id < 0: # Terminal state</pre>
2760
                     return []
2761
                actions = []
2762
                pool = state['pool']
2763
                # Generate all possible 'offer' actions by iterating through all item combinations. ranges = [range(pool.get(item, 0) + 1) for item in ITEMS]
2764
                 for combo in itertools.product(*ranges):
2765
                     quantities = {item: count for item, count in zip(ITEMS, combo)}
actions.append(f"player {player_id} offers {_format_quantities(quantities)}")
2766
                  'agree' is a legal move if at least one offer has been made by the opponent.
2767
                if state['offer_history']:
2768
                     actions.append(f"player {player_id} agrees")
2769
                return actions
2770
            def get_observations(state: State) -> list[PlayerObservation]:
                    'Returns the observation for each player, containing public and private information."""
2771
                 base_obs = {
2772
                      'current_player': state['current_player'],
                     'pool': state['pool'],
'num_turns': state['num_turns'],
2773
2774
                      'agreement': state['agreement'],
2775
                 is_terminal = get_current_player(state) == -4
2776
2777
                 # Determine the correct previous_offer based on game state
                 terminal_previous_offer = None
2778
                 if is_terminal:
2779
                     # In a terminal state, the "previous offer" is the one that was on the table
                     # before the final, accepted offer was made. This corresponds to the
# second-to-last offer in the history.
2780
                     if len(state['offer_history']) > 1:
2781
                         terminal_previous_offer = state['offer_history'][-2]
2782
                 for i in range(NUM_PLAYERS):
2783
                     obs = base_obs.copy()
obs['values'] = state[f'player_{i}_values']
2784
                     obs['my_player_id'] = i
2785
                     if is_terminal:
2786
                         obs['previous_offer'] = terminal_previous_offer
                     else:
2787
                         # In an active game, the previous offer is the last one made by the opponent.
                         opponent_id = 1 - i
2788
                         obs['previous_offer'] = next((
2789
                             offer for offer in reversed(state['offer_history']) if offer['player'] == opponent_id
                         ), None)
2790
                     observations.append(obs)
2791
                return observations
2792
             def resample_history(obs_action_history: list[tuple[PlayerObservation, Action | None]], player_id: int) -> list[Action]:
2793
                    "Stochastically samples one of many potential histories given a single player's perspective.
                 first_obs = obs_action_history[0][0]
2794
                 # Opponent's values are private and must be sampled randomly to create a possible history.
2795
                opponent_id = 1 - player_id

opponent_values = {item: random.randint(0, MAX_ITEM_VALUE) for item in ITEMS}
2796
2797
                 values = [{}, {}]
                 values[player_id] = first_obs['values']
2798
                 values[opponent_id] = opponent_values
2799
                 # Reconstruct the 'chance' action that started the game
2800
                 chance_action = (
                     f"{_format_quantities(first_obs['pool'])};"
2801
                     f"{_format_quantities(values[0])}
                     f"{_format_quantities(values[1])}'
2802
                 # Collect all known offers (own and opponent's) from the observation history.
2804
                known_offers = []
seen_turns = set()
2805
                 for obs, action in obs_action_history:
                     # Opponent's offers are seen in the 'previous_offer' field.
prev_offer = obs.get('previous_offer')
2806
2807
                     if prev_offer and prev_offer['num_turn'] not in seen_turns:
                         known_offers.append(prev_offer)
```

```
2808
                          seen_turns.add(prev_offer['num_turn'])
2809
                      # Own offers are reconstructed from the actions taken.
2810
                      if action and 'offers' in action:
                          # An offer action increments the turn number for the *next* state's observation. # The offer itself is recorded with this new turn number.
2811
2812
                           turn = obs['num_turns'] + 1
                          if turn not in seen_turns:
    quantities = _parse_quantities(action.split(' offers ')[1])
2813
                               known_offers.append({'num_turn': turn, 'player': player_id, 'quantities': quantities})
2814
                               seen_turns.add(turn)
2815
                  # Reconstruct the sequence of actions in chronological order.
2816
                 known_offers.sort(key=lambda x: x['num_turn'])
resampled_actions = [chance_action]
2817
                  {\tt resampled\_actions.extend(}
                      f"player {o['player']} offers {_format_quantities(o['quantities'])}"
2818
                      for o in known_offers
2819
                  # Add the final action if it was not an offer (e.g., 'agrees').
                 final_action = obs_action_history[-1][1]
if final_action and 'offers' not in final_action:
2821
                      resampled\_actions.append(final\_action)
2823
                 return resampled_actions
2824
             def value_function(state: dict[str, Any], player_id: int) -> float:
2825
                 ""Returns the value estimate for player_id in state.
2826
               For terminal states the function returns the true return. For ongoing play
2827
               the function should return a value estimate that reflect the winning potential
2828
               of the player with given player_id.
2829
               # 1 Handle Terminal States
               if get_current_player(state) == -4:
                 if player_id < 0 or player_id >= NUM_PLAYERS:
                    # For non-players like 'terminal' (-4), return 0.0
2831
                  # For active players, return the actual reward achieved.
                 return get_rewards(state)[player_id]
2833
               # --- Heuristic for Non-Terminal States ---
2834
               # 2. Basic Information
2835
               my_values = state.get(f'player_{player_id}_values')
2836
                 This can happen in the initial 'chance' state before values are assigned.
               if not my_values:
                 return 0.0
2838
               opponent_id = 1 - player_id
               pool = state['pool']
offer_history = state['offer_history']
2840
               current_turn_player = get_current_player(state)
2841
                # 3. Calculate Total Potential Value
               # The maximum value this player could get if they received all items. total\_my\_value\_ = sum(pool.get(item, 0) * my\_values.get(item, 0) for item in ITEMS)
2842
2843
               if total_my_value == 0:
                 return \ 0.0 \ \# \ If nothing in the pool is valuable, expected outcome is 0.
2844
               # 4. Define Baseline "Fair" Expectation
# A simple assumption that the player aims for about half the total value.
fair_value_estimate = total_my_value / 2.0
2845
2846
2847
               # 5. Core Heuristic Logic based on current negotiation status
               heuristic_value = fair_value_estimate # Default to fair split expectation
2848
               if not offer_history:
2849
                  # First turn, no offers yet. The best estimate is a fair split.
                 heuristic_value = fair_value_estimate
               else:
2851
                  last_offer = offer_history[-1]
                 if current_turn_player == player_id:
    # It's my turn to act.
2852
                    if last_offer['player'] == opponent_id:
2853
                      # Opponent made the last offer. I can agree or counter
# Calculate the value of their offer to me.
2854
                      their_proposed_share = last_offer['quantities']
2855
                      my_share_if_agree = {
                           item: pool.get(item, 0) - their_proposed_share.get(item, 0)
2856
                           for item in ITEMS
                      value_on_table = sum(
2858
                          my_share_if_agree.get(item, 0) * my_values.get(item, 0) for item in ITEMS
2859
2860
                      # My position's value is the better of what I can get now
                      # versus my general expectation from continued negotiation.
                    heuristic_value = max(value_on_table, fair_value_estimate)
else: # last_offer['player'] == player_id
```

```
2862
                            # The last offer was mine, but the opponent didn't agree.
2863
                           # My turn again means my previous offer was implicitly rejected.
# Fall back to the baseline expectation.
2864
                            heuristic_value = fair_value_estimate
                      elif current_turn_player == opponent_id:
# It's the opponent's turn. They are considering my last offer.
if last_offer['player'] == player_id:
    my_proposed_share = last_offer['quantities']
# The value of the state is the value of the offer I'm hoping they accept.
2865
2866
2867
                            value_of_my_offer = sum(
2868
                                 \label{eq:my_proposed_share.get} $$  \mbox{my_proposed\_share.get(item, 0)} *  \mbox{my_values.get(item, 0)} $$  \mbox{for item in ITEMS} 
2869
2870
                           heuristic_value = value_of_my_offer
2871
                   \# 6. Apply Time Pressure Discount \# As turns run out, the risk of getting 0 from a forced agreement increases. \# This discounts the potential future value accordingly.
2872
                   if state['num_turns'] >= MAX_TURNS:
    return 0.0 # Game is over or will be forced to 0 reward on next action.
2873
2874
                   turns_left = MAX_TURNS - state['num_turns']
# A sqrt factor makes the discount less severe in early turns.
pressure_factor = (turns_left / MAX_TURNS) ** 0.5
2875
2876
2877
                   return float(heuristic_value * pressure_factor)
2878
2879
2880
2881
2882
2883
2884
2885
2886
2887
2888
2889
2890
2891
2892
2893
2894
2895
2896
2897
2898
2899
2900
2901
2902
2903
2904
2905
2906
2907
2908
2909
2910
2911
2912
2913
2914
2915
```

```
2916
                 I.3 Bargaining (imperfect information, closed deck)
2917
2918
                 import copy
2919
                 import itertools
                 import rand
2920
                 from typing import Any, Dict. List. Tuple, Optional
                 # Type definitions
2922
                 Action = str
                 State = Dict[str, Any]
2923
                PlayerObservation = Dict[str, Any]
2924
                 # Game constants
2925
                 NUM PLAYERS = 2
                 ITEMS = sorted(['X', 'Y', 'Z'])
2926
                 MAX_TURNS = 10
POOL_VALUES = range(1, 6)
2927
                 ITEM_VALUES = range(0, 7)
                 # Create a fixed, reproducible set of possible game scenarios for the chance node.
2929
                 _CHANCE_OUTCOMES = []
                 _chance_rng = random.Random(0)
2930
                 for _ in range(20):
                      pool = {item: _chance_rng.choice(POOL_VALUES) for item in ITEMS)
p0_values = {item: _chance_rng.choice(ITEM_VALUES) for item in ITEMS}
p1_values = {item: _chance_rng.choice(ITEM_VALUES) for item in ITEMS}
2931
2932
2933
                       \begin{aligned} & \text{pool\_str} = ",".j \text{oin}(f"\{k\}=\{v\}" \text{ for } k, \text{ } v \text{ in sorted(pool.items()))} \\ & \text{p0\_str} = ",".j \text{oin}(f"\{k\}=\{v\}" \text{ for } k, \text{ } v \text{ in sorted(p0\_values.items()))} \\ & \text{p1\_str} = ",".j \text{oin}(f"\{k\}=\{v\}" \text{ for } k, \text{ } v \text{ in sorted(p1\_values.items()))} \\ & \text{\_CHANCE\_OUTCOMES.append}(f"pool:\{pool\_str\};p0\_values:\{p0\_str\};p1\_values:\{p1\_str\}")} \end{aligned} 
2934
2935
2936
                 def _parse_offer_action(action: Action) -> Tuple[int, Dict[str, int]]:
    """Parses an offer action string into player ID and quantities."""
2937
                    parts = action.split()
                   player_id = int(parts[1])
quantities = {item: int(q) for item, q in zip(ITEMS, parts[3].split(','))}
2938
2939
                    return player_id, quantities
2940
                def _calculate_reward(bundle: Dict[str, int], values: Dict[str, int]) -> float:
    """Calculates the total value of a bundle of items for a player."""
    return sum(bundle.get(item, 0) * values.get(item, 0) for item in ITEMS)
2941
2942
                 def _reconstruct_offer_action(offer: Dict[str, Any]) -> Action:
2943
                    """Reconstructs an offer action string from an offer dictionary."""
quantities_str = ",".join(str(offer['quantities'].get(item, 0)) for item in ITEMS)
2944
                    return f"player {offer['player']} offers {quantities_str}"
2945
                 def apply_action(state: State, action: Action) -> State:
2946
                        Returns the new state after an action has been taken."""
                   new_state = copy.deepcopy(state)
2947
                    if state.get('current_player') == 'chance':
2948
                         # Initialize the game state from the chance node action.
2949
                         parts = action.split(';')
                         parts = action.split(';')
new_state['pool'] = {p.split('=')[0]: int(p.split('=')[1]) for p in parts[0].split(':')[1].split(',')}
new_state['player_0_values'] = {p.split('=')[0]: int(p.split('=')[1]) for p in parts[1].split(':')[1].split(',')}
new_state['player_1_values'] = {p.split('=')[0]: int(p.split('=')[1]) for p in parts[2].split(':')[1].split(',')}
new_state['current_player'] = 0
2950
2951
                         return new_state
2952
2953
                    if 'agrees' in action:
                         # An agreement is reached. The game becomes terminal. last_offer = new_state['offer_history'][-1]
2954
                         offerer_id = last_offer['player']
2955
                         offerer_bundle = last_offer['quantities']
                         accepter_bundle = {item: new_state['pool'][item] - offerer_bundle.get(item, 0) for item in ITEMS}
2956
2957
                         agreement = [{}, {}]
agreement[offerer_id] = offerer_bundle
2958
                          agreement[1 - offerer_id] = accepter_bundle
2959
                         new state['agreement'] = agreement
2960
                          new_state['current_player'] = None
                    elif 'offers' in action:
    # An offer is made. Increment turn count and switch player.
2961
                         player_id, quantities = _parse_offer_action(action)
new_state['num_turns'] += 1
offer = {'num_turn': new_state['num_turns'], 'player': player_id, 'quantities': quantities}
2962
2963
                         new_state['offer_history'].append(offer)
2964
                         if new_state['num_turns'] >= MAX_TURNS:
2965
                               new_state['current_player'] = None # End game if turn limit reached.
                         else:
2966
                               new_state['current_player'] = 1 - player_id
2967
                   return new_state
2968
                 def get_current_player(state: State) -> int:
    """Returns current player, with -1 for chance and -4 for terminal."""
2969
                    if state.get('current_player') == 'chance':
```

```
2970
                               return -
2971
                       if state.get('current_player') is None or state['agreement'] or state['num_turns'] >= MAX_TURNS:
                              return -
2972
                        return state['current_player']
2973
                     \begin{tabular}{ll} \be
2974
                        return {-1: 'chance', -4: 'terminal'}.get(player_id, str(player_id))
2975
                    def get_rewards(state: State) -> list[float]:
2976
                        """Returns rewards. Rewards are 0 if the game ends due to the turn limit."""
if state['num_turns'] >= MAX_TURNS or not state['agreement']:
2977
                               return [0.0] * NUM_PLAYERS
2978
                        p0_reward = _calculate_reward(state['agreement'][0], state['player_0_values'])
2979
                        p1_reward = _calculate_reward(state['agreement'][1], state['player_1_values'])
2980
                        return [float(p0_reward), float(p1_reward)]
2981
                    def get_legal_actions(state: State) -> list[Action]:
                          ""Returns legal actions that can be taken in current state."""
                        player = get_current_player(state)
2983
                        if player == -4:
2984
                               return []
                        if player == -1:
2985
                              return _CHANCE_OUTCOMES
2986
                        actions = []
                        if state['num_turns'] > 0:
2987
                              actions.append(f"player {player} agrees")
2988
                        # Generate all possible offer combinations based on the item pool.
2989
                        pool = state['pool']
                        quantity_ranges = [range(pool.get(item, 0) + 1) for item in ITEMS]
2990
                        quantities in itertools.product(*quantity_ranges):
    q_str = ",".join(map(str, quantities))
    actions.append(f"player {player} offers {q_str}")
2991
2992
                       return actions
2993
                    def get_observations(state: State) -> list[PlayerObservation]:
2994
                        """Returns the observation for each player.
observations = []
2995
                        player_at_turn = get_current_player(state)
2996
                        is_terminal = (player_at_turn == -4)
2997
                        for i in range(NUM_PLAYERS):
2998
                              previous_offer = None
                              # In a terminal state with an agreement, the "previous offer" is the one before the accepted one if is_terminal and state['agreement'] and len(state['offer_history']) > 1:
                                     previous_offer = state['offer_history'][-2]
3000
                               elif state['offer_history']:
3001
                                     previous_offer = state['offer_history'][-1]
3002
                                      = {
'my_player_id': i,
'pool': state['pool'],
'values': state[f'player_(i)_values'],
'num_turns': state['num_turns'],
'agreement': state['agreement'],
3003
3004
3005
                                      'previous_offer': previous_offer,
'current_player': str(player_at_turn) if player_at_turn >= 0 else None,
3006
3007
                              observations.append(obs)
3008
                        return observations
3009
                    def resample_history(obs_action_history: list[tuple[PlayerObservation, Action | None]], player_id: int. last_is_terminal: bool) ->
                   test[Action]:
    list[Action]:
    """Stochastically sample one of many potential histories of actions for all players."""
3010
3011
                        first\_obs = obs\_action\_history[0][0]
3012
                        opponent\_values = \{'X': 3, 'Y': 3, 'Z': 4\} # Assume fixed opponent values for reproducibility.
3013
                        p_{vals} = [{}, {}]
                        p_vals[player_id] = first_obs['values']
3014
                        p_vals[1 - player_id] = opponent_values
3015
                        3016
3017
                        yield f"pool:{pool_str};p0_values:{p0_str};p1_values:{p1_str}
3018
                        # 2. Reconstruct the interleaved game actions from the player's perspective.
3019
                        last opponent turn vielded = 0
                        my_last_action = None
3020
                        for obs, action in obs_action_history:
                              if action:
3021
                                     my_last_action = action
3022
                              if obs.get('previous_offer'):
3023
                                      offer = obs['previous_offer']
                                      # Only yield opponent offers that haven't been yielded yet to avoid duplication.
```

```
3024
                          if offer['player'] != player_id and offer['num_turn'] > last_opponent_turn_yielded:
3025
                               yield _reconstruct_offer_action(offer)
                              last_opponent_turn_yielded = offer['num_turn']
3026
3027
                     if action:
                         vield action
3028
                          if 'agrees' in action:
                              return
3029
                # 3. Deduce the final hidden actions if the game ended with an agreement not initiated by the player.
3030
                if last is terminal:
3031
                     last_obs, last_action = obs_action_history[-1]
                     if last_action is None and last_obs['agreement']:
    agreement = last_obs['agreement']
    _, my_last_quantities = _parse_offer_action(my_last_action)
3032
3033
                          # Case A: Opponent agreed to my last offer. My bundle in the agreement matches my last offer.
3034
                         # Case A. Opponent agreed to my last over. In would In the agreement matches my
iff my_last_quantities == agreement[player_id]:
    yield f"player {1 - player_id} agrees"
# Case B: Opponent made a counter-offer, which I would have implicitly agreed to.
3035
3036
                              opponent_id = 1 - player_id
3037
                              opponent_bundle = agreement[opponent_id]
quantities_str = ",".join(str(opponent_bundle.get(item, 0)) for item in ITEMS)
yield f"player (opponent_id) offers (quantities_str)"
3038
3039
                              yield f"player {player_id} agrees"
3040
              from typing import Any, Dict, List
3041
3042
              def value_function(state: dict[str, Any], player_id: int) -> float:
3043
                  ""Returns the value estimate for player_id in state.
3044
                For terminal states the function returns the true return. For ongoing play
3045
                the function should return a value estimate that reflect the winning potential
                of the player with given player_id.
3046
                \mbox{\# Game constants} and helper functions defined in local scope for self-containment.
3047
                ITEMS = sorted(['X', 'Y', 'Z'])
                MAX_TURNS = 10
3048
                NUM_PLAYERS = 2
3049
                \label{lem:def_calculate_reward(bundle: Dict[str, int], values: Dict[str, int]) -> float: """Calculates the total value of a bundle of items for a player."""
3050
                  return sum(bundle.get(item, 0) * values.get(item, 0) for item in ITEMS)
3051
3052
                def _get_current_player_internal(state_dict: Dict[str, Any]) -> int:
                   """Determines the current player or if the state is terminal.""
3053
                  current_player = state_dict.get('current_player')
agreement = state_dict.get('agreement')
num_turns = state_dict.get('num_turns', 0)
3054
3055
                  is\_terminal = (
                       Let minal = \sqrt{\ } current_player is None or (agreement and isinstance(agreement, list) and len(agreement) > 0) or
3056
3057
                       num_turns >= MAX_TURNS
3058
3059
                  \quad \textbf{if} \ is\_terminal: \\
                     return -4 # Terminal node code
3060
                   if current_player == 'chance'
                                 # Chance node code
3061
                  return int(current_player)
3062
                def _get_rewards_internal(state_dict: Dict[str, Any]) -> List[float]:
3063
                  """Calculates rewards for all players in a terminal state.
agreement = state_dict.get('agreement')
3064
                   num_turns = state_dict.get('num_turns', 0)
3065
                    No reward if the game ends due to turn limit or no agreement is made
3066
                  if num_turns \geq MAX_TURNS or not agreement or (isinstance(agreement, list) and len(agreement) == 0):
                     return [0.0] * NUM_PLAYERS
3067
                  p0_reward = _calculate_reward(agreement[0], state_dict['player_0_values'])
3068
                  p1_reward = _calculate_reward(agreement[1], state_dict['player_1_values'])
3069
                  return [float(p0_reward), float(p1_reward)]
3070
                # --- Main value function logic begins ---
3071
                current_player_code = _get_current_player_internal(state)
3072
                # 1. Handle Terminal States: Return the exact final reward.
3073
                if current_player_code == -4:
3074
                  \label{eq:continuous} \textbf{if player\_id} \, < \, \textbf{0:} \quad \textit{\# MCTS may query the value for the terminal node itself.}
3075
                   return _get_rewards_internal(state)[player_id]
3076
                # 2. Handle Non-Terminal States: Return a heuristic-based value estimate.
3077
                my_values = state[f'player_{player_id}_values']
```

```
3078
                pool = state['pool']
3079
                offer_history = state.get('offer_history', [])
3080
                \# Heuristic for the start of the game (no offers yet). \# A neutral assumption is that the player can achieve half of their maximum possible value.
3081
                if not offer_history:
                   my_total_pool_value = _calculate_reward(pool, my_values)
3082
                   return my_total_pool_value / 2.0
3083
                last_offer = offer_history[-1]
3084
                # Case A: It's the opponent's turn. This means I made the last offer.
3085
                 # The value of my last offer is a good estimate of my current potential, as it reflects my aspiration.
3086
                if current_player_code != player_id:
   my_bundle_in_my_last_offer = last_offer['quantities']
3087
                   my_value_of_my_offer = _calculate_reward(my_bundle_in_my_last_offer, my_values)
return float(my_value_of_my_offer)
3088
                # Case B: It's my turn. The opponent made the last offer.
# My potential lies between what they offered and what I last asked for.
3089
3090
                   He Calculate the value of their offer to me. This is a concrete value I can achieve by accepting. offered_bundle_to_opponent = last_offer['quantities']
3091
3092
                       item: pool.get(item, 0) - offered_bundle_to_opponent.get(item, 0) for item in ITEMS
3093
3094
                   value_of_their_offer_to_me = _calculate_reward(implied_bundle_to_me, my_values)
3095
                   \ensuremath{\text{\#}} Find my last offer to gauge my own aspiration level.
3096
                   my_aspiration = -1.0
                   for offer in reversed(offer_history):
3097
                     if offer['player'] == player_id:
    my_bundle_in_my_last_offer = offer['quantities']
    my_aspiration = _calculate_reward(my_bundle_in_my_last_offer, my_values)
3098
3099
                       break
3100
                   \# If I haven't made an offer yet, my aspiration defaults to the initial 50/50 baseline
                   if my_aspiration < 0:</pre>
3101
                     my_total_pool_value = _calculate_reward(pool, my_values)
                     my_aspiration = my_total_pool_value / 2.0
3102
                   # The heuristic is the midpoint between their offer and my aspiration, representing a likely compromise point.
3103
                   \label{eq:heuristic_value} \mbox{heuristic\_value} = \mbox{(value\_of\_their\_offer\_to\_me + my\_aspiration)} \mbox{/ } 2.0 \\ \mbox{return float(heuristic\_value)}
3104
3105
3106
3107
3108
3109
3110
3111
3112
3113
3114
3115
3116
3117
3118
3119
3120
3121
3122
3123
3124
3125
3126
3127
3128
3129
3130
3131
```