# Observation Adaptation via Annealed Importance Resampling for Partially Observable Markov Decision Processes

## Yunuo Zhang, Baiting Luo, Ayan Mukhopadhyay, Abhishek Dubey

Vanderbilt University {yunuo.zhang, baiting.luo, ayan.mukhopadhyay, abhishek.dubey}@vanderbilt.edu

#### Abstract

Partially observable Markov decision processes (POMDPs) are a general mathematical model for sequential decisionmaking in stochastic environments under state uncertainty. POMDPs are often solved online, which enables the algorithm to adapt to new information in real time. Online solvers typically use bootstrap particle filters based on importance resampling for updating the belief distribution. Since directly sampling from the ideal state distribution given the latest observation and previous state is infeasible, particle filters approximate the posterior belief distribution by propagating states and adjusting weights through prediction and resampling steps. However, in practice, the importance resampling technique often leads to particle degeneracy and sample impoverishment when the state transition model poorly aligns with the posterior belief distribution, especially when the received observation is highly informative. We propose an approach that constructs a sequence of bridge distributions between the state-transition and optimal distributions through iterative Monte Carlo steps, better accommodating noisy observations in online POMDP solvers. Our algorithm demonstrates significantly superior performance compared to stateof-the-art methods when evaluated across multiple challenging POMDP domains.

### Introduction

Partially observable Markov decision processes (POMDPs) (Åström 1965) provide a general mathematical framework for modeling decision-making problems under uncertainty, where the true state of the environment is not fully observable and actions have probabilistic outcomes (Kaelbling, Littman, and Cassandra 1998). These models have been successfully applied to various real-world scenarios, including time-critical UAV search and rescue operations where efficient path planning must balance computational constraints with effective decision-making (Zhang et al. 2024). However, POMDPs face challenges to solve exactly due to the "curse of dimensionality" and the "curse of history", which make the computation of optimal policies non-scalable (Papadimitriou and Tsitsiklis 1987). To address these computational and scalability issues, online planning algorithms have emerged as a prominent approach. Instead of comput-

Copyright © 2025, Association for the Advancement of Artificial Intelligence (www.aaai.org). All rights reserved.

ing a full policy offline, online planning interleaves planning and execution by focusing computational resources on the current belief state and the immediate decisions to be made (Ross et al. 2008).

A critical component of online planning algorithms is the belief update process where belief denotes a distribution over the possible states. At each node in the search tree, the algorithm needs to update the belief state based on simulated actions and observations. However, performing exact belief updates is computationally infeasible in large state spaces due to the high dimensionality of the belief space (Rodriguez, Parr, and Koller 1999). To manage this, online planning algorithms approximate the belief state using sampling-based methods. They represent the belief state as a collection of sampled states, or particles, rather than as explicit probability distributions (Silver and Veness 2010). During the planning process, belief updates are performed by propagating these particles through the state-transition and observation models, using methods such as direct sampling and sequential importance sampling (Doucet and Johansen 2009). However, these approximation methods face a significant issue: the variance of their estimations increases exponentially with the search depth. This variance escalation is primarily due to the accumulation of uncertainty from observations, which provide imperfect information about the true current state. As the search depth increases, the divergence between the sampled belief and the target posterior can grow substantially, potentially limiting the effectiveness of these planning algorithms, especially in deeper searches where precise belief representation is crucial.

Annealed Importance Sampling (AIS) (Neal 1998) was originally developed to address similar challenges in Bayesian inference, where sampling from complex posterior distributions is often difficult. AIS creates a sequence of intermediate distributions that gradually transition from an easy-to-sample proposal distribution to the target posterior, while Annealed Importance Resampling (AIR) extends AIS by incorporating a resampling step at each temperature level (Herbst and Schorfheide 2019). AIR helps overcome the large divergence between proposal and target distributions by breaking the sampling problem into more manageable steps. However, adapting AIR for POMDP planning faces several challenges. First, it depends on smooth transition paths which aren't guaranteed in the highly dynamic be-

lief spaces of POMDPs. Second, AIR merely directs selection without addressing the underlying divergence between proposal and target distributions. Finally, AIR struggles with large belief spaces due to high computational costs, making it impractical for real-time POMDP planning where decisions must be made quickly. These challenges compound one another to make applying AIR to POMDP Planning Difficult.

Recognizing this, we propose AIROAS (Annealed Importance Resampling for Observation Adaptation Search), a novel approach that maintains the advantages of AIR while specifically addressing the barriers to its effective use in POMDP planning. AIROAS combines tree search with particle-based belief representation and enhanced annealed importance resampling techniques tailored for online POMDP planning. Our approach addresses the three key challenges through: (1) sigmoid-based tempering that strategically coordinates with tree search by concentrating computational effort at belief nodes with high uncertainty, creating smoother transitions precisely where belief updates are most prone to degeneracy; (2) a target inefficiency ratio mechanism that dynamically guides tree exploration by triggering resampling only in belief nodes with degenerate particle distributions, creating a synergistic relationship between search selection and particle refinement; (3) integration of AIR with excess uncertainty criterion (Ye et al. 2017; Wu et al. 2021) that pinpoints belief nodes with the highest proposal/target distribution mismatch, focusing computational effort on uncertain, high-value regions to overcome the computational cost barrier of traditional AIR in large belief spaces. These innovations enable AIROAS to provide more accurate state estimations at greater search depths while remaining computationally efficient for real-time decision-making in complex POMDP environments.

### **Background**

### **POMDPs**

Partially observable Markov decision processes (POMDPs) are a framework for modeling decision-making under uncertainty, defined by the tuple  $(S,A,T,R,\Omega,O)$ . Here, S is the set of states, A the set of actions, T(s'|s,a) the transition function, R(s,a) the reward function,  $\Omega$  the set of observations, and O(o|s',a) the observation function. In POMDPs, at time t, the agent maintains a belief state  $b_t$ , a probability distribution over states. This distribution is updated after action  $a_t$  and observation  $o_{t+1}$  using the equation:  $b_{t+1}(s') = \eta O(o_{t+1}|s',a_t) \sum_{s \in S} T(s'|s,a_t)b_t(s)$ , where  $\eta$  is a normalizing factor. This belief update process enables the agent to estimate its current state despite partial observability, facilitating informed decision-making in uncertain environments.

#### **Particle Filters**

Particle filters (Arulampalam et al. 2002), also known as sequential Monte Carlo (SMC) methods, are computational algorithms used to estimate the state of a dynamic system when it is observed through noisy measurements. In the

context of Partially observable Markov decision processes (POMDPs), particle filters offer an efficient way to approximate the belief state by maintaining a set of weighted samples (particles) that represent the probability distribution over possible states. The widely-known bootstrap particle filter algorithm (Gordon, Salmond, and Smith 1993) works by recursively propagating these particles through the state transition model and updating their weights based on new observations. In POMDPs, at each time step t, the particle filter aims to approximate the optimal posterior distribution:

$$p(s_t|o_t, s_{t-1}, a) \propto p(o_t|s_t, a) \cdot p(s_t|s_{t-1}, a)$$
 (1)

where  $s_t$  is the current state,  $o_t$  is the current observation,  $s_{t-1}$  is the previous state, and a is the action taken.

A crucial step in bootstrap particle filter is resampling (Kitagawa 1996), which addresses the problem of particle degeneracy. Over time, some particles may have negligible weights, leading to poor representation of the state space. Resampling involves randomly drawing particles with replacement from the current set, with probabilities proportional to their weights, and then resetting all weights to 1/N, where N is the number of particles. However overly frequent resampling can lead to sample impoverishment, where the particle set loses diversity and fails to adequately represent the full state space. To determine when to resample, the Effective Sample Size (ESS) (Forthofer, Lee, and Hernandez 2007) is often used:

$$ESS = \frac{1}{\sum_{i=1}^{N} (w_i^n)^2}$$
 (2)

where  $w_i^n$  are the normalized weights. When the ESS falls below a predefined threshold (typically N/2), resampling is triggered.

Although the bootstrap particle filter is relatively easy to implement, its accuracy can be significantly compromised when there is a substantial discrepancy between the optimal posterior distribution described in equation (1) and the state-transition distribution. This limitation becomes particularly pronounced in scenarios where the one-step received observation is highly informative about the true state, leading to a sharply peaked posterior distribution. In such cases, the bootstrap filter's reliance on the state-transition model for proposal generation may lead to particles being proposed in regions where the posterior has significant mass but the proposal density is low, resulting in many particles having negligible weights. Consequently, this can exacerbate the problem of sample degeneracy, necessitating more frequent resampling and potentially leading to sample impoverishment.

### **Annealed Importance Sampling**

Importance Sampling (IS) is a technique for estimating properties of a target distribution by sampling from a different, easier-to-sample proposal distribution and reweighting the samples (Kloek and van Dijk 1978). However, its major limitation lies in its inefficiency when the proposal distribution differs significantly from the target distribution — in such cases, only a few samples that happen to fall in the high-probability regions of the target distribution receive

high importance weights, while most samples have negligible weights. This leads to high variance in the estimates and poor sample efficiency.

Annealed Importance Sampling (AIS) constructs a sequence of intermediate distributions that gradually bridge the gap between the proposal and target distributions (Neal 1998). Given a target density p(x) and a proposal density q(x), AIS defines a sequence of intermediate distributions:

$$\pi_k(x) = p(x)^{\beta_k} q(x)^{(1-\beta_k)}$$
 (3)

where  $0=\beta_0\leq\cdots\leq\beta_K=1$  represents a sequence of inverse temperatures. The method proceeds by first drawing an initial sample  $x_0\sim\pi_0(x)$ , then evolving this sample through a series of transition kernels:

$$x_k \sim T_k(x_k|x_{k-1}) \text{ for } k = 1, \dots, K$$
 (4)

where each  $T_k$  is constructed to leave  $\pi_k$  invariant. The final importance weight is computed as the product of ratios between successive distributions:

$$w = \frac{\pi_1(x_1)}{\pi_0(x_1)} \cdots \frac{\pi_K(x_K)}{\pi_{K-1}(x_K)}$$
 (5)

This gradual transition through intermediate distributions helps AIS overcome the limitations of standard importance sampling by maintaining better overlap between successive distributions, resulting in more reliable estimates.

#### **Related Work**

Monte Carlo tree search (MCTS) has demonstrated success in solving large POMDPs online. POMCP (Silver and Veness 2010) pioneered this approach by combining a UCT-based tree search with particle filtering for belief updates, but faces challenges with continuous observation spaces. DESPOT and its variant AR-DESPOT (Ye et al. 2017) improve upon POMCP by focusing the search on a fixed set of scenarios and using dual bounds to guide exploration, making it more robust for large discrete problems. POMCPOW (Sunberg and Kochenderfer 2018) extends these ideas to continuous state-action-observation spaces by incorporating progressive widening and weighted particle filtering.

A common thread among these approaches is their reliance on particle filtering techniques, specifically bootstrap particle filtering or Sequential Importance Resampling (SIR), to update belief states. While these filtering methods are computationally efficient, they can sometimes lead to particle degeneracy issues, especially in continuous observation spaces or when observations are unlikely under the current belief.

Most recently, AdaOPS (Wu et al. 2021) attempts to address this challenge using KLD-sampling, which dynamically adapts the number of particles based on the Kullback-Leibler divergence between the true and approximated distributions. However, KLD-sampling exhibits two major limitations: First, it requires discretizing the state space into bins for divergence calculation, which becomes computationally prohibitive especially in high-dimensional

spaces (Li, Sun, and Sattar 2013). Second, and more critically, while it adjusts particle quantities, it does not modify the particle values themselves, leading to potential underestimation of distribution variance and failure to capture multimodal aspects of the target distribution. In contrast, Annealed Importance Sampling (AIS) offers several compelling advantages: it constructs a sequence of intermediate distributions that gradually bridge the prior and posterior distributions, effectively maintaining particle diversity while preventing degeneracy. This annealing process enables particles to adaptively migrate toward regions of high posterior probability, making it particularly effective when dealing with concentrated observation likelihoods or significant disparities between prior and posterior distributions. The gradual transition through intermediate distributions allows AIS to better capture the full structure of multimodal distributions and provide more accurate representations of distribution tails compared to KLD-sampling approaches.

Previous work has explored adding tempering iterations to particle filters, where the optimal posterior distribution is constructed adaptively through a sequence of Monte Carlo steps. For example, M. Johansen (2015) developed a blocktempered particle filter that uses bridge distributions to gradually adapt particle values, while Herbst and Schorfheide (2019) proposed a tempered particle filter that sequentially reduces inflated measurement error variance in Dynamic Stochastic General Equilibrium (DSGE) models. However, these tempering approaches have been primarily studied in the context of state estimation and system identification, but not yet explored within the POMDP planning literature. To our knowledge, we are the first to investigate using bridge distributions to connect target and proposal distributions in particle filtering specifically for POMDP planning, combining the strengths of both sequential Monte Carlo methods and POMDP planning algorithms.

### **Approach**

In this section, we present AIROAS (Annealed Importance Resampling for Observation Adaptation Search), a novel approach for POMDP planning that combines tree search with particle-based belief representation and annealed importance resampling. The complete procedure is detailed in algorithm 1 and fig. 1. AIROAS constructs a search tree that alternates between belief nodes and action nodes. Each belief node approximates the optimal posterior distribution using a set of weighted particles, as illustrated in fig. 1, and maintains both upper and lower bounds of the optimal value. Using a sequence of bridging distributions, the algorithm refines the value bounds and adjusts particle states and weights to better approximate the optimal posterior distribution described in equation (1). At each timestep, starting from the current belief  $b_0$  as the root node, AIROAS expands the belief tree by exploring various paths from the initial states (Line 3 in algorithm 1). A key characteristic of this structure is that sibling nodes share identical particle states but maintain distinct weight distributions. This configuration persists until annealed importance resampling is applied to mutate the particle values, as we will discuss below.

### Algorithm 1: AIROAS

```
Require: initial belief \bar{b}_0, maximum depth MAX_DEPTH
 1: while Time Allowed and l(\bar{b}_0) < u(\bar{b}_0) do
 2:
          while depth(\bar{b}) < MAX_DEPTH do
 3:
 4:
               if b is a leaf node then
                    \bar{b} \leftarrow \text{Annealed Importance Resampling}(\bar{b})
 5:
                    ExpandAndBackup(\bar{b})
 6:
 7:
               end if
               \bar{b} \leftarrow \text{Selection}(\bar{b})
 8:
 9:
          end while
10:
          if depth(\bar{b}) \geq MAX_DEPTH then
11:
               u(\bar{b}) \leftarrow l(\bar{b})
12:
          end if
13: end while
14: return arg \max_{a \in \mathcal{A}} l(\bar{b}_0, a)
```

### Algorithm 2: Selection( $\bar{b}$ )

```
1: a^* = \arg\max_{a \in A} u(\bar{b}, a)

2: o^* \leftarrow \arg\max_{o \in O_{\bar{b}, a}} \hat{p}(o|\bar{b}, a) \text{EU}(\tau(\bar{b}, a^*, o))

3: \bar{b} \leftarrow \tau(\bar{b}, a^*, o^*)
```

#### Selection

Similar to Wu et al. (2021); Ye et al. (2017); Smith and Simmons (2004), during online planning, AIROAS maintains both upper and lower bounds,  $u(\bar{b},a)$  and  $l(\bar{b},a)$ , for each action node a at a belief node  $\bar{b}$ . These bounds estimate the optimal value that can be achieved by taking action a at belief  $\bar{b}$ . The action selection follows an optimistic strategy at each belief node  $\bar{b}$ , AIROAS chooses the action  $a^*$  that maximizes the upper bound(Line 1 in algorithm 2):

$$a^* = \operatorname*{arg\,max}_{a \in A} u(\bar{b}, a) \tag{6}$$

After selecting an action, we need to choose which observation branch to explore next. This choice is guided by the probability-weighted excess uncertainty (EU) criterion (Ye et al. 2017; Smith and Simmons 2004) (Line 2 in algorithm 2). For each possible observation o after taking action  $a^*$ , we compute:

$$\hat{p}(o|\bar{b}, a^*) \cdot EU(\tau(\bar{b}, a^*, o)) \tag{7}$$

where  $\hat{p}(o|\bar{b},a^*)$  is an estimation of the probability  $p(o|\bar{b},a^*)$  i.e estimates the probability of receiving observation o after taking action  $a^*$  at belief  $\bar{b}, \tau(\bar{b},a^*,o)$  represents the updated belief after executing  $a^*$  and receiving observation o, and  $EU(\cdot)$  measures the excess uncertainty at a belief node (Wu et al. 2021). The excess uncertainty quantifies the degree of value uncertainty at each belief node:

$$EU(\bar{b}) = [u(\bar{b}) - l(\bar{b})] - \frac{\xi[u(\bar{b}_0) - l(\bar{b}_0)]}{\gamma^{d(\bar{b})}}$$
(8)

where  $u(\bar{b})$  and  $l(\bar{b})$  are the upper and lower bounds at belief  $\bar{b}$ ,  $\bar{b}_0$  is the root belief,  $d(\bar{b})$  denotes the depth of belief

### Algorithm 3: ExpandAndBackup( $\bar{b}$ )

```
1: for a \in \mathcal{A} do
           Create Action Node (\bar{b}, a)
 2:
 3:
           b' \leftarrow \text{Empty Particles Set}
           \mathcal{O} \leftarrow Empty Observation Set
 4:
 5:
           for s \in \bar{b} do
                 (s\prime,o,r)\sim\mathcal{G}(s,a)
 6:
                 \mathcal{O} \leftarrow \mathcal{O} \cup o
 7:
                 weight(st) \leftarrow 1
 8:
 9:
                 b' \leftarrow b' \cup s'
10:
           end for
           for o \in \mathcal{O} do
11:
12:
                 \tilde{b} \leftarrow \text{DeepCopy}(b')
                 for s \in \tilde{b} do
13:
                       weight(s) = p(o|s, a) \cdot weight(s)
14:
15:
                 end for
                 Create Belief Node \tilde{b}
16:
17:
           end for
18: end for
19: BackUp(\bar{b})
```

 $\bar{b}$  in the tree,  $\gamma$  is the discount factor, and  $\xi \in (0,1)$  controls the target uncertainty reduction at the root. The observation  $o^*$  that maximizes the weighted excess uncertainty is selected for exploration, effectively focusing the search on belief states with high uncertainty and high probability of being reached (Wu et al. 2021).

### **Expand And Backup**

When encountering a leaf node, we first perform annealed importance resampling to adjust the particle states and weights (detailed in the next section). Then, AIROAS expands this leaf node by creating child nodes for each possible action a. For each action node, the expansion process involves propagating the parent belief node's particles forward using a simulator  $\mathcal{G}$  that generates state transitions (Line 6 in algorithm 3):

$$(s', o, r) \sim \mathcal{G}(s, a)$$

The expansion continues by creating new belief nodes for each unique observation encountered during particle propagation. Specifically, after propagating each particle through  $\mathcal G$ , we obtain a set of updated particle states, initially assigning each particle a weight of 1. For each observation o generated during this process, we create a new belief node where the particle weights are updated according to the observation likelihood p(o|s,a) — the probability of receiving observation o given the particle state s and action a. This results in a tree structure where each action node branches into multiple belief nodes based on possible observations, with each belief node containing weighted particles that represent the updated belief state.

After expanding belief nodes, we perform a backup operation to update the bounds of their ancestor nodes. As in Ye et al. (2017); Wu et al. (2021), this backup process is implemented by recursively applying the Bellman equation, which

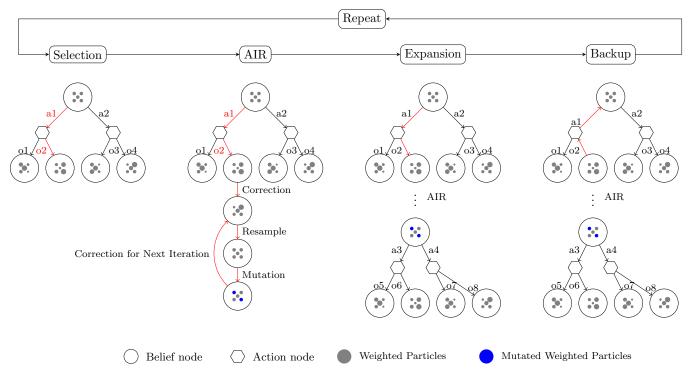


Figure 1: Process of AIROAS Tree Search. AIROAS constructs a search tree that alternates between belief nodes and action nodes. AIR represents Annealed Importance Resampling.

calculates the optimal value of a belief state in two parts. First, we consider the expected immediate reward for taking an action in the current belief state, calculated by looking at each possible state, weighing its reward by how likely we think we are in that state according to our current belief. Second, we consider the long-term value by looking at all possible observations we might receive after taking an action. For each possible observation, we calculate how likely we are to see that observation, determine what our new belief state would be after seeing it, consider the value of that resulting belief state, and weight this value by the probability of getting that observation. The total value is the sum of these immediate and future components, and we choose the action that maximizes this total value. In practice, since we can't compute exact values over continuous states and observations, we approximate this calculation using discrete sums over our particle-based belief representation.

#### **Annealed Importance Resampling**

As discussed above, the belief of each belief node is generated using a state-transition function from its parent belief node. Suppose the leaf node represents belief at timestep t and its parent belief node represents belief at timestep t-1, while our simulator  $\mathcal G$  implements the state-transition function  $p(s_t|s_{t-1},a)$  during expansion, our ultimate objective is to have the final particle distribution approximate the optimal posterior distribution  $p(s_t|o_t,s_{t-1},a)$  referenced in eq. (1).

The primary challenge in this context stems from the impossibility of directly sampling from the optimal posterior

#### Algorithm 4: Annealed Importance Resampling $(\bar{b})$

```
Require: b from current belief node, sequence of tempering
     parameters 0 = \beta_0 \le \cdots \le \beta_K = 1, target inefficiency
     ratio r^*, transition kernels T_1(s, s'), \ldots, T_K(s, s')
 1: for \beta_k in \{\beta_0, \ldots, \beta_K\} do
          if k == 0 then
 2:
               Continue
 3:
 4:
          end if
 5:
          weights(b) \leftarrow Update_Weights(b, \beta_k, \beta_{k-1})
          if In\_eff(\bar{b}) \le r^* then
 6:
 7:
               break
 8:
          else
 9:
               r^* \leftarrow \operatorname{In\_eff}(\overline{b})
10:
          end if
          Resample b
11:
          \bar{b} \leftarrow \text{Mutation}(\bar{b})
12:
13: end for
```

distribution. Previous works (Silver and Veness 2010; Ye et al. 2017; Sunberg and Kochenderfer 2018; Wu et al. 2021) approximates sampling from this optimal posterior distribution using Sequential Importance Resampling, treating the state-transition distribution as proposal distribution and the optimal posterior distribution as target distribution. However, this approach suffers from particle degeneracy issues, particularly when observations are highly informative about the true state, leading to a sharply peaked posterior distribution and a significant divergence between the optimal pos-

terior and the proposal distribution derived from the current belief (Wu et al. 2021).

To address the issue, Annealed Importance Sampling (Neal 1998) introduced in Background offers an effective solution by constructing intermediate bridging distributions between the proposal distribution and target distribution, gradually guiding particles toward high-probability regions of the target distribution by constructing a sequence of intermediate distributions parameterized by tempering parameters  $0 = \beta_0 \le \beta_1 \le \cdots \le \beta_K = 1$ .

Building upon these principles, we integrate Annealed Importance Resampling (AIR), which extends AIS with a resampling step, into our AIROAS algorithm to address the particle degeneracy issue inherent in POMDP planning. Our approach adapts this concept to the POMDP planning context by designing intermediate distributions that account for observations. By gradually bridging the proposal and target distributions, particles are adaptively reweighted and resampled, allowing them to migrate toward regions of high posterior probability. This gradual transition through intermediate distributions not only addresses degeneracy but also enables our approach to capture the multimodal characteristics of the optimal posterior distribution more effectively.

These tempering parameters are generated using a sigmoid function applied to a linearly spaced sequence between  $10^{-3}$  and 1. Specifically, for a sequence  $\{x_1, x_2, \cdots x_K\}$ linearly spaced between  $10^{-3}$  and 1, each tempering parameter  $\beta_i$  is computed as:

$$\beta_i = \frac{1}{1 + e^{-10(x_i - 0.5)}} \tag{9}$$

This sigmoid transformation ensures a smooth progression of  $\beta$  values from near zero to one, with denser sampling in the middle range.

The intermediate distributions are characterized by their probability densities:

$$\pi_k = p(s_t|o_t, s_{t-1}, a)^{\beta_k} p(s_t|s_{t-1}, a)^{(1-\beta_k)}$$
 (10)

When k = K, this formulation exactly matches the optimal posterior distribution. By progressively transforming particles sampled from the prior distribution  $p(s_t|s_{t-1},a)$ through the sequence  $\pi_1, \pi_2, \dots, \pi_K$ , we can effectively approximate particle states drawn from the optimal posterior distribution.

It is worth noting that while our current experiments use sequential execution for fair comparison with baseline methods, the annealing approach naturally supports parallelization across particles and tempering steps. This computational advantage offers significant potential for further acceleration, and we will provide detailed parallel implementation as part of the released code.

The overall procedure of annealed importance resampling is described in algorithm 4.

### Weights

To update weights of the particles at iteration k and timestep t, for every  $\beta_k$  in the sequence  $\{\beta_0, \dots, \beta_K\}$ , we transform the particle approximation from  $\pi_{k-1}$  to  $\pi_k$  by recalculating the weights. Assume the particle set at iteration k and Algorithm 5: Mutation( $\bar{b}$ )

- 1: for particle  $s_{t(k-1)}^{j}$  in  $\bar{b}$  do
- $T_k(s_{k-1}^j, \cdot) \leftarrow \mathcal{N}(s_{t(k-1)}^j, \sigma^2 \cdot I)$
- 3:
- 4:
- $\begin{aligned} s_{tk}^j &\sim T_k(s_{k-1}^j, \cdot) \\ \text{Compute the acceptance probability p} \\ s_{tk}^j &\leftarrow s_{tk}^j \text{ for probability p, otherwise } s_{tk}^j \leftarrow s_{t(k-1)}^j \end{aligned}$

timestep t is  $S_{tk}$ , for each particle  $s_{tk}^{\jmath}$  in the particle set  $S_{tk}$ at iteration k, the weight is updated as:

$$w_{tk}^{j} \propto \frac{p(o_t \mid s_{tk}^{j}, a_t)^{\beta_k}}{p(o_t \mid s_{tk}^{j}, a_t)^{\beta_{k-1}}} \cdot w_{t(k-1)}^{j}$$
(11)

This update reflects the incremental adjustment of weights based on the annealing parameter  $\beta_k$ . Specifically, at iteration k, the importance weight of each particle  $s_{tk}^{\jmath}$  is updated by multiplying its previous weight  $\boldsymbol{w}_{t(k-1)}^{j}$  with an incremental weight. This incremental weight is computed as the ratio of two observation likelihood terms: the observation density function raised to the current tempering parameter  $\beta_k$  in the numerator, and the same function raised to the previous tempering parameter  $\beta_{k-1}$  in the denominator. Specifically, the incremental weight is  $\frac{p(o_t|s_{tk}^j,a_t)^{\beta_k}}{p(o_t|s_{tk}^j,a_t)^{\beta_{k-1}}}$ , where  $o_t$  is the observation received at timestep t and  $a_t$  is the action taken.

After updating the particle weights, we calculate their inefficiency score in a similar way as in Herbst and Schorfheide (2019). For a particle set containing M particles, the inefficiency score is computed as:

In\_eff(
$$\bar{b}_{tk}$$
) =  $\frac{1}{M} \sum_{i=1}^{M} (\frac{w_{tk}^{j}}{\frac{1}{M} \sum_{i=1}^{M} w_{tk}^{j}})^{2}$  (12)

This score reflects the variance in particle weights. As in Herbst and Schorfheide (2019) and specified in algorithm 4, we maintain a pre-defined target inefficiency ratio  $r^*$ . If the inefficiency score exceeds  $r^*$ , indicating that the variance of the resulting weights remains high, we continue adjusting the particle states and weights. Otherwise, we exit the current iteration.

### Mutation

For each particle  $s_{t(k-1)}^{j}$ , the algorithm employs a Markov transition kernel  $T_k(s_{k-1}^j, s_k^j)$  to evolve the state from time k-1 to k. This transition kernel is constructed to preserve  $\pi_k$  as its invariant distribution, as specified in eq. (4).

The transition mechanism implements a Metropolis-Hastings kernel (Hastings 1970) with a multivariate Gaussian proposal distribution. Given a state vector  $s_{t(k-1)}^{j}$ that can be decomposed into components  $s_{t(k-1)}^{j} \;\;= <$  $x, y, \dots >$ , the proposal distribution is centered at  $s_{t(k-1)}^j$ 

with covariance structure. The forward proposal distribution i.e the transition kernel generates new states according to:

$$s_{tk}^j \sim \mathcal{N}(s_{t(k-1)}^j, \sigma^2 \cdot I)$$

where I denotes the identity matrix and  $\sigma^2$  is a scaling parameter calibrated proportionally to the L1 distance between the current state  $s_{t(k-1)}^j$  and the observation  $o_t$  in the state space. Then we define its reverse transition kernel  $\tilde{T}_k(s_k^j)$  as:

$$\mathcal{N}(s_{tk}^j, \sigma^2 \cdot I)$$

The acceptance rate shown in algorithm 5 is computed as:

$$p_{accept} = \min \left\{ 1, \frac{T_k(s_{t(k-1)}^j, s_{tk}^j) \cdot p(o_t | s_{tk}^j, a_t)^{\beta_k}}{\tilde{T}_k(s_{tk}^j, s_{t(k-1)}^j) \cdot p(o_t | s_{t(k-1)}^j, a_t)^{\beta_k})} \right\}$$
(13)

Then we update the particle state:

$$s_{tk}^{j} = \begin{cases} s_{tk}^{j} \text{with prob. } p_{\text{accept}} \\ s_{t(k-1)}^{j} \text{with prob. } 1 - p_{\text{accept}} \end{cases}$$
 (14)

The iterative process continues until either the particle approximation converges to the optimal posterior distribution, or the variance of the particle weights becomes less than or equal to the threshold  $r^*$ .

### **Experiments**

In this section, we evaluate our method on several domains and demonstrate superior performance over state-of-the arts methods.

#### **Baseline Approaches**

We evaluate our approach against four state-of-the-art baselines: POMCP (Silver and Veness 2010), ARDESPOT (Ye et al. 2017), POMCPOW (Sunberg and Kochenderfer 2018), and AdaOPS (Wu et al. 2021) on four domains: Light Dark (LD), Tag, Laser Tag, and Rock Sample (RS) (Platt et al. 2010; Egorov et al. 2017; Smith and Simmons 2004). For ARDESPOT and AdaOPS, we employ different initialization strategies: in Light Dark, Tag, and Laser Tag domains, we use domain-specific independent bounds based on maximum and minimum achievable discounted rewards, while for RockSample, bounds are initialized using heuristics. For POMCPOW, we optimize the maxUCB parameter—which governs action selection at each node—by testing values in {1.0, 10.0, 20.0} and selecting the best-performing configuration for each domain. The performance comparison across all domains is presented in table 1. Our experimental results for the baselines may differ from previously published results due to modifications in implementation code, differences in computing infrastructure, and alternative bound initialization methods. Specifically, while Wu et al. (2021) uses heuristics for bound initialization in Light Dark, Tag, and Laser Tag, we initialize bounds based on maximum and minimum achievable discounted rewards. However, we maintain consistent initialization methods across all approaches within each domain to ensure fair comparison.

#### **Domains**

We validate our approach using four standard POMDP environments:

- 1. **Light Dark** (Platt et al. 2010): A continuous-state problem where an agent navigates to a target with state-dependent observation noise. Observations become more accurate in "light" regions and less accurate in "darker" regions. We vary step size  $\alpha$  to modulate problem complexity.
- 2. **Tag** (Egorov et al. 2017): An agent aims to tag an opponent by occupying the same grid cell and executing a tag action. The agent observes its own position but can only detect the opponent when they share a cell. The opponent follows a stochastic evasion policy.
- 3. Laser Tag (Egorov et al. 2017): A more challenging Tag variant where the agent must locate and tag a target without knowing either position initially. The agent relies on noisy sensor readings to infer both positions.
- 4. RockSample (Smith and Simmons 2004): A robot explores a grid to collect valuable rocks while avoiding bad ones. Rock quality can only be determined through noisy sensor readings that decrease in accuracy with distance. The robot can move, sample rocks, or use sensors to check rock status.

### **Experiment Settings**

The main configurations for all baselines are detailed in Baseline Approaches. For fair comparison, we ensure AIROAS's maximum allowed time per decision matches those of comparable baselines with similar parameters. While POMCP and POMCPOW use iterations rather than time limits, we maintain consistency by keeping the number of iterations identical during testing. AIROAS uses the same bound initialization methods as ARDESPOT and AdaOPS. For AIROAS-specific configurations, we generate 100 tempering parameters  $\beta_k$  ranging from 0 to 1 using the sigmoid transformation described in eq. (9). We evaluate AIROAS using different target inefficiency ratios  $r^* \in \{2.0, 3.0, 5.0, 10.0\}$  and select the best-performing value for each domain.

Regarding computational efficiency, our performance analysis shows that SIR requires 0.12s per iteration, while AIR takes 0.14s. In worst-case scenarios (ratio=2.0, 100 tempering parameters), annealing adds 0.13s overhead. However, these computational costs are justified by AIR's superior particle diversity and estimation accuracy. We also plan to implement parallel annealing to reduce computation time while maintaining performance benefits.

All experiments were conducted on a computer equipped with an Intel(R) Core(TM) i9-13900KS processor with 32 CPUs.

#### **Results**

The average discounted return and its standard error of mean (SEM) are presented in table 1. In all these domains, AIROAS outperforms the other solvers with suitable target inefficiency ratio and tempering parameters.

	$LD (\alpha = 0.5)$	$LD (\alpha = 1.0)$	Tag	Laser Tag	RS(11, 11)	RS(15,15)
$ \mathcal{S} $	$\infty$	$\infty$	870	4830	247,808	7,372,800
$ \mathcal{A} $	3	3	5	5	20	16
$ \mathcal{O} $	$\infty$	$\infty$	30	$\sim 1.5 \times 10^6$	3	3
POMCP	$0.357 \pm 0.22$	$0.691 \pm 0.41$	$-18.530 \pm 0.40$	$-18.112 \pm 0.18$	$12.252 \pm 0.44$	$8.207 \pm 0.53$
ARDESPOT	$-1.24 \pm 0.25$	$0.760 \pm 0.41$	$-13.704 \pm 0.71$	$-15.982 \pm 0.64$	$18.113 \pm 0.64$	$16.111 \pm 0.66$
POMCPOW	$1.555 \pm 0.21$	$1.468 \pm 0.49$	$-16.783 \pm 0.42$	$-15.263 \pm 0.40$	$13.521 \pm 0.67$	$9.141 \pm 0.44$
AdaOPS	$2.231 \pm 0.13$	$2.762 \pm 0.27$	$-9.920 \pm 0.70$	$-13.427 \pm 0.67$	$21.904 \pm 0.59$	$19.091 \pm 0.57$
AIROAS	$\textbf{2.303} \pm \textbf{0.43}$	$\boldsymbol{3.102 \pm 0.36}$	$-8.451 \pm 0.75$	$\bf -12.931 \pm 1.22$	$22.872 \pm 0.51$	$20.246 \pm 0.50$

Table 1: Performance Comparison between AIROAS and baselines.  $LD(\alpha)$  stands for Light Dark with step size  $\alpha$ . RS(n,m) stands for the Rock Sample with  $n \times n$  map and m rocks.  $\infty$  means continuous state (or observation) space. The results represent the average discounted return and its standard error of mean (SEM) for each method and each domain. Higher discounted return indicates better performance.

In the Light Dark Problem, we evaluate performance with step sizes  $\alpha=0.5$  and  $\alpha=1.0$ . The state consists of position coordinates and termination status (terminated or not), while observations provide noisy measurements of the current position. For particle state mutation, we generate a Gaussian distribution over the position components, where the variance is proportional to the distance between the state's position and the received observation. Through this mutation strategy in bridging distributions, combined with weight adjustments, we observe significant mitigation of the sample impoverishment problem, leading to improved performance.

In the Tag problem, the state consists of the robot's position, target's position, and tagging status. While the agent's position is fully observable, the target's position is only observed when both actors occupy the same cell. For particle state mutation, we generate a Gaussian distribution exclusively over the target's position, as the agent's position is known. Given this limited observability structure, we did not anticipate significant performance improvements from Annealed Importance Resampling. The observed improvements may be attributed to the compact state and observation spaces, where the mutation step effectively explores different states within this limited domain.

The Laser Tag problem represents a more challenging variant of the Tag problem. The state space comprises the agent's position, target's position, and terminal signal, with neither position directly observable. The observation consists of noisy range readings from laser sensors in eight different directions. For particle state mutation, we construct a multivariate Gaussian distribution over both agent and target positions. Given these high-dimensional observations, which induce multimodal distributions over the state space, we anticipate Annealed Importance Sampling to demonstrate significant advantages in this domain.

In the RockSample problem, the state space comprises the robot's position and the quality (good or bad) of each rock. With scenarios containing 11 or 15 rocks, this results in an extremely high-dimensional state space with a large number of states. In contrast, the observation structure is simple, providing information about only one rock's quality at a time

(or no observation). For particle state mutation, we focus specifically on mutating the quality of the rock for which the observation provides information, rather than mutating the entire state vector.

### **Conclusion**

We presented AIROAS, a novel online POMDP solver that leverages Annealed Importance Resampling to address observation uncertainty. By constructing bridge distributions between state-transition and optimal distributions, our approach effectively mitigates particle degeneracy in traditional particle filtering. Empirical evaluations across diverse POMDP domains demonstrate consistent performance advantages over state-of-the-art baselines, particularly for domains with complex observation spaces or highly informative observations that create sharply peaked posterior distributions.

While our current approach uses sigmoid transformation for generating tempering parameters, future work could explore adaptive scheduling that dynamically adjusts based on belief uncertainty and parallel implementation to reduce computational overhead. Additionally, incorporating temporal abstraction techniques would enable AIROAS to efficiently handle high-dimensional continuous action spaces through learned macro-actions, extending its applicability to robotic control and autonomous driving domains (Luo et al. 2025). These enhancements could further improve AIROAS's efficiency and effectiveness across both discrete and continuous action spaces.

### Acknowledgements

This material is based upon work sponsored by the National Science Foundation (NSF) under Grant CNS-2238815 and by the Defense Advanced Research Projects Agency (DARPA) under the Assured Neuro Symbolic Learning and Reasoning program. Results presented in this paper were obtained using the Chameleon testbed supported by the National Science Foundation. Any opinions, findings, conclusions, or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the NSF or the DARPA.

### References

- Arulampalam, M. S.; Maskell, S.; Gordon, N. J.; and Clapp, T. 2002. A tutorial on particle filters for online nonlinear/non-Gaussian Bayesian tracking. *IEEE Trans. Signal Process.*, 50: 174–188.
- Doucet, A.; and Johansen, A. 2009. A Tutorial on Particle Filtering and Smoothing: Fifteen Years Later. *Handbook of Nonlinear Filtering*, 12.
- 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. Forthofer, R. N.; Lee, E. S.; and Hernandez, M. 2007. 15 Analysis of Survey Data. In Forthofer, R. N.; Lee, E. S.; and Hernandez, M., eds., *Biostatistics (Second Edition)*, 421–
- Gordon, N. J.; Salmond, D.; and Smith, A. F. M. 1993. Novel approach to nonlinear/non-Gaussian Bayesian state estimation. In *IEE Proceedings F Radar and Signal Processing*.

443. San Diego: Academic Press, second edition edition.

ISBN 978-0-12-369492-8.

- Hastings, W. K. 1970. Monte Carlo Sampling Methods Using Markov Chains and Their Applications. *Biometrika*, 57(1): 97–109.
- Herbst, E.; and Schorfheide, F. 2019. Tempered particle filtering. *Journal of Econometrics*, 210(1): 26–44. Annals Issue in Honor of John Geweke "Complexity and Big Data in Economics and Finance: Recent Developments from a Bayesian Perspective".
- Kaelbling, L. P.; Littman, M. L.; and Cassandra, A. R. 1998. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101(1): 99–134.
- Kitagawa, G. 1996. Monte Carlo Filter and Smoother for Non-Gaussian Nonlinear State Space Models. *Journal of Computational and Graphical Statistics*, 5(1): 1–25.
- Kloek, T.; and van Dijk, H. K. 1978. Bayesian Estimates of Equation System Parameters: An Application of Integration by Monte Carlo. *Econometrica*, 46(1): 1–19.
- Li, T.; Sun, S.; and Sattar, T. 2013. Adapting sample size in particle filters through KLD-resampling. *Electronics Letters*, 49(12): 740–742.
- Luo, B.; Pettet, A.; Laszka, A.; Dubey, A.; and Mukhopadhyay, A. 2025. Scalable Decision-Making in Stochastic Environments through Learned Temporal Abstraction. In *The Thirteenth International Conference on Learning Representations*.
- M. Johansen, A. 2015. On Blocks, Tempering and Particle MCMC for Systems Identification. *IFAC-PapersOnLine*, 48(28): 969–974. 17th IFAC Symposium on System Identification SYSID 2015.
- Neal, R. M. 1998. Annealed Importance Sampling. arXiv:physics/9803008.
- Papadimitriou, C. H.; and Tsitsiklis, J. N. 1987. The Complexity of Markov Decision Processes. *Mathematics of Operations Research*, 12(3): 441–450.

- Platt, R.; Tedrake, R.; Kaelbling, L. P.; and Lozano-Perez, T. 2010. Belief space planning assuming maximum likelihood observations. In *Robotics: Science and Systems*.
- Rodriguez, A.; Parr, R.; and Koller, D. 1999. Reinforcement Learning Using Approximate Belief States. In Solla, S.; Leen, T.; and Müller, K., eds., *Advances in Neural Information Processing Systems*, volume 12. MIT Press.
- Ross, S.; Pineau, J.; Paquet, S.; and Chaib-draa, B. 2008. Online Planning Algorithms for POMDPs. *Journal of Artificial Intelligence Research*, 32: 663–704.
- Silver, D.; and Veness, J. 2010. Monte-Carlo Planning in Large POMDPs. In Lafferty, J.; Williams, C.; Shawe-Taylor, J.; Zemel, R.; and Culotta, A., eds., *Advances in Neural Information Processing Systems*, volume 23. Curran Associates, Inc.
- Smith, T.; and Simmons, R. 2004. Heuristic search value iteration for POMDPs. In *Proceedings of the 20th Conference on Uncertainty in Artificial Intelligence*, UAI '04, 520–527. Arlington, Virginia, USA: AUAI Press. ISBN 0974903906. Sunberg, Z.; and Kochenderfer, M. 2018. Online algorithms for POMDPs with continuous state, action, and observation spaces. arXiv:1709.06196.
- Wu, C.; Yang, G.; Zhang, Z.; Yu, Y.; Li, D.; Liu, W.; and Hao, J. 2021. Adaptive Online Packing-guided Search for POMDPs. In Ranzato, M.; Beygelzimer, A.; Dauphin, Y.; Liang, P.; and Vaughan, J. W., eds., *Advances in Neural Information Processing Systems*, volume 34, 28419–28430. Curran Associates, Inc.
- 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.
- Zhang, Y.; Luo, B.; Mukhopadhyay, A.; Stojcsics, D.; Elenius, D.; Roy, A.; Jha, S.; Maroti, M.; Koutsoukos, X.; Karsai, G.; and Dubey, A. 2024. Shrinking POMCP: A Framework for Real-Time UAV Search and Rescue. In 2024 International Conference on Assured Autonomy (ICAA), 48–57.
- Åström, K. 1965. Optimal control of Markov processes with incomplete state information. *Journal of Mathematical Analysis and Applications*, 10(1): 174–205.