

BetaZero: Belief-State Planning for Long-Horizon POMDPs using Learned Approximations

Anonymous submission

Abstract

Real-world planning problems, including autonomous driving and sustainable energy applications like carbon storage and resource exploration, have recently been modeled as partially observable Markov decision processes (POMDPs) and solved using approximate methods. To solve high-dimensional POMDPs in practice, state-of-the-art methods use online planning with problem-specific heuristics to reduce planning horizons and make the problems tractable. Algorithms that learn approximations to replace heuristics have recently found success in large-scale fully observable domains. The key insight is the combination of online Monte Carlo tree search with offline neural network approximations of the optimal policy and value function. In this work, we bring this insight to partially observed domains and propose *BetaZero*, a belief-state planning algorithm for high-dimensional POMDPs. BetaZero learns offline approximations that replace heuristics to enable online decision making in long-horizon problems. We address several challenges inherent in large-scale partially observable domains; namely challenges of transitioning in stochastic environments, prioritizing action branching with a limited search budget, and representing beliefs as input to the network. To formalize the use of all limited search information we train against a novel Q -weighted policy vector target. We test BetaZero on various well-established benchmark POMDPs found in the literature and a real-world, high-dimensional problem of critical mineral exploration. Experiments show that BetaZero outperforms state-of-the-art POMDP solvers on a variety of tasks.¹

1 Introduction

Optimizing sequential decisions in real-world settings is challenging due to the inherent uncertainty about the true state of the environment. Modeling such problems as partially observable Markov decision processes (POMDPs) has shown recent success in autonomous driving (Wray et al. 2021), robotics (Lauri, Hsu, and Pajarinen 2022), and aircraft collision avoidance (Kochenderfer, Holland, and Chrysanthacopoulos 2012). Solving large or continuous POMDPs require approximations in the form of state-space discretizations or modeling assumptions, e.g., assuming full observability. Although these approximations are necessary when making fast decisions in a short time horizon, scaling these solutions to long-horizon problems is challenging

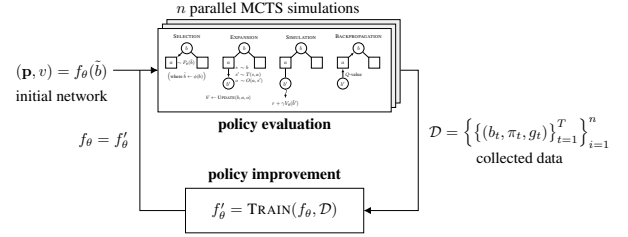


Figure 1: The *BetaZero* POMDP policy iteration algorithm.

(Shani, Pineau, and Kaplow 2013). Recently, POMDPs have been used to model large-scale information gathering problems such as carbon capture and storage (CCS) (Corso et al. 2022; Wang et al. 2023), remediation for groundwater contamination (Wang et al. 2022), and critical mineral exploration for battery metals (Mern and Caers 2023) and are solved using online tree search methods such as DESPOT (Ye et al. 2017) and POMCPOW (Sunberg and Kochenderfer 2018). The performance of these online methods rely on heuristics for action selection (to reduce search tree expansion) and heuristics to estimate the value function (to avoid expensive rollouts and reduce tree search depth). Without heuristics, online methods have difficulty planning for long-term information acquisition to reason about uncertain future events. Thus, algorithms to solve high-dimensional POMDPs need to be designed to learn heuristic approximations to enable decision making in long-horizon problems.

Related work. Algorithms to solve high-dimensional, fully observable Markov decision processes (MDPs) learn approximations to replace problem-specific heuristics. Silver et al. introduced the *AlphaZero* algorithm for large, deterministic MDPs and showed considerable success in games such as Go, chess, shogi, and Atari (Silver et al. 2018; Schrittwieser et al. 2020). The success is attributed to the combination of online Monte Carlo tree search (MCTS) and a neural network that approximates the optimal value function and the offline policy. Extensions of AlphaZero and the model-free variant *MuZero* (Schrittwieser et al. 2020) have already addressed several challenges when applying to broad classes of MDPs. For large or continuous action spaces, Hubert et al. introduced a policy improvement algorithm called *Sampled MuZero* that samples an action set

¹Code: <https://placeholder-url-for-github.com/BetaZero.jl>

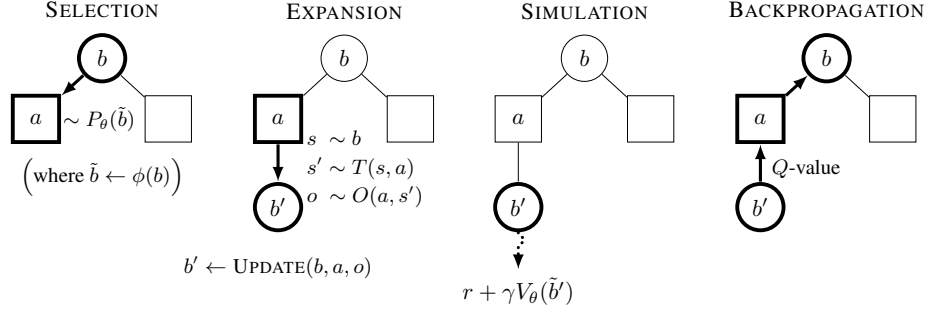


Figure 2: The four stages of MCTS in *BetaZero* using a value V_θ and policy P_θ network (the *policy evaluation* step in fig. 1).

of an *a priori* fixed size every time a node is expanded. Antonoglou et al. introduced *Stochastic MuZero* that extends MuZero to games with stochastic transitions but assumes a finite set of possible next states so that each transition can be associated with a chance outcome. Applying these algorithms to large or continuous spaces remains challenging.

To handle partial observability in stochastic games, Ozair et al. combine VQ-VAEs with MuZero to encode future discrete observations into latent variables. Other approaches incorporate partial observability by inputting action-observation histories directly into the network (Kimura, Sakamoto, and Sogabe 2020; Vinyals et al. 2019). Similarly, Igl et al. introduce a method to learn a belief representation within the network when the agent is only given access to histories. Their work focuses on the *reinforcement learning* (RL) domain and they show that a belief distribution can be represented as a latent state in the learned model. The FORBES algorithm (Chen et al. 2022) builds a normalizing flow-based belief and learns a policy through an actor-critic RL algorithm. Methods to learn the belief are necessary when a prior belief model is not available. When such models *do* exist, as is the case with many geological applications, using the models can be valuable for long-term planning. Hoel et al. apply AlphaGo Zero (Silver et al. 2017) to an autonomous driving POMDP using the most-likely state as the network input but overlook significant belief uncertainty information. Planning under uncertainty is crucial for learning effective POMDP policies.

Planning vs. reinforcement learning. In POMDP *planning*, models of the transitions, rewards, and observations are known. In contrast, in the model-based *reinforcement learning* (RL) domain, these models are learned along with a policy or value function (Sutton and Barto 2018). A difference between these settings is that RL algorithms reset the agent and learn through experience, while planning algorithms, like MCTS, must consider future trajectories from any state. When RL problems have deterministic state transitions, they can be cast as a planning problem by replaying the full state trajectory along a tree path, which may be prohibitively expensive for long-horizon problems. Both settings are closely related and pose interesting research challenges. Specifically, sequential planning over given models in high-dimensional, long-horizon POMDPs remains challenging (Lauri, Hsu, and Pajarinen 2022).

Online POMDP planning. Existing online POMDP planning algorithms rely on problem-specific heuristics for good performance. Sunberg and Kochenderfer introduced the POMCPOW planning algorithm that iteratively builds a particle set belief within the tree, designed for continuous observations. In practice, POMCPOW relies on heuristics for value function estimation and action selection (e.g., work from Mern and Caers). Wu et al. introduced AdaOPS that adaptively approximates the belief through particle filtering and maintains value function bounds that are initialized with heuristics (e.g., solving the MDP or using expert policies). The major limitation of existing solvers is the reliance on heuristics to make long-horizon POMDPs tractable, which may not scale to high-dimensional problems.

Our work avoids heuristics and uses the insight of combining online MCTS with offline neural network approximations to solve high-dimensional POMDPs. We propose the *BetaZero* belief-state planning algorithm that plans in belief space using no problem-specific heuristics, only using information from the search. We address partially observable challenges in large discrete action spaces and continuous state and observation spaces. In stochastic environments, BetaZero uses progressive widening (Couëtoux et al. 2011) to limit belief-state expansion—only using information available in the tree. When planning in belief space, computationally expensive belief updates limit the search budget in practice. Therefore, information from the policy network is used to prioritize branching on promising actions. We introduce a novel Q -weighted policy vector target that formalizes the use of all the information seen during the limited search. To capture state uncertainty, the network receives a parametric approximation of the belief \tilde{b} as input. BetaZero uses the learned policy network P_θ to reduce search breadth and the learned value estimate V_θ to reduce search depth to enable long-horizon online planning (fig. 2).

2 Background

A partially observable Markov decision process (POMDP) is a model for sequential decision making problems where the true state is unobservable. Defined by the tuple $\langle \mathcal{S}, \mathcal{A}, \mathcal{O}, T, R, O, \gamma \rangle$, POMDPs are an extension to the Markov decision process (MDP) used in reinforcement learning and planning with the addition of an observation space \mathcal{O} (where $o \in \mathcal{O}$) and observation model $O(o | a, s')$.

Given a current state $s \in \mathcal{S}$ and taking an action $a \in \mathcal{A}$, the agent transitions to a new state s' using the transition model $s' \sim T(\cdot | s, a)$. Without access to the true state, an observation is received $o \sim O(\cdot | a, s')$ and used to update the belief b over the possible next states s' to get the posterior

$$b'(s') \propto O(o | a, s') \int_{s \in \mathcal{S}} T(s' | s, a) b(s) ds. \quad (1)$$

The non-parametric *particle set* belief can represent a broad range of distributions (Thrun, Burgard, and Fox 2005) and Lim et al. show that optimality guarantees exist in finite-sample particle-based POMDP approximations. Despite choosing to study particle-based beliefs, our work generalizes well to problems with parametric beliefs.

A stochastic POMDP policy $\pi(a | b)$ is defined as the distribution over actions given the current belief b . After taking an action $a \sim \pi(\cdot | b)$, the agent receives a reward r from the environment according to the reward function $R : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ or $R : \mathcal{S} \times \mathcal{A} \times \mathcal{S} \rightarrow \mathbb{R}$ using the next state.

Belief-state MDPs. In *belief-state* MDPs, the POMDP is converted to an MDP by treating the belief as a state (Kochenderfer, Wheeler, and Wray 2022). The reward function then becomes a weighted sum of the state-based reward:

$$R_b(b, a) = \int_{s \in \mathcal{S}} b(s) R(s, a) ds \quad (2)$$

$$\approx \sum_{s \in b} b(s) R(s, a) \quad (3)$$

The belief-state MDP shares the same action space as the POMDP and operates over a belief space \mathcal{B} that is the simplex over the state space \mathcal{S} . The belief-MDP defines a new belief-state transition function $b' \sim T_b(\cdot | b, a)$ as:

$$\begin{aligned} s &\sim b(\cdot) & s' &\sim T(\cdot | s, a) \\ o &\sim O(\cdot | a, s') & b' &\leftarrow \text{UPDATE}(b, a, o) \end{aligned} \quad (4)$$

where the belief update can be done using a particle filter (Gordon, Salmond, and Smith 1993). Therefore, the belief-state MDP is defined by the tuple $\langle \mathcal{B}, \mathcal{A}, T_b, R_b, \gamma \rangle$ with the finite-horizon discount factor $\gamma \in [0, 1)$ that controls the effect that future rewards have on the current action.

The objective to solve belief-MDPs is to find a policy π that maximizes the *value function* from an initial belief b_0 :

$$V^\pi(b_0) = \mathbb{E}_\pi \left[\sum_{t=1}^{\infty} \gamma^t R_b(b_t, a_t) \mid b_t \sim T_b, a_t \sim \pi \right] \quad (5)$$

Instead of explicitly constructing a policy over all beliefs, online planning algorithms estimate the next best action through a planning procedure, often a best-first tree search.

Monte Carlo tree search (MCTS)

Monte Carlo tree search (Coulom 2007; Browne et al. 2012) is an online, recursive, best-first tree search algorithm to determine the approximately optimal action to take from a given root state of an MDP. Extensions to MCTS have been applied to POMDPs through several algorithms: *partially observable Monte Carlo planning* (POMCP) treats

the state nodes as histories h of action-observation trajectories (Silver and Veness 2010), *POMCP with observation widening* (POMCPOW) construct weighted particle sets at the state nodes and extends POMCP to fully continuous domains (Sunberg and Kochenderfer 2018), and *particle filter trees* (PFT) and *information particle filter trees* (IPFT) treat the POMDP as a belief-state MDP and plan directly over the belief-state nodes using a particle filter (Fischer and Tas 2020). All variants of MCTS execute the following four key steps. In this section we use s to represent the state, the history h , and the belief state b and refer to them as “the state”.

1. **Selection.** During *selection*, an action is selected from the children of a state node based on criteria that balances exploration and exploitation and is considered a multi-armed bandit problem (Munos 2014). The *upper-confidence tree* algorithm (UCT) (Kocsis and Szepesvári 2006) is a commonly used criterion that selects an action that maximizes the upper-confidence bound

$$Q(s, a) + c \sqrt{\frac{\log N(s)}{N(s, a)}} \quad (6)$$

where $Q(s, a)$ is the Q -value estimate for state-action pair (s, a) , the visit count for that state-action pair is $N(s, a)$, the total visit count is $N(s) = \sum_a N(s, a)$ for the children $a \in A(s)$, and $c \geq 0$ is an exploration constant. Rosin introduced the *upper-confidence tree with predictor* algorithm (PUCT), modified by Silver et al., where a predictor $P(s, a)$ guides the exploration towards promising branches and selects an action that maximizes

$$Q(s, a) + c \left(P(s, a) \frac{\sqrt{N(s)}}{1 + N(s, a)} \right). \quad (7)$$

When dealing with a small discrete action space, the selection step often first branches on all actions and then selects based on some criterion; as is the case for AlphaZero that first expands on all legal moves (Silver et al. 2016, 2018). For large or continuous action spaces, algorithms attempt to balance branching and selection.

2. **Expansion.** In the *expansion* step, the selected action is taken in simulation and the transition model $T(s' | s, a)$ is sampled to determine the next state s' . When the transitions are deterministic, the child node is always a single state. If the transition dynamics are stochastic, techniques to balance the branching factor such as progressive widening (Couëtoux et al. 2011) and state abstraction refinement (Sokota et al. 2021) have been proposed.
3. **Rollout/Simulation.** In the *rollout* step, also called the *simulation* step due to recursively simulating the MCTS tree expansion, the value is estimated through the execution of a rollout policy until termination or using heuristics to approximate the value function from the given state s' . Silver et al. replace the expensive rollouts done by AlphaGo with a value network lookup in AlphaGo Zero and AlphaZero (Silver et al. 2016, 2017, 2018).
4. **Backpropagation.** Finally, during the *backpropagation* step, the Q -value estimate from the rollout is propagated up the path in the tree as a running average.

Root node action selection. After repeating the four steps of MCTS, the final action is selected from the children $a \in A(s)$ of the root state s and executed in the environment. One way to select the best root node action, referred to as the *robust child* (Schadd 2009; Browne et al. 2012), selects the action with the highest visit count as $\arg \max_a N(s, a)$. Sampling from the normalized counts, exponentiated by an exploratory temperature τ , is also common (Silver et al. 2017). Another method uses the highest estimated Q -value as $\arg \max_a Q(s, a)$. Both criteria have been shown to have problem-based trade-offs (Browne et al. 2012).

Double progressive widening. To handle stochastic state transitions and large or continuous state and action spaces, double progressive widening (DPW) balances between sampling new nodes to expand on or selecting from existing nodes already in the tree (Couëtoux et al. 2011). Two hyperparameters $\alpha \in [0, 1]$ and $k \geq 0$ control the branching factor. If the number of actions tried from state s is less than $kN(s)^\alpha$, then a new action is sampled from the action space and added as a child of node s . Likewise, if the number of expanded states from state-action node (s, a) is less than $kN(s, a)^\alpha$, then a new state is sampled from the transition function $s' \sim T(\cdot | s, a)$ and added as a child of that state-action node. If the state widening condition is not met, then a next state is sampled from the existing children. To encourage widening, let $k \rightarrow \infty$ and $\alpha \rightarrow 1$, e.g., in highly stochastic problems (Sokota et al. 2021), and to discourage widening let $k \rightarrow 1$ and $\alpha \rightarrow 0$ (Moss et al. 2020).

Note that in the following sections we will refer to the belief state as b and the true (hidden) state as s .

3 Proposed Algorithm: BetaZero

We introduce the *BetaZero* POMDP planning algorithm that replaces heuristics with learned approximations of the optimal policy and value function. BetaZero is a belief-space policy iteration algorithm with two offline steps:

1. **Policy evaluation:** Evaluate the current value and policy network through n parallel episodes of MCTS and collect training data: $\mathcal{D} = \{ \{ (b_t, \pi_t, g_t) \}_{t=1}^T \}_{i=1}^n$
2. **Policy improvement:** Improve the estimated value function and policy by retraining the neural network parameters θ with data from the most recent MCTS simulations.

The policy vector over actions $\mathbf{p} = P_\theta(\tilde{b}, \cdot)$ and the value $v = V_\theta(\tilde{b})$ are combined into a single network with two output heads $(\mathbf{p}, v) = f_\theta(\tilde{b})$; we refer to P_θ and V_θ separately.

During *policy evaluation*, training data is collected from the outer POMDP loop. The belief b_t and the tree policy π_t are collected for each time step t . At the end of each episode, the returns $g_t = \sum_{i=t}^T \gamma^{(i-t)} r_i$ are computed from the set of observed rewards for all time steps up to a terminal horizon T . Traditionally, MCTS algorithms use a tree policy π_t that is proportional to the root node visit counts of its children actions $\pi_t(b_t, a) \propto N(b_t, a)^{1/\tau}$. The counts are sampled after exponentiating with a temperature τ to encourage exploration but evaluated online with $\tau \rightarrow 0$ to return the maximizing action (Silver et al. 2017). In certain settings, relying solely on visit counts may overlook crucial information.

Policy vector as Q -weighted counts. When planning over belief states, expensive belief updates occur in the tree search and thus a limited MCTS budget may be used. Therefore, the visit counts may not converge towards an optimal strategy as the budget may be spent on exploration. Danihelka et al. and Czech, Korus, and Kersting suggest using knowledge of the Q -values from search in MCTS action selection. Using only tree information, we incorporate Q -values and train against the policy $\pi_t(b_t, a)$ proportional to

$$\left(\left(\frac{\exp Q(b_t, a)}{\sum_{a'} \exp Q(b_t, a')} \right)^{z_q} \left(\frac{N(b_t, a)}{\sum_{a'} N(b_t, a')} \right)^{z_n} \right)^{1/\tau} \quad (8)$$

which is then normalized to get a valid probability distribution. Equation (8) simply weights the visit counts by the softmax Q -value distribution with parameters $z_q \in [0, 1]$ and $z_n \in [0, 1]$ defining the influence of the values and the visit counts, respectively. If $z_q = z_n = 1$, then the influence is equal and if $z_q = z_n = 0$, then the policy becomes uniform. Once the tree search finishes, the final action is selected as $a \sim \pi_t(b_t, \cdot)$ and returns the $\arg \max$ when $\tau \rightarrow 0$.

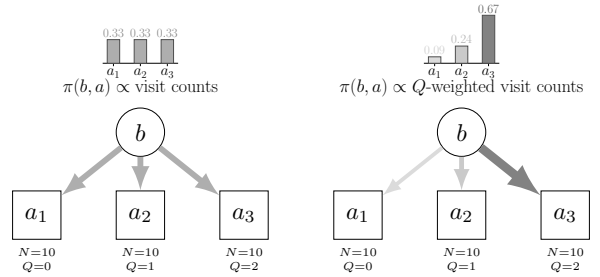


Figure 3: An illustrative example of when using a small MCTS budget with high exploration and collecting policy data based purely on visit counts (left) would perform worse than weighting the counts based on Q -values (right). In this example, MCTS had a budget of 30 iterations and visited each action 10 times. When training on the visit-counts, we miss useful Q -value information and select the worst performing action a_1 . When using both the Q -values and visit counts, we incorporate both what the tree search *focused on* and the *values it found*. This also accounts for uncertainty by weighting the Q -value by the number of samples in the estimate. An ablation study in the appendix tests this idea.

Loss function. Using the latest collected data, the *policy improvement* step retrains the policy network head using the cross-entropy loss $\mathcal{L}_{P_\theta}(\pi_t, \mathbf{p}_t) = -\pi_t^\top \log \mathbf{p}_t$. The value network head is simultaneously trained to fit the returns g_t using mean-squared error (MSE) or mean-absolute error (MAE) to predict the value of the belief b_t . The choice of value loss function \mathcal{L}_{V_θ} depends on the characteristics of the return distribution. In sparse reward problems, MAE is a better choice as the distribution is closer to Laplacian (Hodson 2022). When the reward is distributed closer to Gaussian, then MSE is more suitable (Chai and Draxler 2014). The final loss function combines the value and policy losses with added L_2 -regularization scaled by λ :

$$\ell_{\beta_0} = \mathcal{L}_{V_\theta}(g_t, v_t) + \mathcal{L}_{P_\theta}(\pi_t, \mathbf{p}_t) + \lambda \|\theta\|^2 \quad (9)$$

Algorithm 1: BetaZero action progressive widening.

```

1 function ACTIONSELECTION( $f_\theta, b$ )
2    $\tilde{b} \leftarrow \phi(b)$   $\triangleright$  belief representation eq. (10)
3   if  $|A(b)| \leq kN(b)^\alpha$   $\triangleright$  action progressive widening
4      $a \sim P_\theta(\tilde{b}, \cdot)$   $\triangleright$  prioritized from network
5      $N(b, a) \leftarrow N_0(b, a)$ 
6      $Q(b, a) \leftarrow Q_0(b, a)$   $\triangleright$  bootstrap initial  $Q$ -value
7      $A(b) \leftarrow A(b) \cup \{a\}$   $\triangleright$  add to visited actions  $A(b)$ 
8   end
9   return  $\arg \max_{a \in A(b)} Q(b, a) + c \left( P_\theta(\tilde{b}, a) \frac{\sqrt{N(b)}}{1+N(b, a)} \right)$   $\triangleright$  PUCT

```

Prioritized action widening. Planning in belief space explicitly handles state uncertainty but may incur computational overhead when performing belief updates, therefore we avoid trying all actions at every belief-state node. Action progressive widening has been used successfully in the context of continuous action spaces (Moerland et al. 2018) and large discrete action spaces (Yee et al. 2016). Mern et al. show that prioritizing actions can improve MCTS performance in large discrete action spaces and Browne et al. found action progressive widening to be effective in cases where favorable actions were tried first.

In BetaZero, we use algorithm 1 to select actions through progressive widening and use information from the learned policy network P_θ to sample new actions. This way, we can first focus the expansion on promising actions, then make the final selection based on PUCT. In the appendix, we perform an ablation to measure the effect of using the policy P_θ to prioritize actions when widening the tree.

Stochastic belief-state transitions. A challenge with partially observable domains is handling non-deterministic belief-state transitions in the tree search. The belief-state transition function T_b consists of several stochastic components and the belief—which is a probability distribution over states—is continuous. To address this, we use progressive widening from Couëtoux et al. (algorithm 2). Other methods for state expansion, like state abstraction refinement from

Algorithm 2: BetaZero belief-state progressive widening.

```

1 function BELIEFSTATEEXPANSION( $b, a$ )
2   if  $|B(b, a)| \leq kN(b, a)^\alpha$   $\triangleright$  belief progressive widening
3      $s \sim b(\cdot)$ 
4      $s' \sim T(\cdot | s, a)$ 
5      $o \sim O(\cdot | a, s')$ 
6      $b' \leftarrow \text{UPDATE}(b, a, o)$ 
7      $B(b, a) \leftarrow B(b, a) \cup \{b'\}$   $\triangleright$  add to visited beliefs
8   else
9      $b' \sim B(b, a)$   $\triangleright$  sample from belief-states in the tree
10  end
11   $r \leftarrow R(b, a)$  or  $r \leftarrow R(b, a, b')$ 
12  return  $b', r$ 

```

Sokota et al., rely on problem-specific distance metrics between states to perform a nearest neighbor search. Progressive widening avoids problem-specific heuristics and uses only the information in the search tree to provide artificially limited belief-state branching which is important as the belief updates can be computationally expensive, thus limiting the MCTS search budget in practice.

Parametric belief representation. Although a particle set belief is not parametrically defined, approximating the belief as summary statistics (e.g., mean and std of the state particles) may capture enough information for value and policy estimation to be used during planning (further analyzed in the appendix). Approximating the particle set parametrically is easy to implement and computationally inexpensive. We show that the approximation works well across various problems and, unsurprisingly, using only the mean state is inadequate. We represent the particle set b parametrically as:

$$\phi(b) \stackrel{\text{def}}{=} [\mu(b), \sigma(b)] \quad (10)$$

BetaZero plans over the full belief b in the tree and only converts to the belief representation $\tilde{b} = \phi(b)$ for network evaluations. Other algorithms (e.g., FORBES from Chen et al.) could instead be used to learn this belief representation.

Bootstrapping initial Q -values. When a new state-action node is added to the tree, initial Q -values can be bootstrapped using the estimate from the value network V_θ :

$$Q_0(b, a) \stackrel{\text{def}}{=} R_b(b, a) + \gamma V_\theta(\phi(b')) \text{ where } b' \sim T_b(\cdot | b, a) \quad (11)$$

Bootstrapping occurs in algorithm 1 (line 6) and incurs an additional belief update through the belief-state transition T_b and may be opted only during online execution. The bootstrapped estimate is more robust (Kumar et al. 2019) and can be useful to initialize online search. Note that MuZero also uses bootstrapping (Schrittwieser et al. 2020).

Algorithm 3 details MCTS for BetaZero and the full BetaZero algorithm is in the appendix (algorithms 4–6).

Algorithm 3: BetaZero MCTS simulation.

```

1 function SIMULATE( $f_\theta, b, d$ )  $\triangleright$  where  $P_\theta, V_\theta \leftarrow f_\theta$ 
2   if  $d = 0$  return 0
3   if  $b \notin \mathcal{T}$ 
4      $\mathcal{T} \leftarrow \mathcal{T} \cup \{b\}$   $\triangleright$  add belief-state node to the tree  $\mathcal{T}$ 
5      $N(b) \leftarrow N_0(b)$ 
6      $\tilde{b} \leftarrow \phi(b)$   $\triangleright$  belief representation eq. (10)
7     return  $V_\theta(\tilde{b})$   $\triangleright$  value-estimate lookup
8   end
9    $N(b) \leftarrow N(b) + 1$ 
10   $a \leftarrow \text{ACTIONSELECTION}(f_\theta, b)$   $\triangleright$  algorithm 1
11   $(b', r) \leftarrow \text{BELIEFSTATEEXPANSION}(b, a)$   $\triangleright$  algorithm 2
12   $q \leftarrow r + \gamma \text{SIMULATE}(f_\theta, b', d - 1)$   $\triangleright$  algorithm 3
13   $N(b, a) \leftarrow N(b, a) + 1$ 
14   $Q(b, a) \leftarrow Q(b, a) + \frac{q - Q(b, a)}{N(b, a)}$   $\triangleright$  backpropagation
15  return  $q$ 

```

4 Experiments

Three benchmark problems were chosen to evaluate the performance of BetaZero. Appendices further describe the POMDPs, network architectures, and experiment design.

	$ \mathcal{S} $	$ \mathcal{A} $	$ \mathcal{O} $
LightDark(5 and 10)	$ \mathbb{R} $	3	$ \mathbb{R} $
RockSample(15, 15)	7,372,800	20	3
RockSample(20, 20)	419,430,400	25	3
Mineral Exploration	$ \mathbb{R}^{32 \times 32} $	38	$ \mathbb{R}_{\geq 0} $

Table 1: POMDP state, action, and observation spaces.

In **LIGHTDARK**(y) (Platt Jr. et al. 2010), the goal of the agent is to execute a `stop` action at the origin while receiving noisy observations of its true location. The noise is minimized in the “light” region $y = 5$. We also benchmark against a harder version with the light region at $y = 10$ from Sunberg and Kochenderfer and restrict the agent to only three actions: move up or down by one, or `stop` (removing actions of moving ten units). The modified problem requires information gathering over longer horizons.

ROCKSAMPLE(n, k) (Smith and Simmons 2004) is a scalable information gathering problem where an agent moves in an $n \times n$ grid to observe k rocks to sample only the “good” rocks. Well-established POMDP benchmarks go up to $n = 15$ and $k = 15$; we also test a harder version with $n = 20$ and $k = 20$ to show the scalability of BetaZero, noting that Cai et al. evaluated this in the multi-agent case.

In the real-world **MINERAL EXPLORATION** problem (Mern and Caers 2023) the agent drills over a 32×32 region to determine if a subsurface ore body should be mined or abandoned. The agent receives a continuous ore quality observation at the drill locations to build its belief. Drilling incurs a penalty and if chosen to mine then the agent is rewarded or penalized based on an economic threshold of the extracted ore mass. The problem is challenging due to reasoning over limited observations with sparse rewards.

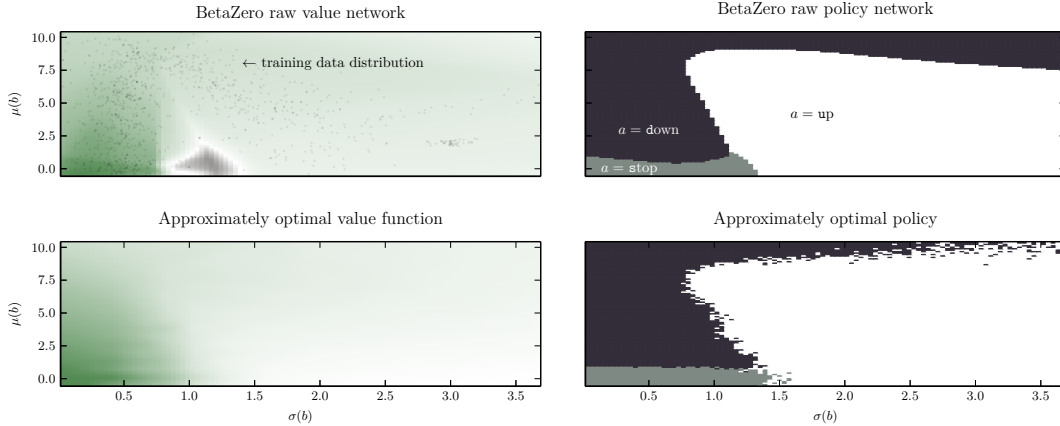


Figure 4: **LIGHTDARK**(10) value and policy plots: BetaZero (top) and value iteration (bottom) over belief mean and std. High uncertainty (horizontal axis) makes the agent localize up near $y = 10$, then moves down and stops at the origin.

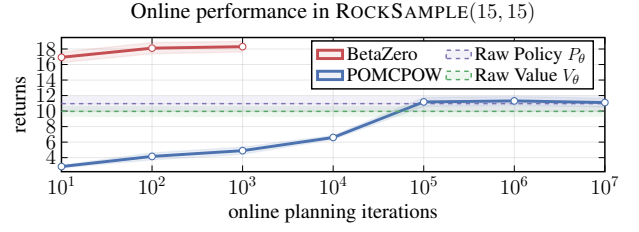


Figure 5: Performance of POMCPOW with heuristics up to 10 million online iterations plateaus, indicating that extending online searches alone misses valuable offline experience.

We baseline BetaZero against several online POMDP algorithms, namely AdaOPS (Wu et al. 2021a), POMCPOW (Sunberg and Kochenderfer 2018), and DESPOT (Ye et al. 2017). In **LightDark**, we solve for an approximately optimal policy using *local approximation value iteration* (LAVI) (Kochenderfer 2015) over a discretized parametric belief space. For a fair comparison, parameters were set to roughly match the total number of simulations experienced.

5 Results and Discussion

Figure 4 compares the raw BetaZero value and policy network with *value iteration* for **LIGHTDARK**(10). Qualitatively, BetaZero learns an accurate optimal policy and value function in areas where training data was collected. Despite unrepresented value function regions (gray), BetaZero remains nearly optimal as those beliefs don’t occur during execution. Out-of-distribution methods could quantify this uncertainty, e.g., an ensemble of networks (Salehi et al. 2022).

Table 2 shows that BetaZero outperforms state-of-the-art algorithms in most cases, with larger improvements when baseline algorithms do not rely on heuristics. Timing results show BetaZero incurs a significant offline penalty, further elaborated in the limitations section below.

In **ROCKSAMPLE**(n, k), BetaZero is comparable to AdaOPS and DESPOT which compute an upper bound us-

	LightDark(5)		LightDark(10)		RockSample(15, 15)		RockSample(20, 20)		Mineral Exploration	
	returns	time [s]	returns	time [s]	returns	time [s]	returns	time [s]	returns	time [s]
BetaZero	4.47 ± 0.28	[2274, 0.014]	16.77 ± 1.28	[2740, 0.331]	20.15 ± 0.71	[5701, 0.477]	13.09 ± 0.55	[7081, 1.109]	10.67 ± 2.25	[22505, 5.126]
Raw Policy P_θ	4.44 ± 0.28	[2274, 0.004]	13.74 ± 1.33	[2740, 0.004]	10.96 ± 0.98	[5701, 0.018]	2.03 ± 0.34	[7081, 0.084]	8.67 ± 2.52	[22505, 0.533]
Raw Value V_θ^*	3.16 ± 0.40	[2274, 0.008]	12.70 ± 1.46	[2740, 0.009]	9.96 ± 0.65	[5701, 0.158]	3.57 ± 0.40	[7081, 0.204]	9.75 ± 2.42	[22505, 1.420]
AdaOPS	3.78 ± 0.27 (3.79 ± 0.07)	[68, 0.089]	5.22 ± 1.77	[81, 0.510]	20.67 ± 0.72 (17.16 ± 0.21)	[7, 2.768]	—	—	3.33 ± 1.95	[5, 0.112]
AdaOPS (fixed bounds)	3.70 ± 0.25	[0, 0.039]	4.98 ± 2.01	[0, 0.573]	13.37 ± 0.71	[0, 1.349]	11.66 ± 0.49	[1, 1.458]	"	"
POMCPOW	3.21 ± 0.38 (3.23 ± 0.11)	[59, 0.189]	0.68 ± 0.41	[70, 1.261]	11.14 ± 0.59 (10.40 ± 0.18)	[0, 0.929]	10.22 ± 0.47	[0, 1.480]	9.43 ± 2.19	[0, 6.728]
POMCPOW (no heuristics)	1.96 ± 0.58	[0, 0.099]	-5.90 ± 5.78	[0, 0.742]	10.17 ± 0.61	[0, 1.485]	4.03 ± 0.44	[0, 5.173]	5.38 ± 2.15	[0, 5.915]
DESPOT	2.37 ± 0.37 (2.50 ± 0.10)	[0, 0.008]	0.43 ± 0.36	[0, 0.046]	18.44 ± 0.69 (15.67 ± 0.20)	[7, 3.822]	—	—	5.29 ± 2.17	[5, 0.283]
DESPOT (fixed bounds)	2.70 ± 0.50	[0, 0.008]	0.49 ± 0.30	[0, 0.025]	4.29 ± 0.45	[0, 5.091]	0.00	[0, 5.179]	"	"
Approx. Optimal	4.09 ± 0.33	[267, 0.037]	14.16 ± 1.39	[260, 0.025]	—	—	—	—	11.90 ± 0.18	N/A

* One-step look-ahead over all actions using only the value network with 5 observations per action. All results report the mean return \pm standard error over 100 seeds. Entries with "—" indicate they failed to run on that domain, entries with " are the same as the ones above, and entries in (parentheses) are from the literature.

Table 2: Results comparing *BetaZero* to various state-of-the-art POMDP solvers (reporting returns and [offline, online] timing).

ing QMDP. QMDP computes the optimal utility of the fully observable MDP over all $k - 1$ rock combinations, which scales exponentially in n . For larger state spaces, like ROCK-SAMPLE(20, 20), the QMDP solution is intractable. Thus, fixed bounds are used assuming an optimistic V_{\max} (Wu et al. 2021b) while BetaZero scales well to these higher dimensional problems. Indicated in table 2, the raw networks alone perform well but outperform when combined with online planning, enabling reasoning with current information.

If online algorithms ran for a large number of iterations, one might expect to see convergence to the optimal policy. In practice, this may be an intractable number as fig. 5 shows POMCPOW has not reached the required number of iterations for RockSample. The advantage of BetaZero is that it can generalize from a more diverse set of experiences.

The inability of existing online algorithms to plan over long horizons is also evident in the mineral exploration

POMDP (fig. 6). POMCPOW ran for one million online iterations without a value estimator heuristic and BetaZero ran online for 100 iterations (using about 850,000 offline simulations). In the figure, the probability of selecting a drilling location is shown as vertical bars for each action, overlaid on the initial belief uncertainty (i.e., the standard deviation of the belief in subsurface ore quality). BetaZero learned to take actions in areas of the belief space with high uncertainty (which matches the domain-specific heuristic developed for the mineral exploration problem from Mern and Caers), while POMCPOW fails to distinguish between the actions and resembles a uniform policy.

Limitations. It is standard for POMDP planning algorithms to assume known models but this may limit the applicability to certain problems where reinforcement learning may be better suited. We chose a simplified belief representation to allow for further research innovations in using other parametric and non-parametric representations. Other limitations include compute resource requirements for training neural networks and parallelizing MCTS simulations. We designed BetaZero to use a single GPU for training and to scale based on available CPUs. Certain POMDPs may not require the training burden, especially when known heuristics perform well. BetaZero is useful for long-horizon, high-dimensional continuous POMDPs but may be unnecessary when offline training is computationally limited. BetaZero is designed for problems where the simulation cost is the dominating factor compared to offline training time.

6 Conclusions

We propose the *BetaZero* belief-state planning algorithm for POMDPs; designed to learn from offline experience to inform online decisions. Planning in belief space explicitly handles state uncertainty and learning offline approximations to replace heuristics enables effective online planning in long-horizon POMDPs. BetaZero can also scale to larger problems where certain heuristics break down. Results suggest that BetaZero can solve large-scale POMDPs and learns to plan in belief space using zero heuristics.

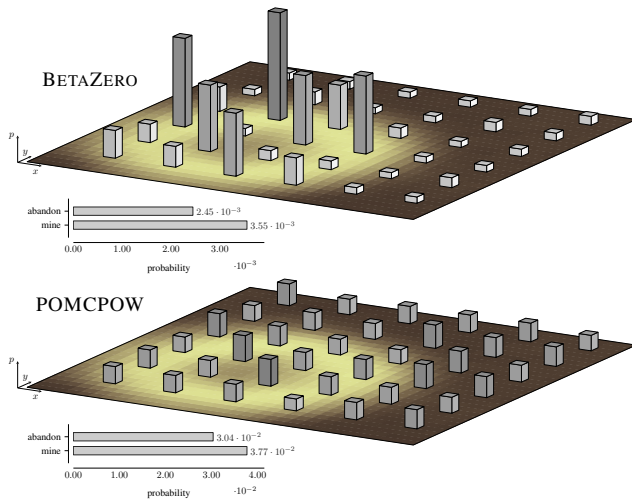


Figure 6: Mineral exploration policies: BetaZero prioritizes uncertainty, matching heuristics from Mern and Caers.

References

- Antonoglou, I.; Schrittwieser, J.; Ozair, S.; Hubert, T. K.; and Silver, D. 2021. Planning in Stochastic Environments with a Learned Model. In *International Conference on Learning Representations*.
- Browne, C. B.; Powley, E.; Whitehouse, D.; Lucas, S. M.; Cowling, P. I.; Rohlfshagen, P.; Tavener, S.; Perez, D.; Samothrakis, S.; and Colton, S. 2012. A Survey of Monte Carlo Tree Search Methods. *IEEE Transactions on Computational Intelligence and AI in Games*, 4(1): 1–43.
- Cai, P.; Luo, Y.; Hsu, D.; and Lee, W. S. 2021. HyP-DESPOT: A hybrid parallel algorithm for online planning under uncertainty. *The International Journal of Robotics Research*, 40(2-3): 558–573.
- Chai, T.; and Draxler, R. R. 2014. Root mean square error (RMSE) or mean absolute error (MAE)?—Arguments against avoiding RMSE in the literature. *Geoscientific Model Development*, 7(3): 1247–1250.
- Chen, X.; Mu, Y. M.; Luo, P.; Li, S.; and Chen, J. 2022. Flow-Based Recurrent Belief State Learning for POMDPs. In *International Conference on Machine Learning*.
- Corso, A.; Wang, Y.; Zechner, M.; Caers, J.; and Kochenderfer, M. J. 2022. A POMDP Model for Safe Geological Carbon Sequestration. *NeurIPS Workshop on Tackling Climate Change with Machine Learning*.
- Couëtoux, A.; Hoock, J.-B.; Sokolovska, N.; Teytaud, O.; and Bonnard, N. 2011. Continuous Upper Confidence Trees. In *Learning and Intelligent Optimization*. Springer.
- Coulom, R. 2007. Efficient Selectivity and Backup Operators in Monte-Carlo Tree Search. In *Computers and Games*.
- Csilléry, K.; Blum, M. G.; Gaggiotti, O. E.; and François, O. 2010. Approximate Bayesian Computation (ABC) in practice. *Trends in Ecology & Evolution*, 25(7): 410–418.
- Czech, J.; Korus, P.; and Kersting, K. 2021. Improving AlphaZero Using Monte-Carlo Graph Search. In *International Conference on Automated Planning and Scheduling (ICAPS)*, volume 31, 103–111.
- Danihelka, I.; Guez, A.; Schrittwieser, J.; and Silver, D. 2022. Policy Improvement by Planning with Gumbel. In *International Conference on Learning Representations*.
- Egorov, M.; Sunberg, Z. N.; Balaban, E.; Wheeler, T. A.; Gupta, J. K.; and Kochenderfer, M. J. 2017. POMDPs.jl: A Framework for Sequential Decision Making under Uncertainty. *Journal of Machine Learning Research*, 18(26): 1–5.
- Fischer, J.; and Tas, Ö. S. 2020. Information Particle Filter Tree: An Online Algorithm for POMDPs with Belief-Based Rewards on Continuous Domains. In *International Conference on Machine Learning*, 3177–3187. PMLR.
- Gordon, N. J.; Salmond, D. J.; and Smith, A. F. 1993. Novel approach to nonlinear/non-Gaussian Bayesian state estimation. In *IEEE for Radar and Signal Processing*, volume 140.
- Hinton, G.; Srivastava, N.; and Swersky, K. 2014. Neural Networks for Machine Learning. Lecture 6a: Overview of Mini-Batch Gradient Descent. http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf.
- Hodson, T. O. 2022. Root-mean-square error (RMSE) or mean absolute error (MAE): when to use them or not. *Geoscientific Model Development*, 15(14): 5481–5487.
- Hoel, C.-J.; Driggs-Campbell, K.; Wolff, K.; Laine, L.; and Kochenderfer, M. J. 2019. Combining Planning and Deep Reinforcement Learning in Tactical Decision Making for Autonomous Driving. *IEEE Transactions on Intelligent Vehicles*, 5(2): 294–305.
- Hubert, T.; Schrittwieser, J.; Antonoglou, I.; Barekatin, M.; Schmitt, S.; and Silver, D. 2021. Learning and Planning in Complex Action Spaces. In *International Conference on Machine Learning*, 4476–4486. PMLR.
- Igl, M.; Zintgraf, L.; Le, T. A.; Wood, F.; and Whiteson, S. 2018. Deep Variational Reinforcement Learning for POMDPs. In *International Conference on Machine Learning*.
- Kimura, T.; Sakamoto, K.; and Sogabe, T. 2020. Development of AlphaZero-Based Reinforcement Learning Algorithm for Solving Partially Observable Markov Decision Process (POMDP) Problem. *Bulletin of Networking, Computing, Systems, and Software*, 9(1): 69–73.
- Kingma, D. P.; and Ba, J. 2014. Adam: A Method for Stochastic Optimization. *arXiv preprint arXiv:1412.6980*.
- Kochenderfer, M. J. 2015. *Decision Making Under Uncertainty: Theory and Application*. MIT Press.
- Kochenderfer, M. J.; Holland, J. E.; and Chrysanthacopoulos, J. P. 2012. Next-Generation Airborne Collision Avoidance System. *Lincoln Laboratory Journal*, 19(1).
- Kochenderfer, M. J.; Wheeler, T. A.; and Wray, K. H. 2022. *Algorithms for Decision Making*. MIT Press.
- Kocsis, L.; and Szepesvári, C. 2006. Bandit Based Monte-Carlo Planning. In *European Conference on Machine Learning*, 282–293. Springer.
- Kumar, A.; Fu, J.; Soh, M.; Tucker, G.; and Levine, S. 2019. Stabilizing Off-Policy Q-Learning via Bootstrapping Error Reduction. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 32.
- Lauri, M.; Hsu, D.; and Pajarinen, J. 2022. Partially Observable Markov Decision Processes in Robotics: A Survey. *IEEE Transactions on Robotics*.
- LeCun, Y.; Bottou, L.; Orr, G. B.; and Müller, K.-R. 2002. Efficient BackProp. In *Neural Networks: Tricks of the Trade*, 9–50. Springer.
- Lim, M. H.; Becker, T. J.; Kochenderfer, M. J.; Tomlin, C. J.; and Sunberg, Z. N. 2023. Optimality Guarantees for Particle Belief Approximation of POMDPs. *Journal of Artificial Intelligence Research*, 77: 1591–1636.
- Mern, J.; and Caers, J. 2023. The Intelligent Prospector v1.0: Geoscientific Model Development and Prediction by Sequential Data Acquisition Planning with Application to Mineral Exploration. *Geoscientific Model Development*, 16(1): 289–313.
- Mern, J.; Yildiz, A.; Bush, L.; Mukerji, T.; and Kochenderfer, M. J. 2021. Improved POMDP Tree Search Planning with Prioritized Action Branching. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(13): 11888–11894.

- Moerland, T. M.; Broekens, J.; Plaat, A.; and Jonker, C. M. 2018. AOC: Alpha Zero in Continuous Action Space. *arXiv preprint arXiv:1805.09613*.
- Moss, R. J.; Lee, R.; Visser, N.; Hochwarth, J.; Lopez, J. G.; and Kochenderfer, M. J. 2020. Adaptive Stress Testing of Trajectory Predictions in Flight Management Systems. *Digital Avionics Systems Conference (DASC)*.
- Munos, R. 2014. From Bandits to Monte-Carlo Tree Search: The Optimistic Principle Applied to Optimization and Planning. *Foundations and Trends® in Machine Learning*, 7(1).
- Ozair, S.; Li, Y.; Razavi, A.; Antonoglou, I.; Van Den Oord, A.; and Vinyals, O. 2021. Vector Quantized Models for Planning. In *International Conference on Machine Learning*, 8302–8313. PMLR.
- Platt Jr., R.; Tedrake, R.; Kaelbling, L.; and Lozano-Perez, T. 2010. Belief Space Planning Assuming Maximum Likelihood Observations. *Robotics: Science and Systems VI*.
- Rosin, C. D. 2011. Multi-Armed Bandits with Episode Context. *Annals of Mathematics and Artificial Intelligence*, 61(3): 203–230.
- Salehi, M.; Mirzaei, H.; Hendrycks, D.; Li, Y.; Rohban, M. H.; and Sabokrou, M. 2022. A Unified Survey on Anomaly, Novelty, Open-Set, and Out-of-Distribution Detection: Solutions and Future Challenges. *Transactions of Machine Learning Research*.
- Schadd, F. C. 2009. Monte-Carlo Search Techniques in the Modern Board Game Thurn and Taxis. *Maastricht University: Maastricht, The Netherlands*.
- Schrittwieser, J. 2020. MuZero Intuition. <https://www.furidamu.org/blog/2020/12/22/muzero-intuition/>.
- Schrittwieser, J.; Antonoglou, I.; Hubert, T.; Simonyan, K.; Sifre, L.; Schmitt, S.; Guez, A.; Lockhart, E.; Hassabis, D.; Graepel, T.; Lillicrap, T.; and Silver, D. 2020. Mastering Atari, Go, Chess and Shogi by Planning with a Learned Model. *Nature*, 588(7839): 604–609.
- Shani, G.; Pineau, J.; and Kaplow, R. 2013. A survey of point-based POMDP solvers. *Autonomous Agents and Multi-Agent Systems*, 27: 1–51.
- Silver, D.; Huang, A.; Maddison, C. J.; Guez, A.; Sifre, L.; van den Driessche, G.; Schrittwieser, J.; Antonoglou, I.; Panneershelvam, V.; Lanctot, M.; Dieleman, S.; Grewe, D.; Nham, J.; Kalchbrenner, N.; Sutskever, I.; Lillicrap, T.; Leach, M.; Kavukcuoglu, K.; Graepel, T.; and Hassabis, D. 2016. Mastering the game of Go with deep neural networks and tree search. *Nature*, 529(7587): 484–489.
- Silver, D.; Hubert, T.; Schrittwieser, J.; Antonoglou, I.; Lai, M.; Guez, A.; Lanctot, M.; Sifre, L.; Kumaran, D.; Graepel, T.; Lillicrap, T.; Simonyan, K.; and Hassabis, D. 2018. A general reinforcement learning algorithm that masters chess, shogi, and Go through self-play. *Science*, 362(6419).
- Silver, D.; Schrittwieser, J.; Simonyan, K.; Antonoglou, I.; Huang, A.; Guez, A.; Hubert, T.; Baker, L.; Lai, M.; Bolton, A.; Chen, Y.; Lillicrap, T.; Hui, F.; Sifre, L.; van den Driessche, G.; Graepel, T.; and Hassabis, D. 2017. Mastering the game of Go without human knowledge. *Nature*, 550(7676).
- Silver, D.; and Veness, J. 2010. Monte-Carlo Planning in Large POMDPs. *Advances in Neural Information Processing Systems (NeurIPS)*, 23.
- Smith, T.; and Simmons, R. 2004. Heuristic Search Value Iteration for POMDPs. In *Conference on Uncertainty in Artificial Intelligence (UAI)*, 520–527.
- Sokota, S.; Ho, C. Y.; Ahmad, Z.; and Kolter, J. Z. 2021. Monte Carlo Tree Search with Iteratively Refining State Abstractions. *Advances in Neural Information Processing Systems (NeurIPS)*, 34: 18698–18709.
- Sunberg, Z. N.; and Kochenderfer, M. J. 2018. Online Algorithms for POMDPs with Continuous State, Action, and Observation Spaces. In *Proceedings of the International Conference on Automated Planning and Scheduling*, volume 28.
- Sutton, R. S.; and Barto, A. G. 2018. *Reinforcement Learning: An Introduction*. MIT Press.
- Thrun, S.; Burgard, W.; and Fox, D. 2005. *Probabilistic Robotics*. MIT Press.
- Vinyals, O.; Babuschkin, I.; Czarnecki, W. M.; Mathieu, M.; Dudzik, A.; Chung, J.; Choi, D. H.; Powell, R.; Ewalds, T.; Georgiev, P.; et al. 2019. Grandmaster Level in StarCraft II using Multi-Agent Reinforcement Learning. *Nature*, 575(7782): 350–354.
- Wang, Y.; Zechner, M.; Mern, J. M.; Kochenderfer, M. J.; and Caers, J. K. 2022. A sequential decision-making framework with uncertainty quantification for groundwater management. *Advances in Water Resources*, 166: 104266.
- Wang, Y.; Zechner, M.; Wen, G.; Corso, A. L.; Mern, J. M.; Kochenderfer, M. J.; and Caers, J. K. 2023. Optimizing Carbon Storage Operations for Long-Term Safety. *arXiv preprint arXiv:2304.09352*.
- Wray, K. H.; Lange, B.; Jamgochian, A.; Witwicki, S. J.; Kobashi, A.; Hagaribommanahalli, S.; and Ilstrup, D. 2021. POMDPs for Safe Visibility Reasoning in Autonomous Vehicles. In *IEEE International Conference on Intelligence and Safety for Robotics (ISR)*, 191–195. IEEE.
- Wu, C.; Yang, G.; Zhang, Z.; Yu, Y.; Li, D.; Liu, W.; and Hao, J. 2021a. Adaptive Online Packing-guided Search for POMDPs. *Advances in Neural Information Processing Systems (NeurIPS)*, 34: 28419–28430.
- Wu, C.; Yang, G.; Zhang, Z.; Yu, Y.; Li, D.; Liu, W.; and Hao, J. 2021b. OpenReview Response: Adaptive Online Packing-guided Search for POMDPs. https://openreview.net/forum?id=0zvTBoQb5PA¬eId=bJTC_0sZPzr.
- Ye, N.; Somani, A.; Hsu, D.; and Lee, W. S. 2017. DESPOT: Online POMDP Planning with Regularization. *Journal of Artificial Intelligence Research*, 58: 231–266.
- Yee, T.; Lisý, V.; Bowling, M. H.; and Kambhampati, S. 2016. Monte Carlo Tree Search in Continuous Action Spaces with Execution Uncertainty. In *International Joint Conference on Artificial Intelligence (IJCAI)*, 690–697.
- Zaheer, M.; Kottur, S.; Ravanbakhsh, S.; Poczos, B.; Salakhutdinov, R. R.; and Smola, A. J. 2017. Deep Sets. *Advances in Neural Information Processing Systems (NeurIPS)*.

Appendix

This section contains material detailing the POMDP environments and experiments, the ablation studies, additional analysis of bootstrapping and DPW, the network architectures, the hyperparameters and tuning, computational resources, information regarding open-source code for reproducibility, and the full BetaZero algorithm pseudocode.

POMDP Environments

This section describes the benchmark POMDPs in detail.

Light dark. The $\text{LIGHTDARK}(y)$ POMDP is a one-dimensional localization problem (Platt Jr. et al. 2010). The objective is for the agent to execute the `stop` action at the goal, which is at ± 1 of the origin. The agent is awarded 100 for stopping at the goal and -100 for stopping anywhere else; using a discount of $\gamma = 0.9$. The agent receives noisy observations of their position, where the noise is minimized in the “light” region defined by y . In the $\text{LIGHTDARK}(5)$ problem used by Wu et al., the noise is a zero-mean Gaussian with standard deviation of $|y - 5|/\sqrt{2} + 10^{-2}$. For the $\text{LIGHTDARK}(10)$ problem used by Sunberg and Kochenderfer, the noise is a zero-mean Gaussian with standard deviation of $|y - 10| + 10^{-4}$. In both problems, we use a restricted action space of $\mathcal{A} = [-1, 0, 1]$ where 0 is the `stop` action. The expected behavior of the optimal policy is first to localize in the light region, then travel down to the goal. The BetaZero policy exhibits this behavior which can be seen in fig. 7 (where circles indicate the final location).

The approximately optimal solution to the light dark problems used *local approximation value iteration* (LAVI) (Kochenderfer 2015) over the discretized belief-state space (i.e., mean and std). The belief mean was discretized between the range $[-3, 12]$ and the belief std was discretized between the range $[0, 5]$; each of length 25. The LAVI solver used 100 generative samples per belief state and ran for 25 value iterations with a Bellman residual of 1×10^{-3} .

Rock sample. In the $\text{ROCKSAMPLE}(n, k)$ POMDP introduced by Smith and Simmons, an agent has full observability of its position on an $n \times n$ grid but has to sense the k rocks to determine if they are “good” or “bad”. The agent knows *a priori* the true locations of the rocks (i.e., the rock locations \mathbf{x}_{rock} are a part of the problem, not the state). The observation noise is a function of the distance to the rock:

$$\frac{1}{2} \left(1 + \exp \left(-\frac{\|\mathbf{x}_{\text{rock}} - \mathbf{x}_{\text{agent}}\|_2 \log(2)}{c} \right) \right) \quad (12)$$

where $c = 20$ is the sensor efficiency. The agent can move in the four cardinal directions, sense the k rocks, or take the action to `sample` a rock when it is located under the agent. The agent receives a reward of 10 for sampling a “good” rock and a penalty of -10 for sampling a “bad” rock. The terminal state is the exit at the right edge of the map, where the agent gets a reward of 10 for exiting.

Mineral exploration. The $\text{MINERAL EXPLORATION}$ POMDP introduced by Mern and Caers is an information gather problem with the goal of deciding whether a subsurface ore body is economical to mine or should be abandoned.

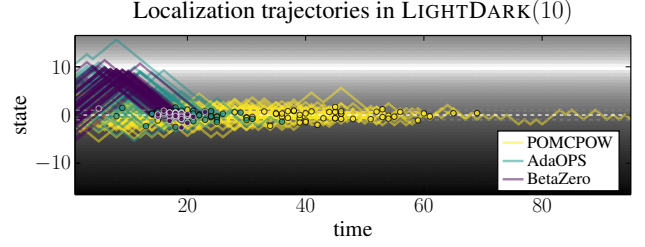


Figure 7: $\text{LIGHTDARK}(10)$ trajectories from 50 episodes. BetaZero (dark blue) learned to first localize in the light region at $y = 10$ before heading to the goal (origin).

The agent can drill every fifth cell of a 32×32 plot of land to determine the ore quality at that location. Therefore, the action space consists of the 36 drill locations and the final decisions to either mine or abandon. The agent receives a small cost for each drill action, a reward proportional to the extracted ore if chosen to mine (which is negative if uneconomical), and a reward of zero if chosen to abandon:

$$R(s, a) = \begin{cases} -c_{\text{drill}} & \text{if } a = \text{drill} \\ \sum \mathbb{1}(s_{\text{ore}} \geq h_{\text{massive}}) - c_{\text{extract}} & \text{if } a = \text{mine} \\ 0 & \text{otherwise} \end{cases}$$

where $c_{\text{drill}} = 0.1$, $h_{\text{massive}} = 0.7$, and $c_{\text{extract}} = 71$. The term $\sum \mathbb{1}(s_{\text{ore}} \geq h_{\text{massive}})$ indicates the cells that have an ore quality value above some massive ore threshold h_{massive} (which are deemed valuable). Figure 8 and fig. 9 show an example of four steps of the mineral exploration POMDP.

Experiment details

Experiment parameters for each problem can be seen in tables 4–6 under the “online” column. For the baseline algorithms, the heuristics follow Wu et al.. Problems that failed to run due to memory limits followed suggestions from Wu et al. to first use the MDP solution and then use a fixed upper bound of $r_{\text{correct}} = 100$ for the light dark problems and the following for the rock sample problems:

$$V_{\text{max}} = r_{\text{exit}} + \sum_{t=1+n-k}^{2k-n} \gamma^{t-1} r_{\text{good}} \quad (13)$$

where $r_{\text{good}} = r_{\text{exit}} = 10$ and the sum computes an optimistic value assuming the rocks are directly lined between the agent and the goal and assuming $n \geq k$ for simplicity.

For problems not studied by Wu et al., we use the same heuristics as their easier counterpart (i.e., $\text{LIGHTDARK}(10)$ uses $\text{LIGHTDARK}(5)$ heuristics and $\text{ROCKSAMPLE}(20, 20)$ uses $\text{ROCKSAMPLE}(15, 15)$ heuristics). For mineral exploration, the baselines used the following heuristics. POMCPOW used a value estimator of $\max(0, R(s, a = \text{mine}))$ and when using “no heuristic” used a random rollout policy to estimate the value. Both AdaOPS and DESPOT used a lower bound computed as the returns if fully drilled all locations, then made the decision to abandon:

$$V_{\text{min}} = - \sum_{t=1}^{T-1} \gamma^{t-1} c_{\text{drill}} \quad (14)$$

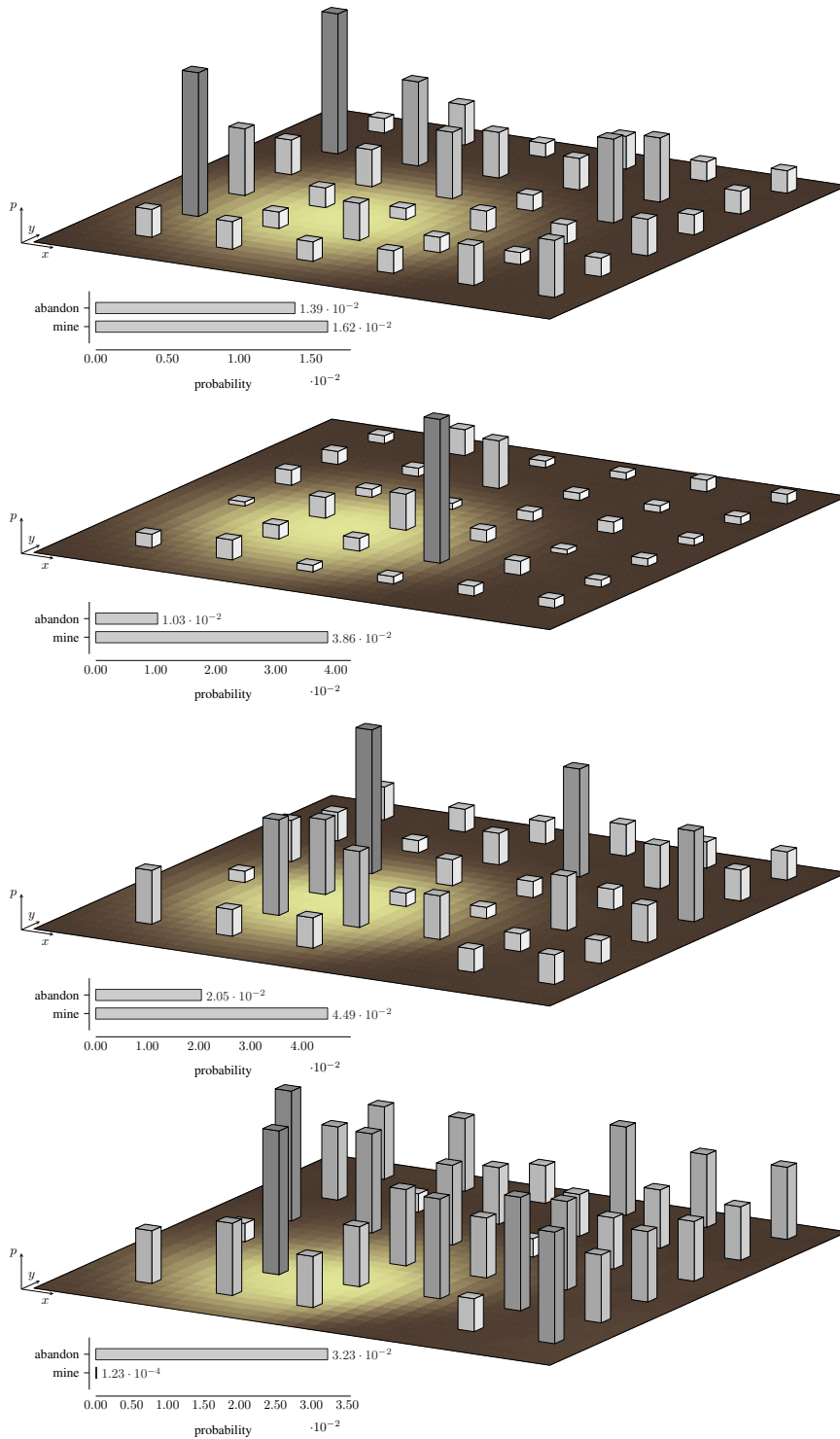


Figure 8: The BetaZero policy shown over belief mean for four steps. BetaZero first prioritizes the edges of the belief mean, corresponding to the belief uncertainty (right-most plots), then explores the outer regions of the subsurface; ultimately gathering information from actions with high mean and std, matching heuristics. At the initial step, abandoning and mining have near-equal probability (bottom left graphs) but by the fourth action, abandoning is much more likely.

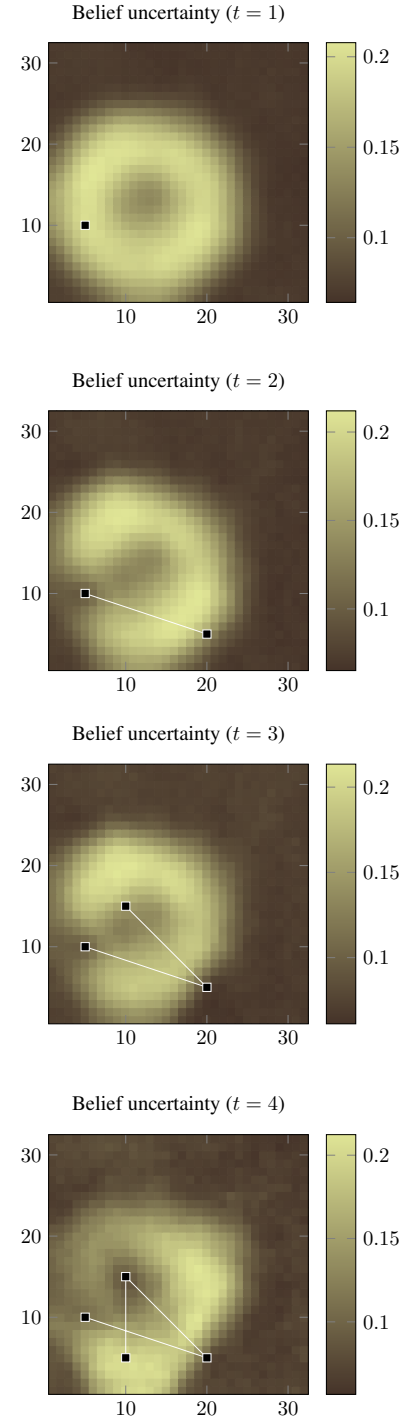


Figure 9: The selected drill actions over belief uncertainty, showing that uncertainty collapses after drilling.

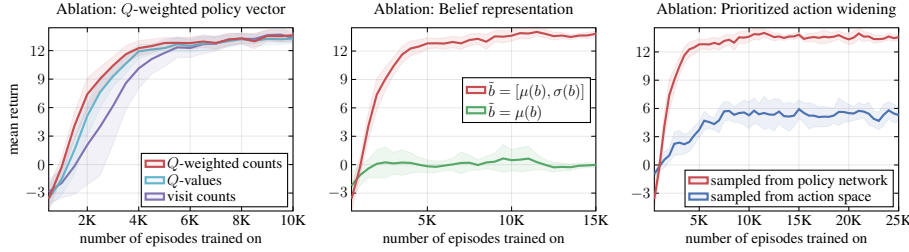


Figure 10: Ablation study in LIGHTDARK(10). (Left) Learning is faster when policy network is trained using Q -weighted visit counts. (Middle) Incorporating belief uncertainty is crucial for learning. (Right) Action widening from the policy network shows significant improvement. The same red curves are shown and one std is shaded from three seeds with exponential smoothing with weight 0.6.

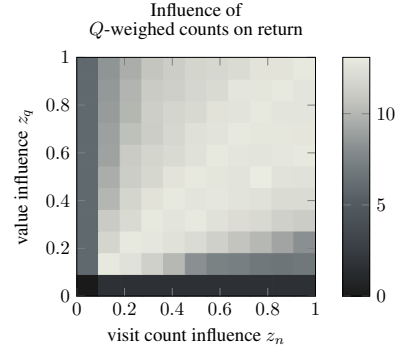


Figure 11: Ablation study in ROCK-SAMPLE(20,20). Combining value and count information leads to the highest return. The diagonal is identical due to the arg max of eq. (8).

The upper bound comes from an oracle taking the correct final action without drilling, computed over 10,000 states:

$$V_{\max} = \mathbb{E}_{s_0 \in \mathcal{S}} \left[\max \left(0, R(s_0, a = \text{mine}) \right) \right] \quad (15)$$

$$\approx \frac{1}{n} \sum_{i=1}^n \max \left(0, R(s_0^{(i)}, a = \text{mine}) \right)$$

Particle filtering. Both BetaZero and the baseline algorithms update their belief with a bootstrap particle filter using a low-variance resampler (Gordon, Salmond, and Smith 1993), with $n_{\text{particles}} \in [500, 1000, 1000]$ for the light dark, rock sample, and mineral exploration problems, respectively. The particle filter follows an update procedure of first reweighting then resampling. In mineral exploration, the observations are noiseless which could quickly result in particle depletion. Therefore, approximate Bayesian computation (ABC) is used to reweight each particle using a Gaussian distribution centered at the observation with a standard deviation of $\sigma_{\text{abc}} = 0.1$ (Csilléry et al. 2010).

The belief representation takes the mean and standard deviation across the $n_{\text{particles}}$. In the light dark problems, this is computed across the 500 sampled y -state values to make up the belief. The initial y -value state distribution—which makes up the initial belief—follows a Gaussian distribution and thus the parametric representation is a good approximation of the belief. For the rock sample problem, the belief is represented as the mean and standard deviation of the good rocks from the 1000 sampled states (appending the true position as it is deterministic). The rock qualities are sampled uniformly in $\{0, 1\}$ indicating if they are “good”, which makes the problem non-Gaussian but the parametric belief approximation can model a uniform distribution by placing the mean at the center of the uniform range and stretching the variance to match the uniform. Lastly, the mineral exploration problem flattens the 1000 subsurface 32×32 maps that each have associated ore quality per-pixel between $[0, 1]$ into two images: a mean and standard deviation image of the ore quality that is stacked and used as input to a CNN. The

initial state distribution for the massive ore quantity closely follows a Gaussian, making the parametric belief approximation well suited. For problems where Gaussian approximations do not capture the belief, the parameters of other distributions could be used as a belief representation or the particles themselves could be used as a belief representation—first passing through an order-invariant layer (Igl et al. 2018). Scaling to larger observation spaces will not be an issue as BetaZero plans over belief states instead of observations.

Additional Analysis

This section briefly describes the ablation studies and additional analyses omitted from the main body of the paper.

Ablation studies

This section tests the effect of each contribution. The influence of value and visit count information when selecting an action is shown in fig. 11. Each cell is the mean return for the ROCKSAMPLE(20,20) problem over 100 online trials; selecting root-node actions via the arg max of eq. (8) given z_q and z_n . The cell at (0,0) corresponds to a uniform policy and thus samples actions instead. Using only the visit counts (bottom cells) or only the values (left cells) to make decisions is worse than using a combination of the two. The effect of the Q -weighting is also shown in the leftmost fig. 10, which suggests that it helps learn faster in LIGHTDARK(10).

Unsurprisingly, using the state uncertainty encoded in the belief is crucial for learning as indicated in the middle of fig. 10. Future work could directly input the particle set into the network, first passing through an order invariant layer (Zaheer et al. 2017), to offload the belief approximation to the network itself. Finally, the rightmost plot in fig. 10 suggests that when branching on actions using progressive widening, it is important to first prioritize the actions suggested by the policy network. Offline learning fails if instead we sample uniformly from the action space.

	LightDark(5)		LightDark(10)		RockSample(15, 15)		RockSample(20, 20)		Mineral Exploration	
	returns	time [s]	returns	time [s]	returns	time [s]	returns	time [s]	returns	time [s]
BetaZero	4.22 ± 0.31	0.014	14.45 ± 1.15	0.34	20.15 ± 0.71	0.48	13.09 ± 0.55	1.11	10.32 ± 2.38	6.27
BetaZero (no bootstrap)	4.47 ± 0.28	0.014	16.77 ± 1.28	0.33	19.50 ± 0.71	0.42	11.00 ± 0.54	0.57	10.67 ± 2.25	4.46

Reporting mean ± standard error over 100 seeds (i.e., episodes); timing is average per episode.

Table 3: Effect of Q -value bootstrapping in online *BetaZero* performance (reporting returns and online timing).

Bootstrapping analysis

When adding a state-action pair (b, a) to the MCTS tree, initializing the Q -values via bootstrapping with the value network may improve performance when using a small MCTS budget. Table 3 shows the results of an analysis comparing BetaZero with bootstrapping $Q_0(b, a) = R_b(b, a) + \gamma V_\theta(\tilde{b}')$ where $\tilde{b}' = \phi(b')$ and without bootstrapping $Q_0(b, a) = 0$. Each domain used the online parameters described in tables 4–6. Results indicate that bootstrapping was only helpful in the rock sample problems and incurs additional compute time due to the belief update done in $b' \sim T_b(b, a)$. Note that bootstrapping was not used during offline training. In problems with high stochasticity in the belief-state transitions, bootstrapping may be noisy during the initial search due to the transition T_b sampling a single belief. Further analysis could investigate the use of multiple belief transitions to better estimate the value; at the expense of additional computation. The value estimate of b could instead be used as the bootstrap but we would expect similar results to the one-step bootstrap as many problems we study have sparse rewards.

Limitations of double progressive widening

Double progressive widening (DPW) is a straightforward approach to handle large or continuous state and action spaces in Monte Carlo tree search. It is easy to implement and only requires information available in the tree search, i.e., number of children nodes and number of node visits. It is known that MCTS performance can be sensitive to DPW hyperparameter tuning and Sokota et al. show that DPW ignores information about the relation between states that could provide more intelligent branching. Sokota et al. introduce *state abstraction refinement* that uses a distance metric between states to determine if a similar state should be added

to the tree; requiring a state transition every time a state-action node is visited. For our work, we want to reduce the number of expensive belief-state transitions in the tree and avoid the use of problem-specific heuristics required when defining distance metrics. Using DPW in BetaZero was motivated by simplicity and allows future work to innovate on the components of belief-state and action branching.

To analyze the sensitivity of DPW, fig. 12–13 show a sweep over the α and k parameters for DPW in LIGHTDARK(10). Figure 12 shows that the light dark problem is sensitive to belief-state widening and fig. 13 indicates that this problem may not require widening on all actions—noting that when $k = 0$, the only action expanded on is the one prioritized from the policy head $a \sim P_\theta(\tilde{b}, \cdot)$. The light dark problems have a small action space of $|\mathcal{A}| = 3$, therefore this prioritization leads to good performance when only a single action is evaluated (left cells in fig. 13 when $k = 0$).

In ROCKSAMPLE(20, 20), fig. 14–15 indicates that this problem benefits from a higher widening factor (top right of the figures) as the action space $|\mathcal{A}| = 25$ is larger and the belief-state transitions operate over a much larger state space. DPW uses a single branching factor throughout the tree search and research into methods that adapt the branching based on learned information would be a valuable direction to explore.

Lim et al. introduce a class of POMDP planning algorithms that use a fixed number of samples to branch on instead of progressive widening. The bottom row of figures 12–15 (where $\alpha = 0$) can be interpreted as a fixed branching factor compared to progressive widening in the other cells. The analysis in the figures shows that there are cases where BetaZero has better performance when using progressive widening (show in the lighter colors).

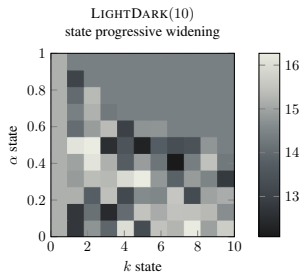


Figure 12: Sensitivity analysis of belief-state progressive widening in LIGHTDARK(10).

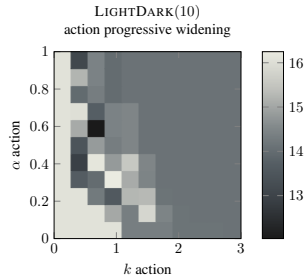


Figure 13: Sensitivity analysis of action progressive widening in LIGHTDARK(10).

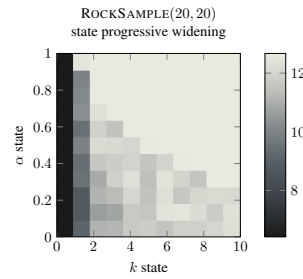


Figure 14: Sensitivity analysis of belief-state progressive widening in ROCKSAMPLE(20, 20).

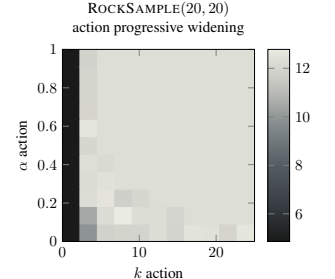


Figure 15: Sensitivity analysis of action progressive widening in ROCKSAMPLE(20, 20).

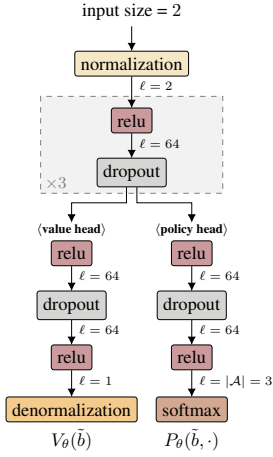


Figure 16: Light dark neural network architecture.

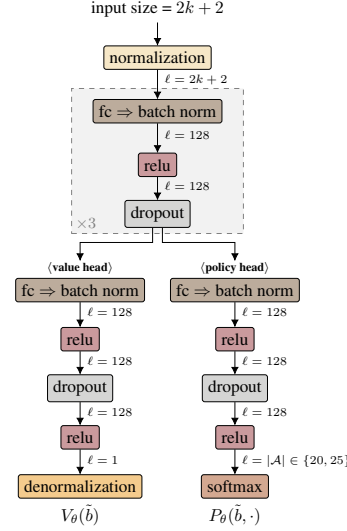


Figure 17: Rock sample neural network architecture.

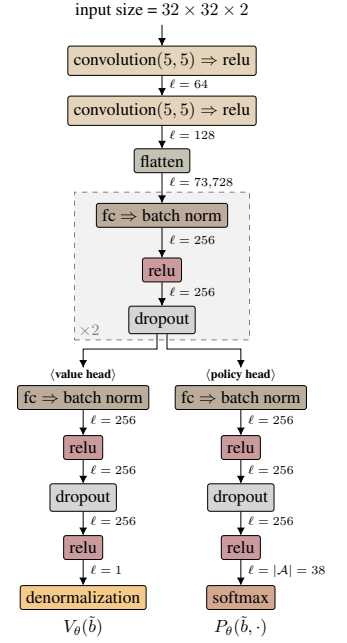


Figure 18: Mineral exploration CNN architecture.

Network Architectures

Figures 16–18 specify the neural network architectures for the three problem domains. The networks were designed to be simple so that future work could focus on incorporating more complicated architectures such as residual networks. Mineral exploration does not normalize the inputs and is the only problem where the input is treated as an image, thus we use a convolutional neural network (CNN). Training occurs on normalized returns and an output denormalization layer is added to the value head to ensure proper magnitude of the predicted values.

Return scaling for output normalization

For general POMDPs, the return can be an unbounded real-value and not conveniently in $[0, 1]$ or $[-1, 1]$; as is often the case with two player games. Schrittwieser et al. use a categorical representation of the value split into a discrete support to make learning more robust (Schrittwieser 2020). We instead simply normalize the target before training as

$$\bar{g}_t = (g_t - \mathbb{E}[G_{\text{train}}]) / \sqrt{\text{Var}[G_{\text{train}}]} \quad (16)$$

where G_{train} is the set of returns used during training; keeping running statistics of all training data. Intuitively, this ensures that the target values have zero mean and unit variance which is known to stabilize training (LeCun et al. 2002). After training, a denormalization layer is added to the normalized output \bar{v} of the value network as

$$v_t = \bar{v} \sqrt{\text{Var}[G_{\text{train}}]} + \mathbb{E}[G_{\text{train}}] \quad (17)$$

to properly scale value predictions when the network is evaluated (which is done entirely internal to the network).

Hyperparameters and Tuning

The hyperparameters used during offline training and on-line execution are described in tables 4–6. Offline training refers to the BetaZero policy iteration steps that collect parallel MCTS data (*policy evaluation*) and then retrain the network (*policy improvement*). The online execution refers to using the BetaZero policy after offline training to evaluate its performance through online tree search. The main difference between these two settings is the final criteria used to select the root node action in MCTS. During offline training of problems with large action spaces (e.g., rock sample and mineral exploration), sampling root node actions according to the Q -weighted visit counts with a temperature τ ensures exploration. To evaluate the performance online, root node action selection takes the maximizing action of the Q -weighted visit counts. During training, we also evaluate a holdout set that uses the $\arg \max$ criteria to monitor the true performance of the learned policy.

The MCTS parameters for the mineral exploration problem were tuned using Latin hypercube sampling based on the lower-confidence bound of the returns. During training, the rock sample problems disabled progressive widening to expand on all actions and transition to a single belief state. Then for online execution, we tuned the DPW parameters as shown in figures 12–14. The problems train with a batch size of 1024 over 80% of 100,000 samples from one round of data collection ($n_{\text{buffer}} = 1$) using p_{dropout} of 0.2, 0.5, 0.7, respectively. Mineral exploration used a value function loss of MAE (otherwise MSE was used), light dark used the Adam optimizer (Kingma and Ba 2014), and the other problems used RMSProp (Hinton, Srivastava, and Swersky 2014) and batch norm regularization with a momentum of 0.7.

	Parameter *	LightDark(5)		LightDark(10)		Description
		Offline	Online	Offline	Online	
BetaZero policy iteration parameters	$n_{\text{iterations}}$	30	—	30	—	Number of offline BetaZero policy iterations.
	n_{data}	500	—	500	—	Number of parallel MCTS data generation episodes per policy iteration.
	bootstrap Q_0	false	false	false	false	Use bootstrap estimate for initial Q -value in MCTS.
Neural network parameters	n_{epochs}	50	—	50	—	Number of training epochs.
	α	1×10^{-4}	—	1×10^{-4}	—	Learning rate.
	λ	1×10^{-5}	—	1×10^{-5}	—	L_2 -regularization parameter.
MCTS parameters	n_{online}	100	1300	100	1000	Number of tree search iterations of MCTS.
	c	1	1	1	1	PUCT exploration constant.
	k_{action}	2.0	2.0	2.0	2.0	Multiplicative action progressive widening value.
	α_{action}	0.25	0.25	0.25	0.25	Exponential action progressive widening value.
	k_{state}	2.0	2.0	2.0	2.0	Multiplicative belief-state progressive widening value.
	α_{state}	0.1	0.1	0.1	0.1	Exponential belief-state progressive widening value.
	d	10	10	10	10	Maximum tree depth.
	τ	0	0	0	0	Exploration temperature for final root node action selection.
	z_q	1	1	1	1	Influence of Q -values in final criteria.
	z_n	1	1	1	1	Influence of visit counts in final criteria.

* Entries with “—” denote non-applicability and “.” denotes they are disabled.

Table 4: *BetaZero* parameters for the LIGHTDARK problems.

	Parameter	RockSample(15, 15)		RockSample(20, 20)		Description
		Offline	Online	Offline	Online	
BetaZero policy iteration parameters	$n_{\text{iterations}}$	50	—	50	—	Number of offline BetaZero policy iterations.
	n_{data}	500	—	500	—	Number of parallel MCTS data generation episodes per policy iteration.
	bootstrap Q_0	false	true	false	true	Use bootstrap estimate for initial Q -value in MCTS.
Neural network parameters	n_{epochs}	10	—	10	—	Number of training epochs.
	α	1×10^{-3}	—	1×10^{-3}	—	Learning rate.
	λ	1×10^{-5}	—	1×10^{-5}	—	L_2 -regularization parameter.
MCTS parameters	n_{online}	100	100	100	100	Number of tree search iterations of MCTS.
	c	50	50	50	50	PUCT exploration constant.
	k_{action}	.	5.0	.	.	Multiplicative action progressive widening value.
	α_{action}	.	0.9	.	.	Exponential action progressive widening value.
	k_{state}	.	1.0	1.0	1.0	Multiplicative belief-state progressive widening value.
	α_{state}	.	0.0	0.0	0.0	Exponential belief-state progressive widening value.
	d	15	15	4	4	Maximum tree depth.
	τ	1.0	0	1.5	0	Exploration temperature for final root node action selection.
	z_q	1	0.4	1	0.5	Influence of Q -values in final criteria.
	z_n	1	0.9	1	0.8	Influence of visit counts in final criteria.

Table 5: *BetaZero* parameters for the ROCKSAMPLE problems.

	Parameter	Offline	Online	Description
BetaZero policy iteration parameters	$n_{\text{iterations}}$	20	—	Number of offline BetaZero policy iterations.
	n_{data}	100	—	Number of parallel MCTS data generation episodes per policy iteration.
	bootstrap Q_0	false	false	Use bootstrap estimate for initial Q -value in MCTS.
Neural network parameters	n_{epochs}	10	—	Number of training epochs.
	α	1×10^{-6}	—	Learning rate.
	λ	1×10^{-4}	—	L_2 -regularization parameter.
MCTS parameters	n_{online}	50	50	Number of tree search iterations of MCTS.
	c	57	57	PUCT exploration constant.
	k_{action}	41.09	41.09	Multiplicative action progressive widening value.
	α_{action}	0.57	0.57	Exponential action progressive widening value.
	k_{state}	37.13	37.13	Multiplicative belief-state progressive widening value.
	α_{state}	0.94	0.94	Exponential belief-state progressive widening value.
	d	5	5	Maximum tree depth.
	τ	1.0	0	Exploration temperature for final root node action selection.
	z_q	1	1	Influence of Q -values in final criteria.
	z_n	1	1	Influence of visit counts in final criteria.

Table 6: *BetaZero* parameters for the MINERAL EXPLORATION problem.

Compute Resources

BetaZero was designed to use a single GPU to train the network and parallelize MCTS evaluations across available CPUs. Evaluating the networks on the CPU is computationally inexpensive due to the size of the neural networks (see fig. 16–18). This design was chosen to enable future research without a computational bottleneck. A single NVIDIA A100 was used with 80GB of memory on an Ubuntu 22.04 machine with 500 GB of RAM. Parallel data collection processes were run on 50 processes split evenly over two separate Ubuntu 22.04 machines: 1) with 40 Intel Xeon 2.3 GHz CPUs, and 2) with 56 Intel Xeon 2.6 GHz CPUs. Algorithm 5 (line 3) show where CPU parallelization occurs. In practice, the MCTS data generation simulations are the bottleneck of the offline component of BetaZero and not the network training—thus, parallelization is useful.

Open-Sourced Code and Experiments

The BetaZero algorithm has been open sourced and incorporated into the Julia programming language POMDPs.jl ecosystem (Egorov et al. 2017). Fitting into this ecosystem allows BetaZero to access existing POMDP models and can easily be compared to various POMDP solvers. The user constructs a `BetaZeroSolver` that takes parameters for policy iteration and data generation, parameters for neural network architecture and training, and parameters for MCTS (described in the tables above). The user may choose to define a method that inputs the belief b and outputs the belief

representation \tilde{b} used by the neural network (where the default computes the belief mean and std).

All experiments, including the experiment setup for the baseline algorithms with their heuristics, are included for reproducibility. Code to run MCTS data collection across parallel processes is also included. The code and experiments presented in this work are available online: <https://placeholder-url-for-github.com/BetaZero.jl>

BetaZero Algorithm

Algorithms 4–6 detail the full BetaZero policy iteration algorithm, including data collection and Q -weighted MCTS. Descriptions of parameters ψ are listed in tables 4–6.

Algorithm 4: BetaZero offline policy iteration.

Require: $\mathcal{P} \stackrel{\text{def}}{=} \langle \mathcal{S}, \mathcal{A}, \mathcal{O}, T, R, O, \gamma \rangle$: POMDP
Require: ψ : Parameters (includes $n_{\text{iterations}}$, n_{data} , n_{online} , and d)

```

1 function BETAZERO( $\mathcal{P}, \psi$ )
2    $f_\theta \leftarrow \text{INITIALIZE\_NETWORK}(\psi)$ 
3    $P_\theta, V_\theta \leftarrow f_\theta$  ▷ where  $(p, v) \leftarrow (P_\theta(\tilde{b}), V_\theta(\tilde{b}))$ 
4   for  $i \leftarrow 1$  to  $n_{\text{iterations}}$ 
5      $\mathcal{D} \leftarrow \text{COLLECT\_DATA}(\mathcal{P}, f_\theta, \psi)$  ▷ policy evaluation
6      $f_\theta \leftarrow \text{TRAIN}(f_\theta, \mathcal{D})$  ▷ policy improvement
7   end
8   return  $\beta_0^\pi(\mathcal{P}, f_\theta)$  ▷ BetaZero online policy (uses alg. 6)

```

Algorithm 5: Collect MCTS data offline for policy evaluation.

```

1 function COLLECTDATA( $\mathcal{P}, f_\theta, \psi$ )
2    $\mathcal{D} = \emptyset$ 
3   parallel for  $i \leftarrow 1$  to  $n_{\text{data}}$  ▷ parallelize MCTS runs across available CPUs
4     for  $t \leftarrow 1$  to  $T$ 
5        $a_t \leftarrow \text{MONTECARLOTREESearch}(\mathcal{P}, f_\theta, b_t, \psi)$  ▷ select next action through online planning
6        $\mathcal{D}_i^{(t)} \leftarrow \mathcal{D}_i^{(t)} \cup \{(b_t, \pi_{\text{tree}}^{(t)}, g_t)\}$  ▷ collect belief and policy data (placeholder for returns)
7        $s_{t+1} \sim T(\cdot | s_t, a_t)$ 
8        $o_t \sim O(\cdot | a_t, s_{t+1})$ 
9        $b_{t+1} \leftarrow \text{UPDATE}(b_t, a_t, o_t)$ 
10       $r_t \leftarrow R(s_t, a_t) \text{ or } R(s_t, a_t, s_{t+1})$  } transition the original POMDP
11    end
12     $g_t \leftarrow \sum_{k=t}^T \gamma^{(k-t)} r_k$  ▷ compute returns from observed rewards
13  end
14  return  $\mathcal{D}$ 

```

Algorithm 6: Monte Carlo tree search algorithm using Q -weighed visit counts.

```

1 function MONTECARLOTREESearch( $\mathcal{P}, f_\theta, b, \psi$ )
2    $\mathcal{M} \leftarrow \langle \mathcal{B}, \mathcal{A}, T_b, R_b, \gamma \rangle$  converted from the POMDP  $\mathcal{P}$  ▷ plan using the belief-state MDP
3   for  $i \leftarrow 1$  to  $n_{\text{online}}$ 
4      $\text{SIMULATE}(f_\theta, b, d)$  ▷ MCTS simulated planning to a depth  $d$  (algorithm 3)
5   end
6    $\pi_{\text{tree}}(b, a) \propto \left( \frac{\exp Q(b, a)}{\sum_{a'} \exp Q(b, a')} \right)^{z_q} \left( \frac{N(b, a)}{\sum_{a'} N(b, a')} \right)^{z_n}^{1/\tau}$  ▷  $Q$ -weighted visit counts eq. (8)
7    $\pi_{\text{tree}}(b, a_i) \leftarrow \pi_{\text{tree}}(b, a_i) / \sum_j \pi_{\text{tree}}(b, a_j)$  ▷ normalize to get a valid probability distribution
8   return  $a \sim \pi_{\text{tree}}(b, \cdot)$  ▷ sample root node action (let  $\tau \rightarrow 0$  to get arg max)

```
