

# Probabilistic Loss and its Online Characterization for Simplified Decision Making Under Uncertainty

Andrey Zhitnikov\* and Vadim Indelman†

\*Technion Autonomous Systems Program †Department of Aerospace Engineering

Technion - Israel Institute of Technology, Haifa 32000, Israel

andreyz@campus.technion.ac.il, vadim.indelman@technion.ac.il

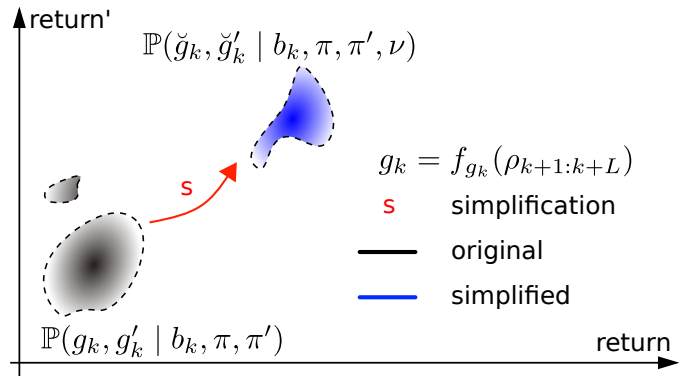
**Abstract**—It is a long-standing objective to ease the computation burden incurred by the decision making process. Identification of this mechanism’s sensitivity to simplification has tremendous ramifications. Yet, algorithms for decision making under uncertainty usually lean on approximations or heuristics without quantifying their effect. Therefore, challenging scenarios could severely impair the performance of such methods. In this paper, we extend the decision making mechanism to the whole by removing standard approximations and considering all previously suppressed stochastic sources of variability. On top of this extension, our key contribution is a novel framework to simplify decision making while assessing and controlling online the simplification’s impact. Furthermore, we present novel stochastic bounds on the return and characterize online the effect of simplification using this framework on a particular simplification technique - reducing the number of samples in belief representation for planning. Finally, we verify the advantages of our approach through extensive simulations.

## I. INTRODUCTION

Autonomous online decision making is a fundamental aspect of intelligence. In a partially observable setting, which is common in real world scenarios, there is no direct access to the state. Instead, the robot has to maintain a belief over the state, and reason about its evolution while accounting for different sources of uncertainty within the decision making stage. The renowned framework to do so is the Partially Observable Markov Decision Process (POMDP) [18]. A crucial element defining the robot’s behavior is the reward operator.

Solving a POMDP, i.e., calculating the “right decision” in terms of an optimal action sequence or policy, involves anticipating every imaginable turn of future events and computing the *returns* based on the corresponding rewards. One common example of the return is the future cumulative reward. The abundance of such possibilities, represented by a belief tree, induces a *distribution over return* for every possible action plan. Therefore, choosing an optimal action (policy) is exceptionally computationally demanding [25]. Since a direct comparison of the distributions is not possible, a decision-maker must project the distribution of return, per a possible reactive future action sequence (policy) [21], on some comparable space. Examples for such a projection are the expectation operator and risk-aware measures [6],[5],[36], such as conditional-value-at-risk (CVaR) [26].

There is a large body of algorithms to approximate decision making under uncertainty. Classical offline methods are based on  $\alpha$ -vectors [24] or value iteration [2]. More



**Fig. 1:** This figure shows alteration of the distribution of joint returns  $g_k$  and  $g'_k$  of two candidate policies  $\pi$  and  $\pi'$  as a result of simplification. Color intensity denotes distribution values. This is a conceptual illustration, i.e., we do not imply higher/lower rewards or change of support due to simplification.

recently, online methods became successful. These include, for example, POMCP [28] and its various extensions (e.g., [32]), an algorithm designed for large POMDP and based on Monte Carlo tree search. Another popular algorithm, DESPOT [30] [37], focuses on the set of randomly sampled scenarios over the belief tree, avoiding drawbacks of the UCT [22] algorithm used in POMCP. Presently, it has been further improved to tackle high dimensional spaces [15]. Kurniawati and Yadav [23] present an interesting adaptive algorithm for dynamic environments.

Standard POMDP formulations consider state-dependent rewards. POMDP with *belief-dependent rewards* received much less attention, although these rewards are essential in numerous problems, such as information gathering, autonomous navigation, and active sensing. Information theoretic rewards are especially significant for belief space planning (BSP) [17], [14]. Araya et al. [1] introduced  $\rho$ -POMDP and extended the exact  $\alpha$ -vectors method and a family of point based approximation algorithms to considering convex belief-dependent reward functions. Later [13] extended their work further to Lipschitz reward functions. Spaan et al. [31] proposed to augment action space with information-reward actions. Dressel and Kochenderfer [7] proposed an extension of SARSOP [24] to specific forms of belief-dependent rewards.

Previous techniques necessitate specific forms of reward operators. One standing-out approach [11], [12] presents incremental reuse of calculations between different planning sessions. Notably, that approach is formulated for general belief-dependent rewards.

We take a different path, which is to simplify the original decision making problem. In other words, instead of approximating the problem, we substitute it with a simpler one. If the order of policies with respect to the original and simplified problems' objective is preserved, such substitution does not affect the decision making quality. Moreover, if, utilizing the simplified problem, we can find online bounds over the returns or objective function of the original problem, it is possible to account for the simplification loss. Replacement of various parts of the decision making problem to ease the computation burden while preserving the precedence of objectives for potential action plans recently appeared in the literature under the name action consistency [10],[9],[19]. Yet, these works considered a limited setting of a specific projection operator, Gaussian distributions, and maximum likelihood observations. Importantly, they formulated action consistency in an absolute way. However, an absolute action consistent simplification is not easy to achieve in complex general scenarios: to find such action consistent simplification, one has to preserve action trends for all conceivable realizations of the future for every potential sequence of actions, even if considering a specific projection operator (expectation). Unfortunately, this is exceptionally challenging without the assumption of maximum likely observation.

In this paper, we introduce a more general framework that allows to reason leniently about simplifications that mostly, or partially, preserve action ordering, i.e., for some of the possible future return realizations of different actions. Given a user-provided threshold on the loss incurred by simplification, one can assess the probability of suffering loss larger than the threshold and thus provide performance guarantees.

We focus on the distribution of the returns. This distribution conveys all the information about the decision making problem. Our goal is to examine how the simplification method influences the performance of the decision-maker. The simplification impacts the joint distribution of the returns, as illustrated in Fig. 1. A simplification technique could affect the dependency between returns marginals [16], as well as the marginals themselves [9].

Ultimately, we shall perform our analysis in an online setting, meaning, without accessing the original problem. In other words, we are permitted to use only the ingredients of the simplified problem to quantify the simplification effect. Therefore, we intend to quantify distribution over loss online.

Furthermore, our study's center is belief-dependent general rewards. In the most general formulation, reward calculation and belief update are stochastic methods. Typically the sources of the stochasticity are sample approximations. To our knowledge, all existing works adopting the POMDP setting do not consider these stochastic aspects. Conventionally the reward approximation and the belief update are considered determin-

istic. We relax this assumption to account for all distortions contracted due to simplification. Note that, stochasticity of the reward operator naturally cancels the assumptions of convexity or Lipschitz continuity.

To summarize, our key contributions are as follows. (a) We extend  $\rho$ -POMDP to probabilistic  $\rho$ -POMDP ( $\mathbb{P}\rho$ -POMDP) by relaxing the assumption that the reward operator and the belief update are deterministic; (b) We introduce the general concept of  $\mathbb{P}\text{Loss}$ ; (c) We provide an online characterization of the  $\mathbb{P}\text{Loss}$ ; (d) Finally, we exemplify our framework on a particular simplification technique, which is reducing the number of samples for planning.

## II. NOTATIONS AND PROBLEM FORMULATION

Let us denote by  $\mathbb{P}$  the probability density function and by  $P$  the probability. By lowercase letter we denote a random vector or its realization. For two random variables  $x$  and  $y$ , we say that they are equal  $x = y$  if they are equal as functions on their measurable space. Further, to shorten notations, we shall often use  $\square_{k+}$  to denote  $\square_{k+1:k+L}$ , where  $L$  is the planning horizon. By  $\equiv$  we denote identity.

### A. $\rho$ -POMDP

Let  $k$  be an arbitrary time step.  $\rho$ -POMDP [1] is a tuple

$$\langle \mathcal{X}, \mathcal{A}, \mathcal{Z}, T, O, \rho, \gamma, b_0 \rangle, \quad (1)$$

where  $\mathcal{X}, \mathcal{A}, \mathcal{Z}$  are state, action, and observation spaces with  $x_k \in \mathcal{X}, a_k \in \mathcal{A}, z_k \in \mathcal{Z}$  the momentary state, action, and observation, respectively,  $T(x_k, a_k, x_{k+1}) = \mathbb{P}_T(x_{k+1}|a_k, x_k)$  is the stochastic transition model from the past momentary state  $x_k$  to the next  $x_{k+1}$  through action  $a_k$ ,  $O(z_k, x_k) = \mathbb{P}_Z(z_k|x_k)$  is the stochastic observation model,  $\rho(b_{k+1}, a_k)$  is a scalar belief and previous action dependent reward operator,  $\gamma \in [0, 1]$  is the discount factor, and  $b_0$  is the belief about the initial state (prior). Throughout this paper we assume that  $\gamma = 1$ .

### B. Belief Space Planning

The posterior belief at time instant  $k$  is given by

$$b_k(x_k) \approx \mathbb{P}(x_k|b_0, a_{0:k-1}, z_{1:k}), \quad (2)$$

The belief is an efficient way of storing all relevant information that is obtainable so far. The usual assumption is that the belief is a sufficient statistic for decision making objective [3]. However, in practice, the belief requires some representation. In general, this representation is not perfect, e.g., parametric or sampled form; thus, in (2), we used the  $\approx$  sign. In a real life scenario

$$b_k = \psi(\psi(\dots \psi(b_0, a_0, z_1), a_{k-2}, z_{k-1}), a_{k-1}, z_k), \quad (3)$$

where  $\psi$  is a method for updating the belief.

Denote by  $\pi_\ell$  policy at time step  $\ell$  such that  $\pi_\ell(b_\ell) = a_\ell$  maps belief to the action. It is noteworthy that policy  $\pi(b)$  is a random function of the belief in general. For simplicity we assume that policy is deterministic. However, our development is not constrained to deterministic policies. By  $\pi \triangleq \pi_{k:k+L-1}$

we denote a vector of policies for  $L$  time steps starting from time step  $k$ .

To behave optimally, the robot shall choose a policy maximizing the objective, which in its general form, is

$$J^L(b_k, \pi) = \varphi \left( \mathbb{P}(\rho_{k+1:k+L} | b_k, \pi_{k:k+L-1}), g_k \right) \quad (4)$$

s.t.  $b_\ell = \psi(b_{\ell-1}, \pi_{\ell-1}(b_{\ell-1}), z_\ell)$ ,

where  $L$  is a planning horizon,  $\rho_\ell$  is a random reward,  $\varphi$  is a projection operator, and  $g_k \triangleq f_{g_k}(\rho_{k+})$  is the return [33]. The return is some known function of the realization of  $\rho_{k+1:k+L}$ ; as discussed in [6], e.g., it could correspond to the cumulative reward  $g_k = \sum_{\ell=1}^L \rho_{k+\ell}$ . If (4) admits Bellman form, it can be written as

$$J^L(b_k, \pi_{k:k+L-1}) = \quad (5)$$

$$J^1(b_k, \pi_k) + \int_{b_{k+1}} \mathbb{P}(b_{k+1} | b_k, \pi_k) J^{L-1}(b_{k+1}, \pi_{k+1:k+L-1})$$

s.t.  $b_\ell = \psi(b_{\ell-1}, \pi_{\ell-1}(b_{\ell-1}), z_\ell)$ .

For example, a common choice for  $\varphi$  is expectation over the distribution of future rewards given all data available [6]. However, our formulation considers a general projection operator.

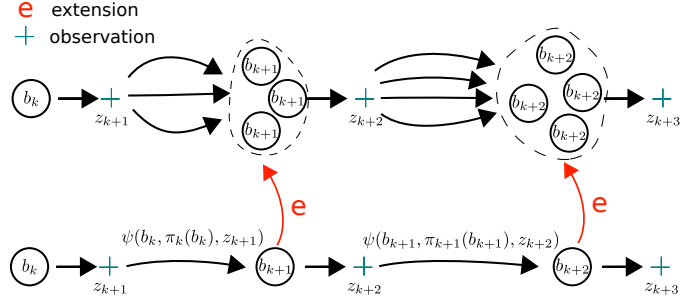
$\psi$  is a general method for propagating the belief with action and updating it with the received observation. Sometimes the belief  $b_{\ell-1}$  has a simple parametric form  $\theta_{\ell-1}$ , where  $\theta_{\ell-1}$  is the vector of parameters, e.g., a Gaussian belief. In this case, belief update  $\psi$  can be deterministic, and is denoted by  $\psi_{dt}(\theta_{\ell-1}, \pi_{\ell-1}(\theta_{\ell-1}), z_\ell)$ , where the subscript  $dt$  stands for deterministic. In more general and challenging scenarios the belief  $b_{\ell-1}$  is given by a set of weighted samples  $\{(w_{\ell-1}^i, x_{\ell-1}^i)\}_{i=1}^N$ . Therefore,  $\psi$  is a stochastic method, e.g., a particle filter [34]. Applying multiple times  $\psi$  on the same input will yield different sets of samples approximating the same distribution of the posterior belief. We denote the stochastic  $\psi$  by  $\psi_{st}(b_{\ell-1}, \pi_{\ell-1}(b_{\ell-1}), z_\ell)$ . Thus,  $\psi_{st}$  is a random function of the previous belief, an action and the observation. Note also another common situation where  $b_{\ell-1}$  is parameterized, but there is no closed form update. In this case,  $\psi$  is also a stochastic method. Another form to formulate the above is that the distribution

$$B(b_{\ell-1}, \pi_{\ell-1}(b_{\ell-1}), z_\ell, b_\ell) \triangleq \mathbb{P}_B(b_\ell | b_{\ell-1}, \pi_{\ell-1}, z_\ell), \quad (6)$$

is not a Dirac delta function. This aspect was disregarded so far, to the best of our knowledge. Note that in a belief-MDP formulation, the assumption is that  $B$  is a Dirac delta function. We emphasize relation to belief-MDP in Appendix A.

Similar arguments also hold for the momentary reward operator of the belief and the previous action. In its pure theoretical form, the momentary reward is a deterministic operator of the posterior belief and possibly an action. For example, a common immediate reward is of the form

$$\rho_{dt}(b) = \mathbb{E}_{x \sim b} [f(b(x), x)] = \int_x b(x) f(b(x), x) dx, \quad (7)$$



**Fig. 2:** Illustration of one branch of the extended belief tree. In a conventional setting (bottom), under the policy  $\pi$ , a specific realization of observations  $z_{k+1:k+3}$  defines the beliefs along the way. In our extended setting (top), that is not the case, as discussed in text. It is customary to choose the same beliefs used to build the tree to obtain reward distribution or samples from the reward. We decoupled beliefs from the tree and beliefs from the reward calculation. By the red arrow, we denote our extension (red e).

where usually  $f(b(x), x) = -\log b(x)$  or some reward on the state  $f(b(x), x) = r(x)$ , producing differential entropy or mean distance to goal. Unfortunately, an analytical expression for the reward operator  $\rho_{dt}(\cdot)$  is available in only limited scenarios, e.g., if the belief is modeled as Gaussian and the reward is differential entropy. The representation of the beliefs in (6) dictates practical reward operators. Sometimes the deterministic operator can be constructed on top of a particular belief representation. E.g., (6) outputs a set of weighted samples and (7) is adapted to be a deterministic operator of this output [4]. However, it is not always possible. In extremely challenging situations the reward includes modification of the representation of the belief. This could introduce an additional source of stochasticity. We extend (7) to

$$R(b_\ell, \pi_{\ell-1}(b_{\ell-1}), \rho_\ell) \triangleq \mathbb{P}_R(\rho_\ell | b_\ell, \pi_{\ell-1}(b_{\ell-1})), \quad (8)$$

embracing these possibilities. To our knowledge, we are the first who treat these aspects as random.

### C. Problem Formulation

Our goal is to generalize the concept of absolute action consistency to probabilistic, more lenient. Given some simplification method, we want to affirm *online* what is the loss and provide probabilistic guarantees. Our approach to analysis and reasoning shall be agnostic to the choice of the projection operator  $\varphi$ .

## III. APPROACH

### A. Probabilistic $\rho$ -POMDP

To capture the complete simplification impact, we shall account for all potential sources of variability. We remove conventional approximations by extending (1) to a probabilistic reward model  $R$  (8) and probabilistic belief update  $B$  (6), and introduce

$$M = \langle \mathcal{X}, \mathcal{A}, \mathcal{Z}, T, O, R, \gamma, b_k, B \rangle, \quad (9)$$

which we name probabilistic  $\rho$ -POMDP ( $\mathbb{P}\rho$ -POMDP). The rationale behind these conditional distributions ( $R$  and  $B$ ) is to capture additional sources of stochasticity, such as stochastic belief update, stochastic calculation of a given reward operator or simply not knowing the operator reward in explicit analytic form. These previously overlooked sources of stochasticity impact the likelihood of the observations

$$\mathbb{P}(z_{k+1:k+L}|b_k, \pi), \quad (10)$$

as well as the joint reward distribution  $\mathbb{P}(\rho_{k+}|b_k, \pi, z_{k+}) \equiv \mathbb{P}(\rho_{k+1:k+L}|b_k, \pi_{k:k+L-1}, z_{k+1:k+L})$  given a realization of future observations. The latter can be factorized as

$$\prod_{\ell=k+1}^{k+L} \int_{b_\ell} \mathbb{P}_R(\rho_\ell|b_\ell, \pi_{\ell-1}) \mathbb{P}_B(b_\ell|b_{\ell-1}, \pi_{\ell-1}, z_\ell), \quad (11)$$

which is Dirac's delta function in the regular setting of POMDP and  $\rho$ -POMDP. If  $B$  is a Dirac function, a sample from (10) uniquely defines the corresponding posterior beliefs  $b_{k+1:k+L}$ . This, therefore, corresponds to the classical belief tree. In contrast, our  $\mathbb{P}\rho$ -POMDP (9), corresponds to an *extended* belief tree, which, due to (6), allows many samples of the beliefs  $b_{k+1:k+L}$  for each sample of  $z_{k+1:k+L}$  from (10). We illustrate this in Fig. 2.

### B. Simplification Formulation

To formally define the simplification procedure, we augment the  $\mathbb{P}\rho$ -POMDP tuple (9) with a simplification operator  $\nu$ ,

$$M_\nu = \langle \mathcal{S}, \mathcal{A}, \mathcal{Z}, T, O, R, \gamma, b_k, B, \nu \rangle, \quad \nu \triangleq \nu_k, \dots, \nu_{k+L}. \quad (12)$$

This general operator defines any possible modification of the original problem defined by (9) alongside with (4) to a new, simpler to solve, problem. The operator  $\nu$  can be for example, sparsification of the initial belief  $b_k$  [9], substitution of the operator differential entropy by a simpler operator, e.g., trace of covariance matrix, discarding the normalizer in the differential entropy operator [27], replacing the reward by its topological signature [19].

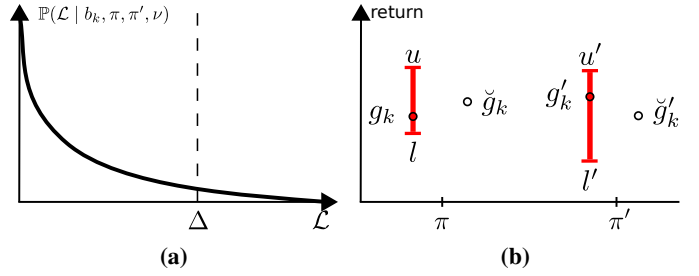
To distinct a simplified reward from the original reward, we denote the former by  $\check{\rho}$  instead of  $\rho$ ; similarly, we denote the simplified belief by  $\check{b}$  instead of  $b$ . Note the operator  $\nu$  can be stochastic, as discussed below.

Specifically, belief simplification is described by the distribution

$$\mathbb{P}(\check{b}_\ell|b_\ell; \nu_\ell^b). \quad (13)$$

In general, the distribution (13) over the simplified belief  $\check{b}_\ell$  corresponds to a stochastic simplification operator  $\nu_\ell^b$ . This is the case, for example, when  $b_\ell$  is represented by a set of  $N$  weighted samples and  $\nu_\ell^b$  is the operation of subsampling  $n$  samples according to weights; i.e., applying this operation on  $b_\ell$  multiple times leads to different sets of  $n$  samples, each representing another realization of  $\check{b}_\ell$  from (13). For a deterministic operator  $\nu_\ell^b$ , (13) is a Dirac function.

Further, there are several cases of how a simplification affects belief update (6) from time  $\ell - 1$  to  $\ell$ .



**Fig. 3:** Illustration of (a) the distribution of loss, and (b) the online bounds of the return.

- 1) Without any simplification we have  $\mathbb{P}_B(b_\ell|b_{\ell-1}, \pi_{\ell-1}, z_\ell)$  from (6).
- 2) Given a simplified belief  $\check{b}_{\ell-1}$ , while keeping the original stochastic belief update  $\psi_{st}$ , we have  $\mathbb{P}_B(\check{b}_\ell|\check{b}_{\ell-1}, \pi_{\ell-1}, z_\ell)$ , where each realization of  $\check{b}_\ell$  is obtained via  $\psi_{st}$ . Thus, given  $\check{b}_{\ell-1}$ , this distribution is not a function of  $\nu$ .
- 3) We can also simplify the belief update operator,  $\psi_{st}$ , to  $\check{\psi}_{st}$ . Denoting the corresponding simplification operator  $\nu_\ell^\psi$ , this yields  $\mathbb{P}_{\check{B}}(\check{b}_\ell|\check{b}_{\ell-1}, \pi_{\ell-1}, z_\ell; \nu_\ell^\psi)$ .
- 4) Finally, one can decide at time  $\ell$  to apply simplification on the belief (determined by  $\nu_\ell^b$ ) via (13). The corresponding belief update can be written as  $\mathbb{P}_{\check{B}}(\check{b}_\ell|\check{b}_{\ell-1}, \pi_{\ell-1}, z_\ell; \nu_\ell^b, \nu_\ell^\psi) = \int_{\check{b}_\ell} \mathbb{P}(\check{b}_\ell|\check{b}_\ell; \nu_\ell^b) \mathbb{P}_{\check{B}}(\check{b}_\ell|\check{b}_{\ell-1}, \pi_{\ell-1}, z_\ell; \nu_\ell^\psi)$ , where  $\check{b}_\ell$  is the integration variable.

We combine these cases and write

$$\check{B}(\check{b}_{\ell-1}, \pi_{\ell-1}, z_\ell, \check{b}_\ell; \nu) \triangleq \mathbb{P}_{\check{B}}(\check{b}_\ell|\check{b}_{\ell-1}, \pi_{\ell-1}, z_\ell; \nu_\ell^b, \nu_\ell^\psi). \quad (14)$$

Similarly, reward simplification could be, in general, stochastic, leading to the distribution

$$\mathbb{P}(\check{\rho}_\ell|\rho_\ell; \nu_\ell^\rho). \quad (15)$$

Thus, given a simplified belief  $\check{b}_\ell$  and  $\check{b}_{\ell-1}$ , and recalling (8), the distribution over  $\check{\rho}_\ell$  is

$$\mathbb{P}_{\check{R}}(\check{\rho}_\ell|\check{b}_\ell, \pi_{\ell-1}(\check{b}_{\ell-1}); \nu) = \int_{\check{\rho}_\ell} \mathbb{P}(\check{\rho}_\ell|\check{\rho}_\ell; \nu_\ell^\rho) \mathbb{P}_R(\check{\rho}_\ell|\check{b}_\ell, \pi_{\ell-1}(\check{b}_{\ell-1})),$$

which we denote as the simplified reward model,

$$\check{R}(\check{b}_\ell, \pi_{\ell-1}(\check{b}_{\ell-1}), \check{\rho}_\ell; \nu) \triangleq \mathbb{P}_{\check{R}}(\check{\rho}_\ell|\check{b}_\ell, \pi_{\ell-1}(\check{b}_{\ell-1}); \nu). \quad (16)$$

Consequently, the models (14) and (16) impact (11), and lead to the following simplified joint reward distribution given a realization of future observations

$$\mathbb{P}(\check{\rho}_{k+}|b_k, \pi, z_{k+}, \nu) = \int_{\check{b}_k} \mathbb{P}(\check{b}_k|b_k; \nu_k^b) \prod_{\ell=k+1}^{k+L} \int_{\check{b}_\ell} \mathbb{P}_{\check{R}}(\check{\rho}_\ell|\check{b}_\ell, \pi_{\ell-1}; \nu) \mathbb{P}_{\check{B}}(\check{b}_\ell|\check{b}_{\ell-1}, \pi_{\ell-1}, z_\ell; \nu). \quad (17)$$

The formulations of (14) and (16) assume that inputs to  $\check{B}$  and  $\check{R}$  are most simplified versions at appropriate time instant  $(\check{\rho}_\ell, \check{b}_\ell, \check{b}_{\ell-1})$ . This simplification approach uses only

the observations from the belief tree. In the sequel, we explain why it is advantageous. In this setting, in addition to maintaining/updating  $b_\ell$ , we also have to maintain/update the simplified version  $\check{b}_\ell$ .

Alternatively, the update of simplified belief  $\check{b}_\ell$  can be avoided. Such simplification builds upon samples from  $B$ , which are already present at the belief tree. Thus,

$$\mathbb{P}(\check{\rho}_{k+}|b_k, \pi, z_{k+}, \nu) = \int_{\check{b}_k} \mathbb{P}(\check{b}_k|b_k; \nu_k^b) \prod_{\ell=k+1}^{k+L} \int_{b_\ell} \int_{\check{b}_\ell} \quad (18)$$

$$\mathbb{P}_{\check{R}}(\check{\rho}_\ell|\check{b}_\ell, \pi_{\ell-1}(\check{b}_{\ell-1}); \nu) \mathbb{P}(\check{b}_\ell|b_\ell; \nu) \mathbb{P}_B(b_\ell|b_{\ell-1}, \pi_{\ell-1}, z_\ell).$$

To approximate (18) one can use original beliefs from the extended belief tree or sample again if needed.

In other words, in this paper we assume that operator  $\nu$  affects exclusively (11). However, the measurements are sampled as in the original problem as in (10).

### C. Probabilistic Loss ( $P_{Loss}$ )

In the previous section, we defined a simplification procedure that results in a corresponding new decision making problem that should be easier to solve. Now we stipulate on the quality of the simplification for two candidate policies  $\pi$  and  $\pi'$ .

From  $\mathbb{P}(\rho_{k+}|b_k, \pi, z_{k+}, \nu)$  and  $\mathbb{P}(\check{\rho}_{k+}|b_k, \pi, z_{k+}, \nu)$  we arrive at the distribution of the original as well as simplified returns  $\mathbb{P}(g_k|b_k, \pi)$  and  $\mathbb{P}(\check{g}_k|b_k, \pi, \nu)$  for each evaluated candidate policy. To quantify the impact of the simplification procedure, we shall consider the *joint* distribution  $\mathbb{P}(g_k, g'_k, \check{g}_k, \check{g}'_k|b_k, \pi, \pi', \nu)$ . Our goal is to examine how the simplification procedure alters the joint distribution  $\mathbb{P}(g_k, g'_k|b_k, \pi, \pi')$  towards  $\mathbb{P}(\check{g}_k, \check{g}'_k|b_k, \pi, \pi', \nu)$ . These two marginal distributions are illustrated in Fig. 1.

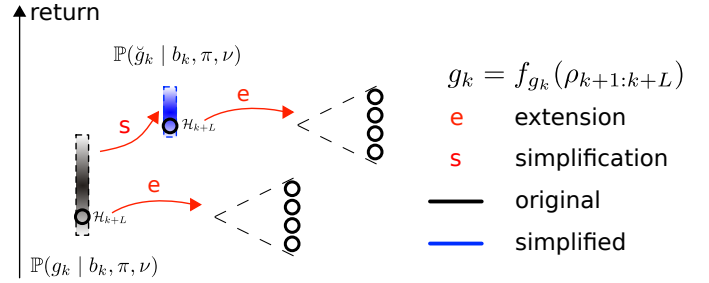
Let us define the following random variable, which we shall refer to as "loss"

$$\mathcal{L} \triangleq f_{\mathcal{L}}(g_k, g'_k, \check{g}_k, \check{g}'_k) = \begin{cases} \max\{g'_k - g_k, 0\} & \text{if } \check{g}_k - \check{g}'_k > 0, \\ \max\{g_k - g'_k, 0\} & \text{if } \check{g}_k - \check{g}'_k < 0, \\ 0 & \text{else.} \end{cases} \quad (19)$$

With (19) we aim to capture a complete impact of a simplification onto the decision making problem. Specifically, this definition captures for each possible realization of  $g_k, g'_k, \check{g}_k, \check{g}'_k$  the absolute difference between the original returns  $\Delta = |g'_k - g_k|$  in case action trend was not preserved on this realization. Meaning, at this realization, the optimal actions of original and simplified problems would differ. Given a sample  $(g_k, g'_k, \check{g}_k, \check{g}'_k)$ , the simplification is action consistent at this sample if the sign of the difference of the returns is preserved. In other words, the same action would be identified as optimal with the original and simplified returns; else we must account for the loss (19).

Our object of interest is the distribution density of  $\mathcal{L}$  given all the information available at our disposal,

$$\mathbb{P}(\mathcal{L}|b_k, \pi, \pi', \nu). \quad (20)$$



**Fig. 4:** Our extended setting permits variability of the reward given the present and a realization of the future. On the contrary, in a conventional setting, (28) is always a Dirac delta function.

We denote this distribution by Probabilistic Loss ( $P_{Loss}$ ), as it generalizes the concept of absolute action consistency to probabilistic. See illustration in Fig. 3a. E.g., if (20) is the Dirac delta function  $\delta(\mathcal{L})$ , the simplification method is absolute action consistent for every possible operator projection  $\varphi$ .

Moreover, for any  $\Delta$ , its cumulative distribution function (CDF)  $\mathbb{P}(\mathcal{L} \leq \Delta|b_k, \pi, \pi', \nu)$  provides probability to suffer loss at most  $\Delta$ . Similarly, the tail distribution function (TDF)  $\mathbb{P}(\mathcal{L} > \Delta|b_k, \pi, \pi', \nu)$  provides probability to suffer loss greater than  $\Delta$ . We shall revisit and discuss these aspects further in Section III-G.

### D. Decomposition of Returns

The source of distribution (20) is

$$\mathbb{P}(g_k, g'_k, \check{g}_k, \check{g}'_k|b_k, \pi, \pi', \nu), \quad (21)$$

i.e., the joint distribution over original and simplified returns of both policies. This distribution decomposes via marginalization over future observations  $z_{k+} \equiv z_{k+1:k+L}$  and  $z'_{k+} \equiv z'_{k+1:k+L}$  as

$$\int_{z_{k+}} \int_{z'_{k+}} \mathbb{P}(g_k, g'_k, \check{g}_k, \check{g}'_k|b_k, \pi, \pi', \nu, z_{k+}, z'_{k+}) \cdot \mathbb{P}(z_{k+}, z'_{k+}|b_k, \pi, \pi') dz_{k+} dz'_{k+}, \quad (22)$$

which, according to (6), (8) and (14)-(16), decomposes to

$$\int_{z_{k+}} \int_{z'_{k+}} \mathbb{P}(g_k, \check{g}_k|\mathcal{H}_{k+L}, \nu) \mathbb{P}(g'_k, \check{g}'_k|\mathcal{H}'_{k+L}, \nu) \cdot \mathbb{P}(z_{k+}, z'_{k+}|b_k, \pi, \pi') dz_{k+} dz'_{k+}, \quad (23)$$

where  $\mathcal{H}_{k+L} \triangleq \{b_k, \pi, z_{k+}\}$  and  $\mathcal{H}'_{k+L} \triangleq \{b_k, \pi', z'_{k+}\}$ . Note that the belief  $b_k$  is shared by both histories.

In other words, the simplification operator  $\nu$  independently affects each realization of the future. Given two such realizations  $(\mathcal{H}_{k+L}, \mathcal{H}'_{k+L}, \nu)$ , the pairs of original and simplified returns are statistically independent of all other rewards. This crucial observation will be significant in the sequel.

### E. Online Bound on Probabilistic Loss (PbLoss)

The distribution defined by (20) requires access to (21) which we do not have in an online setting. To circumvent the requirement of accessing  $g_k$  and  $g'_k$ , we propose to substitute them by online lower and upper bounds  $l, u$  and  $l', u'$ , respectively. These bounds should be accessible without knowledge of original returns.

Let us consider a sampled return realization  $(g_k, g'_k, \check{g}_k, \check{g}'_k) \sim \mathbb{P}(g_k, g'_k, \check{g}_k, \check{g}'_k | b_k, \pi, \pi', \nu)$  from (21). As in an online setting we do not actually have access to the original returns  $(g_k, g'_k)$ , we strive to bound the latter,

$$l \leq g_k \leq u, \quad l' \leq g'_k \leq u', \quad (24)$$

where, for now, we assume (24) holds for any sample of  $(g_k, g'_k, \check{g}_k, \check{g}'_k)$ ; for example, these could be analytically-derived bounds. This setting is illustrated in Fig. 3b. However, in Section III-F we also discuss a more general setting where we allow (24) to be violated with probability larger than zero.

Using these bounds we are able to define online a bound on loss (19) *without* accessing the original problem ( $R$  and  $B$ ),

$$\bar{\mathcal{L}} \triangleq f_{\bar{\mathcal{L}}}(g_k, l, u, \check{g}'_k, l', u') = \begin{cases} \max\{u' - l, 0\} & \text{if } \check{g}_k - \check{g}'_k > 0, \\ \max\{u - l', 0\} & \text{if } \check{g}_k - \check{g}'_k < 0, \\ 0 & \text{else.} \end{cases} \quad (25)$$

Note that sometimes we can find bounds over the returns by applying the same function  $f_{g_k}$  on the bounds on the momentary rewards (returns when  $L = 1$ ), e.g., in case of cumulative reward  $u = \sum_{\ell=k+1}^{k+L} u_\ell$  and  $l = \sum_{\ell=k+1}^{k+L} l_\ell$ . However, this is not always possible, e.g., if  $g_k$  deviates from the sum of momentary rewards or in the case of Bellman form (5). Sometimes it is, therefore, better to work with momentary bounds.

In an online setting, we are interested in the distribution density of  $\bar{\mathcal{L}}$ ,

$$\mathbb{P}(\bar{\mathcal{L}} | b_k, \pi, \pi', \nu), \quad (26)$$

which we denote by Probabilistic Bound on Loss (PbLoss).

As we discuss in Section III-G, PbLoss characterizes the impact of a simplification in an online setting; thus, it enables to determine online if a candidate simplification is acceptable given a user-specified criteria. The decision to either accept or decline a (candidate) simplification is guided by probabilistic guarantees, as provided by our approach.

### F. Online Stochastic Bounds

Our extension allows  $R$  and  $B$ , as well as  $\check{R}$  and  $\check{B}$  to be any distributions. They can remain Dirac functions, e.g., if belief update and the reward calculation have a closed form

$$\mathbb{P}(\rho_\ell | b_{\ell-1}, a_{\ell-1}, z_\ell) = \delta(\rho_\ell - \rho_{dt}(\psi_{dt}(\theta_{\ell-1}, a_{\ell-1}, z_\ell))). \quad (27)$$

In such a case, conditioned on  $(\mathcal{H}_{k+L}, \mathcal{H}'_{k+L}, \nu)$ , the returns  $(g_k, g'_k, \check{g}_k, \check{g}'_k)$  are deterministic.

However, in the more general case, following our extension, there is a joint distribution of original and simplified returns given a realization of the future and the present,

$$\mathbb{P}(g_k, \check{g}_k | \mathcal{H}_{k+L}, \nu), \quad (28)$$

as illustrated in Fig. 4. Since (28) is no longer a Dirac function, we can use knowledge about this distribution to design bounds, which will hold with *some* probability.

In section IV, we show that it is possible to harness the structure of (28) to design the mentioned more lenient online bounds. Our framework permits to detach the process of estimation of the bounds from the realization of the reward and truly use all accessible information in a simplified problem. For example, one way to design probabilistic bounds is to find online a random variable  $\epsilon$  such that the probability

$$\mathbb{P}(|g_k - \check{g}_k| \leq \epsilon | \mathcal{H}_{k+L}, \nu) \quad (29)$$

is bounded from below. The corresponding probabilistic lower and upper bounds will be  $l = \check{g}_k - \epsilon$  and  $u = \check{g}_k + \epsilon$ , respectively. We, therefore, refer to  $l$  and  $u$  as random variables. In our setting, even if the bounds actually bound with very low probability, it is still possible to analyze the quality of the simplification. Moreover, the analytical bounds, designed in a conventional setting, can be used in our extended setting without any revision. In our extended environment, they will bound with probability one.

Having introduced the novel stochastic bounds, we proceed to the formulation of the constraints, that these bounds shall fulfill to be meaningful. Our goal is to formulate conditions, which assure that PbLoss (26) is connected to PLoss (20) and can be used online to analyze the quality of the simplification.

The following conditional

$$\mathbb{P}(g_k, g'_k, \check{g}_k, \check{g}'_k, l, u, l', u' | \mathcal{H}_{k+L}, \mathcal{H}'_{k+L}, \nu), \quad (30)$$

encloses all the variables situated in the problem. Moreover, following Section III-D, (30) decomposes into

$$\mathbb{P}(g_k, \check{g}_k, l, u | \mathcal{H}_{k+L}, \nu) \mathbb{P}(g'_k, \check{g}'_k, l', u' | \mathcal{H}'_{k+L}, \nu). \quad (31)$$

For calculating PbLoss (26) we will need samples from  $\mathbb{P}(\check{g}_k, l, u | \mathcal{H}_{k+L}, \nu)$  and  $\mathbb{P}(\check{g}'_k, l', u' | \mathcal{H}'_{k+L}, \nu)$ . Let  $[\cdot]$  be the Iverson bracket and  $\alpha \in [0, 1)$ . For every possible sample  $(\check{g}_k, \check{g}'_k)$  we do not know which sample  $(g_k, g'_k)$  one could obtain in the original problem. However, if the bounds are designed such that

$$\mathbb{P}(g_k, l, u | \mathcal{H}_{k+L}, \nu) \quad , \quad \mathbb{P}(g'_k, l', u' | \mathcal{H}'_{k+L}, \nu) \quad (32)$$

render

$$(1 - \alpha) \leq P([l \leq g_k \leq u] = 1 | \mathcal{H}_{k+L}, \nu) \quad (33)$$

and

$$(1 - \alpha) \leq P([l' \leq g'_k \leq u'] = 1 | \mathcal{H}'_{k+L}, \nu), \quad (34)$$

we can bound online CDF and TDF of PLoss using PbLoss, as we show in Section III-G.

---

**Algorithm 1** Online characterization of the simplification

---

**Input:** Two candidate policies  $\pi, \pi'$ . Initial belief  $b_k$ . Samplers from  $\mathbb{P}(\check{g}_k, l, u | \mathcal{H}_{k+L}, \nu)$  and  $\mathbb{P}(\check{g}'_k, l', u' | \mathcal{H}'_{k+L}, \nu)$ . Sample  $b_k$  or take the initial samples from inference. Obtain  $\mathbb{P}(z_{k+1:k+L}, z'_{k+1:k+L} | b_k, \pi, \pi')$  and create two belief policy trees.

**for all** sample pairs  $(z_{k+1:k+L}, z'_{k+1:k+L})$  **do**

Obtain sample  $(\check{g}_k, l, u, \check{g}'_k, l', u')$ .

Calculate  $f_{\bar{\mathcal{L}}}(\check{g}_k, l, u, \check{g}'_k, l', u')$  according to (25).

**end for**

$\{f_{\bar{\mathcal{L}}}(\check{g}_k, l, u, \check{g}'_k, l', u')\}$  represents the set of samples of  $\bar{\mathcal{L}}$ .

**Output:**  $\forall \Delta$  empirically calculated  $P(\bar{\mathcal{L}} > \Delta | b_k, \pi, \pi', \nu)$

as  $\frac{\text{number of samples of } \bar{\mathcal{L}} \text{ satisfying } \bar{\mathcal{L}} > \Delta}{\text{number of all samples of } \bar{\mathcal{L}}}$ .

---

We note that, in general, (33) and (34) could each have its own  $\alpha$ . All developments and proofs can be adjusted easily to this setting.

Considering the above,  $\text{PbLoss}$  (26) is based on

$$\mathbb{P}(\check{g}_k, l, u, \check{g}'_k, l', u' | b_k, \pi, \pi', \nu). \quad (35)$$

To summarize, there are three types of online reward bounds:

- 1) Deterministic bounds. These analytical bounds exist in case of a closed form belief update  $\psi_{dt}$  and a deterministic operator reward  $\rho_{dt}(b)$  from 7, e.g., belief is a Gaussian and the reward is differential entropy. In this case, even in our extended setting  $R$  and  $B$  remain Dirac functions.
- 2) Stochastic bounds that hold with probability one. These are also analytical bounds. In our extended setting  $R$  and  $B$  are no longer Dirac functions. However, these bounds hold for any realization of sample approx., as stated around (24).
- 3) Stochastic bounds that hold at least with probability  $1 - \alpha$ . They exist only in our extended setting when  $R$  and  $B$  are not Dirac functions.

### G. Characterization of $\text{P}_{\text{Loss}}$ Online

In this section, we show how  $\text{PbLoss}$  can be used in an online setting to characterize  $\text{P}_{\text{Loss}}$  (which is unavailable online). In turn, this enables to provide online probabilistic performance guarantees for a considered simplification (represented by operator  $\nu$ ), or to decide if it is adequate given a user-specified criteria.

Specifically, recall  $\text{P}_{\text{Loss}}$  CDF and TDF, i.e., probability to suffer loss at most, or greater, than  $\Delta \in \mathbb{R}$ , respectively,

$$\text{P}_{\text{Loss}} \text{ CDF: } P(\mathcal{L} \leq \Delta | b_k, \pi, \pi', \nu) \quad (36)$$

$$\text{P}_{\text{Loss}} \text{ TDF: } P(\mathcal{L} > \Delta | b_k, \pi, \pi', \nu). \quad (37)$$

We now aim to bound  $\text{P}_{\text{Loss}}$  CDF (36) from below, and  $\text{P}_{\text{Loss}}$  TDF (37) from above by utilizing  $\text{PbLoss}$ .

We now consider  $\text{P}_{\text{Loss}}$  TDF and express  $P(\mathcal{L} > \Delta | b_k, \pi, \pi', \nu)$  as

$$P(\mathcal{L} > \Delta, \bar{\mathcal{L}} \geq \mathcal{L} | b_k, \pi, \pi', \nu) + P(\mathcal{L} > \Delta, \bar{\mathcal{L}} < \mathcal{L} | b_k, \pi, \pi', \nu).$$

The first term can be written via chain rule as

$$P(\mathcal{L} > \Delta | \bar{\mathcal{L}} \geq \mathcal{L}, b_k, \pi, \pi', \nu) P(\bar{\mathcal{L}} \geq \mathcal{L} | b_k, \pi, \pi', \nu). \quad (38)$$

Performing chain rule similarly also on the second term and recalling that  $P(\bar{\mathcal{L}} \geq \mathcal{L} | \cdot) + P(\bar{\mathcal{L}} < \mathcal{L} | \cdot) = 1$ , allows to express  $\text{P}_{\text{Loss}}$  TDF as

$$P(\mathcal{L} > \Delta | b_k, \pi, \pi', \nu) = P(\mathcal{L} > \Delta | \bar{\mathcal{L}} \geq \mathcal{L}, b_k, \pi, \pi', \nu) \lambda + P(\mathcal{L} > \Delta | \bar{\mathcal{L}} < \mathcal{L}, b_k, \pi, \pi', \nu) (1 - \lambda), \quad (39)$$

where

$$\lambda \triangleq P(\bar{\mathcal{L}} \geq \mathcal{L} | b_k, \pi, \pi', \nu) \equiv P([\bar{\mathcal{L}} \geq \mathcal{L}] = 1 | b_k, \pi, \pi', \nu). \quad (40)$$

While  $\lambda$  from (40) is unavailable, we can bound it from below using (33) and (34) as follows.

**Theorem 1** (Probability that bound bounds). *Fix  $\alpha \in \mathbb{R}$ . Assume that (33) and (34) hold. Then:*

$$P([\bar{\mathcal{L}} \geq \mathcal{L}] = 1 | b_k, \pi, \pi', \nu) \geq (1 - \alpha)^2. \quad (41)$$

Proof: see Appendix B.

Now we show that given the event  $\bar{\mathcal{L}} \geq \mathcal{L}$ ,  $\text{P}_{\text{Loss}}$  TDF is bounded from above by  $\text{PbLoss}$  TDF.

**Theorem 2** (Conditional TDF Lower bound).  $\forall \Delta \in \mathbb{R}$ ,

$$P(\mathcal{L} > \Delta | \bar{\mathcal{L}} \geq \mathcal{L}, b_k, \pi, \pi', \nu) \leq P(\bar{\mathcal{L}} > \Delta | \bar{\mathcal{L}} \geq \mathcal{L}, b_k, \pi, \pi', \nu).$$

Proof: see Appendix B.

Finally, we characterize  $\text{P}_{\text{Loss}}$  as follows.

**Theorem 3** (Upper and Lower bounds). *Denote  $\beta(\Delta) \triangleq \min \left\{ 1, \frac{P(\bar{\mathcal{L}} > \Delta | b_k, \pi, \pi', \nu)}{(1 - \alpha)^2} + 2\alpha - \alpha^2 \right\}$ , so*

$$P(\mathcal{L} > \Delta | b_k, \pi, \pi', \nu) \leq \beta(\Delta) \quad (42)$$

and

$$P(\mathcal{L} \leq \Delta | b_k, \pi, \pi', \nu) \geq 1 - \beta(\Delta). \quad (43)$$

Proof: see Appendix B.

We can say that a simplification procedure, in a worst case scenario, will render decision making sub optimal at most  $\Delta$  with probability at least  $1 - \beta(\Delta)$ . Moreover, since  $0 \leq \mathcal{L}$ , setting  $\Delta = 0$  in Alg. 1 we can assess the probability to be absolute action consistent in worst case scenario (for any  $\varphi$ ).

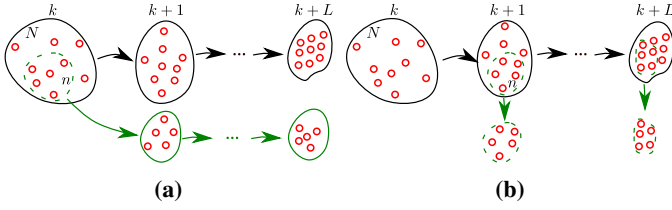
### H. Calculating $\text{P}_{\text{Loss}}$ Offline and $\text{PbLoss}$ Online

We discuss the offline calculation of the  $\text{P}_{\text{Loss}}$  in Appendix C.

So far, we did not explain how to calculate  $\text{PbLoss}$  (26). One approach is to sample  $(\check{g}_k, l, u, \check{g}'_k, l', u')$  from (35) and evaluate  $\bar{\mathcal{L}}$  for each such sample via (25). Then,  $\text{PbLoss}$  is represented by  $\{f_{\bar{\mathcal{L}}}(\check{g}_k, l, u, \check{g}'_k, l', u')\}$ .

Generating samples from (35) involves marginalizing over future measurements  $z_{k+} \equiv z_{k+1:k+L}$  and  $z'_{k+} \equiv z'_{k+1:k+L}$ . Similar to (23),  $\mathbb{P}(\check{g}_k, l, u, \check{g}'_k, l', u' | b_k, \pi, \pi', \nu)$  decomposes to

$$\int_{z'_{k+}}^{z_{k+}} \mathbb{P}(\check{g}_k, l, u | \mathcal{H}_{k+L}, \nu) \mathbb{P}(\check{g}'_k, l', u' | \mathcal{H}'_{k+L}, \nu) \cdot \mathbb{P}(z_{k+}, z'_{k+} | b_k, \pi, \pi') dz_{k+} dz'_{k+}, \quad (44)$$



**Fig. 5:** Potential simplification techniques: **(a)** Choosing a subset of samples only at time  $k$ ; **(b)** Choosing a subset of samples at each time  $\ell$ .

In practice,  $\mathbb{P}(z_{k+1}, z'_{k+1} | b_k, \pi, \pi')$  corresponds to two extended belief policy trees, starting from the same root ( $b_k$ ) and having the same rule for choosing rollouts. The specific way of obtaining samples from

$$\mathbb{P}(\check{g}_k, l, u | \mathcal{H}_{k+L}, \nu) \quad , \quad \mathbb{P}(\check{g}'_k, l', u' | \mathcal{H}'_{k+L}, \nu) \quad (45)$$

depends on the operator  $\nu$ . In the next section, we elaborate on these aspects, considering a specific simplification operator.

#### IV. SPECIFIC SIMPLIFICATION

##### A. Simplification Technique

In this section, we exemplify our technique on a specific simplification method. Assume the belief  $b_k$  is represented by a set of  $N$  weighted samples  $\{(w_k^i, x_k^i)\}_{i=1}^N$ . Our simplification operator  $\nu$  provides a way to choose a subset of  $n$  samples from the original  $N$  samples. For example, subsampling according to weights. We denote by (a) simplification affecting  $b_k$  and producing  $\check{b}_k = \{(w_k^j, x_k^j)\}_{j=1}^n$  ( $\nu$  is only applied at time  $k$  with fixed seed) (17) and by (b) simplification affecting  $b_\ell$  and producing  $\check{b}_\ell = \{(w_\ell^j, x_\ell^j)\}_{j=1}^n$  for  $\ell \in [k+1, k+L]$  (18). We take  $\psi_{st}$  as an off-the-shelf particle filter, which produces the same number of samples as the input. The reward operator is approximated by  $N$  and  $n$  samples in the original and simplified setting, respectively. To make a clear connection to our general framework, in this section, we denote  $\rho_\ell = r_\ell^N$  and  $\check{\rho}_\ell = r_\ell^n$ . The two simplification methods are illustrated in Fig. 5.

##### B. Online Bounds on Sample Based Reward

To present development for (33), we take inspiration from confidence intervals [35]. Let us introduce the following model

$$\begin{pmatrix} g_k \\ \check{g}_k \end{pmatrix} | \mathcal{H}_{k+L}, \nu \sim \mathcal{N} \left( \begin{pmatrix} \mu \\ \mu \end{pmatrix} \begin{pmatrix} \text{se}^2(N) & \text{cov} \\ \text{cov} & \text{se}^2(n) \end{pmatrix} \right), \quad (46)$$

where  $\text{se}$  is the standard error and  $\text{cov}$  is the covariance. Online we do not have access to these quantities. The standard error depends on the number of samples  $N$  and  $n$  respectively, dwindling as the number of samples increases. We assume that each marginal is distributed around the same mean value  $\mu$ . Denote  $y = g_k - \check{g}_k$ . It is known that  $y$  is a zero mean Gaussian with the following variance

$$\text{var}(y) = \text{se}^2(N) + \text{se}^2(n) - 2\text{cov}. \quad (47)$$

The derivation is in Appendix D. Let  $z = \frac{y}{\sqrt{\text{var}(y)}} \sim \mathcal{N}(0, 1)$  and

$$z_{\alpha/2} = \Phi^{-1}(1 - \alpha/2), \quad (48)$$

where  $\Phi$  is a Cumulative Distribution Function (CDF) of a standard normal variable [35] so  $\mathbb{P}(z > z_{\alpha/2}) = \alpha/2$  and

$$P(-z_{\alpha/2} \leq z \leq z_{\alpha/2}) = 1 - \alpha. \quad (49)$$

In other words

$$P(|y| \leq z_{\alpha/2} \sqrt{\text{var}(y)} | \mathcal{H}_{k+L}, \nu) = 1 - \alpha. \quad (50)$$

Using the facts  $\text{se}(N) \leq \text{se}(n)$  and  $\text{cov} \leq \text{se}(N)\text{se}(n)$  we obtain that in the case (a) of the simplification ( $\text{cov} = 0$ )

$$\text{var}(y) = \text{se}^2(N) + \text{se}^2(n) \leq 2\text{se}^2(n). \quad (51)$$

Let us elaborate on why the assumption that  $\check{b}_k$  is given alongside  $b_k$  and  $\nu$  (fixing the seed of subsampling operator) nullifies the covariance between the returns. According to (17) the only source of correlation between returns is (13) at time  $k$ . By fixing the seed, we made (13) a Dirac function. Therefore, conditioning on  $b_k$  and  $\nu$  is equivalent to conditioning on  $\check{b}_k$ .

In case of simplification (b),

$$\text{var}(y) = \text{se}^2(N) + \text{se}^2(n) - 2\text{cov} \leq 4\text{se}^2(n). \quad (52)$$

Thus, from (50) we obtain for both simplification possibilities

$$(a) P(|g_k - \check{g}_k| \leq z_{\alpha/2} \sqrt{2}\text{se}(n) | \mathcal{H}_{k+L}, \nu) \geq 1 - \alpha. \quad (53)$$

$$(b) P(|g_k - \check{g}_k| \leq z_{\alpha/2} 2\text{se}(n) | \mathcal{H}_{k+L}, \nu) \geq 1 - \alpha. \quad (54)$$

##### C. Estimation of the Variance

As we do not have access to  $\text{se}(n)$  in (53) and (54), it has to be estimated. The simplest way to do that is to repeatedly sample simplified returns  $m$  times from (17) in case of simplification (a) or from (18) in case of simplification (b). Note that a possible bias of the particle filter and the estimation of standard error make (50) only asymptotically correct. However, when dealing with a sufficient amount of samples  $N$  and  $n$ , these deviations from (46) are negligible. Even with repeated re-sampling we will reduce computational complexity, as we analyze in Section V-B

Moreover, since we recalculate the simplified reward  $m$  times, we could improve the final simplified return. However, this is out of the scope of this paper. To conclude, the bounds for both simplification methods are

$$(a) u = \check{g}_k + z_{\alpha/2} \sqrt{2}\hat{\text{se}}_m \quad l = \check{g}_k - z_{\alpha/2} \sqrt{2}\hat{\text{se}}_m \quad (55)$$

$$(b) u = \check{g}_k + z_{\alpha/2} 2\hat{\text{se}}_m \quad l = \check{g}_k - z_{\alpha/2} 2\hat{\text{se}}_m. \quad (56)$$

These bounds asymptotically hold with probability at least  $1 - \alpha$ .



#### D. Implementation Details and Computational Complexity

Now we describe steps in calculating  $\text{PbLoss}$ . First, we need to construct two extended belief policy trees appropriate to the two candidate policies (see Fig. 2). Second, we shall apply the simplification and calculate simplified returns and bounds. From now let us assume that  $\psi_{st}$  has a low-variance re-sampler [34]. The entire belief update process complexity is  $\mathcal{O}(N)$ . Since the extended belief tree does not undergo simplification, it is common to the original and simplified problems. Therefore, we discuss it in the Appendix E.

Now we analyze the speedup in running time as a result of simplification. As a momentary reward, we take the differential entropy estimator from [29], [4]. This selection makes the complexity of calculating the momentary reward to be  $\mathcal{O}(N^2)$ . Note it is customary to choose resampled or weighted belief for reward calculation. The resampled belief has identical weights. However, the effects of this choice are negligible.

For the bounds calculation in case of simplification (a) we need to apply particle filter with  $n$  samples (17) and in case of simplification (b) with  $N$  (18))  $L$  times for each return. Since its complexity is linear in the number of samples, the expected speedup is

$$\frac{N^2}{n^2 \cdot m}. \quad (57)$$

We obtained this speedup in all our simulations.

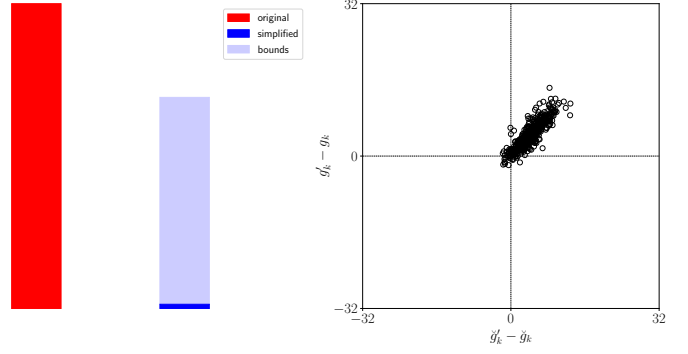
### V. RESULTS

#### A. Illustration via a Toy Example

We start with a toy example. Consider an underwater robot going as low as possible into the sea. Let robot's state at time instant  $k$  be  $x_k \in \mathbb{R}$  (altitude under the sea level). The theoretical distribution of  $x_k$  is out of reach. However, it is not Gaussian. Assume that the robot's belief about its state  $x_k$  at time instant  $k$  is represented by a set of  $N$  *i.i.d* samples, i.e.,  $b_k = \{x_k^i\}_{i=1}^N$ , produced by a perfect particle filter characterized by  $\psi_{st}$ . Further, consider the theoretical reward function  $\rho_{dt}(b) = \mathbb{E}_{x \sim b}[x]$ . Maximizing this reward will take the robot as deep into the sea as possible. The exact calculation of this reward is out of reach due to belief's sampled-based representation.

Therefore, the theoretical reward is replaced by sample mean, i.e.  $\rho_{k+1} = \frac{1}{N} \sum_{i=1}^N x_{k+1}^i$ .

Suppose the robot considers, as simplification, to use only  $n < N$  samples, such that  $\check{\rho}_{k+1} = \frac{1}{n} \sum_{j=1}^n x_{k+1}^j$ . This simplification could be used in two flavors. The first is to subsample  $b_k$ . The simplified version,  $\check{b}_k$ , serves as an input to a particle filter, producing  $\check{b}_{k+1}$  (17). The second is to subsample  $b_{k+1}$  to produce  $\check{b}_{k+1}$  (18). In general,  $\check{b}_{k+1}$  would differ in these two cases. As in previous sections, we denote these two genuine different simplification techniques by (a) and (b). We recall that in the case of (a), we assume a fixed seed of the sampler.



**Fig. 6:** Results for scenario 1 - probabilistic action consistency: (left) Demonstration of runtimes of the total number of the returns for a given extended belief tree where  $N = 1500$  and  $n = 175$ . Note that this illustration agrees with (57); (right) action consistency of the samples of the return.

As an example for stochastic bounds that hold with probability one, one could consider the following analytical bounds.

$$\min_i \{x_{k+1}^i\} \leq \frac{1}{N} \sum_{i=1}^N x_{k+1}^i \leq \max_i \{x_{k+1}^i\}. \quad (58)$$

We can define lower and upper bounds as  $l = \min_i \{x_{k+1}^i\}$  and  $u = \max_i \{x_{k+1}^i\}$ . These bounds bound any sample of the original reward  $\frac{1}{N} \sum_{i=1}^N x_{k+1}^i$  constructed from a corresponding set of  $N$  state samples. In other words, this is an example of bounds that hold with probability one. Of course, these bounds prone to be affected by outliers (one outlier can take  $u$  extremely far), and far better analytical bounds could be developed.

Another possibility is to consider the structure of (28) to define more lenient bounds. From the central limit theorem [8],[35] the following holds asymptotically

$$\rho_{k+1} | \mathcal{H}_{k+1} \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{N}\right), \check{\rho}_{k+1} | \mathcal{H}_{k+1}, \nu \sim \mathcal{N}\left(\mu, \frac{\sigma^2}{n}\right). \quad (59)$$

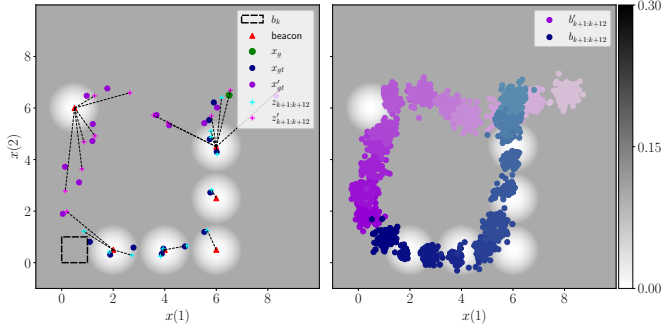
In the case of simplification (a), all samples besides  $b_k$  and  $\check{b}_k$  are independent (fixed seed). Therefore (28) is a multiplication of marginals. We arrive to (53) with  $\text{se}(n) = \frac{\sigma^2}{n}$  (similarly  $\text{se}(N) = \frac{\sigma^2}{N}$ ). Simplification (b) produces (54). Note that  $\sigma$  is unknown and has to be estimated, as in Section IV-C.

#### B. Autonomous Navigation with Light Beacons

We exemplify our method on the problem of autonomous navigation to a goal with light beacons, which can be used for localization. In all our simulations in this section, the return  $g_k$  is a cumulative reward, and, as a representative projection operator  $\varphi$ , we chose the expected value. In this study, the simplification is of type (a).

For simplicity, assume we have a linear motion model  $T$ , where  $x \in \mathbb{R}^2$  as well as  $a \in \mathbb{R}^2$

$$x_{k+1} = x_k + a_k + w_k \quad w_k \sim \mathcal{N}(0, \Sigma_w), \quad (60)$$



**Fig. 7:** Results for scenario 1 - probabilistic action consistency: Illustration of one realization of the future in a simulated scenario considering two possible action sequences. We start from  $b_k$  represented by samples uniformly distributed on a unit square. We demonstrated two sequences of observations alongside ground truth state samples, and the closest beacons produced these observations from the left. From the right, we plotted two sequences of the beliefs produced by these two histories. We show 100 most probable samples of each belief.

where  $\Sigma_w = w \cdot I$  ( $w$  is a given parameter) and action  $a_k \in \mathcal{A}$ , and where the action space  $\mathcal{A}$  is the space of motion primitives.

a) *Characterizing Probabilistic Action Consistency:* The observation model  $O$  is as follows,  $z \sim \mathcal{N}(x, \Sigma_v(x))$ , where the spatially-varying covariance matrix is

$$\Sigma_v(x) = v(x) \cdot I, \quad v(x) = w \cdot \min\{1, \|x - x^*\|_2^2\}, \quad (61)$$

where  $x^*$  is the location of the light beacon closest to  $x$ . The noise has a constant variance  $w$ . Without losing generality, we assume  $b_k$  at planning time is uniformly distributed in a unit square. We set  $L = 12$  and compare two action sequences:  $a_{k+1:k+12}$  is six times  $(1, 0)^T$  and after that six times  $(0, 1)^T$ . In the action sequence  $a'_{k+1:k+12}$  we switched the order of actions such that the robot performs six times  $(0, 1)^T$  and after that six times  $(1, 0)^T$ .

One realization of a possible future in terms of measurements and corresponding posterior beliefs is illustrated in Fig. 7. It is clearly seen that proximity to a beacon improves localization. Note the robot is always able to avoid a dead reckoning scenario as it always gets an observation from the closest beacon. This corresponds to a non-Gaussian noise, promoting the usage of particles-based belief representation. We hope that this setting conveys a real world scenario where an ambulating robot is equipped with long and short range sensors. The close range sensors are activated when the robot is inside a unit circle around the beacon. When the robot is outside a unit circle from the closest beacon, the beacon is detectable only by the long range sensors, which are less sensitive.

We present results of simplification (a) for  $w = 0.1$ ,  $N = 1500$ ,  $m = 50$ ,  $\alpha = 0.01$ ,  $z_{\alpha/2} = 2.56$ , and the total number of observations is 500. For each sample of  $z_{k+1:k+L}$ , we sampled  $b_{k+1:k+L}$  once.

$n$	$P(\mathcal{L} > 0.5\check{\Delta}^* \cdot)$	$\beta(0.5\check{\Delta}^*)$	$\check{\Delta}^*$	$P(\mathcal{L} > \check{\Delta}^* \cdot)$	$\beta(\check{\Delta}^*)$
175	0.0	0.33	4.14	0.0	0.11
150	0.01	0.43	4.04	0.0	0.17
125	0.01	0.43	4.21	0.0	0.2
100	0.0	0.56	4.08	0.0	0.29
75	0.01	0.64	4.01	0.0	0.39
50	0.02	0.83	3.72	0.01	0.63
25	0.07	1.0	3.34	0.03	0.94

**TABLE I:** Results for scenario 1 - probabilistic action consistency: Online characterization for  $N = 1500$ ,  $\alpha = 0.01$ ,  $z_{\alpha/2} = 2.56$ .

As we see in the left part of Fig. 6 we gained speedup as expected (57) for  $n = 175$ . We added measurements of all running times in our simulations to Appendix G.

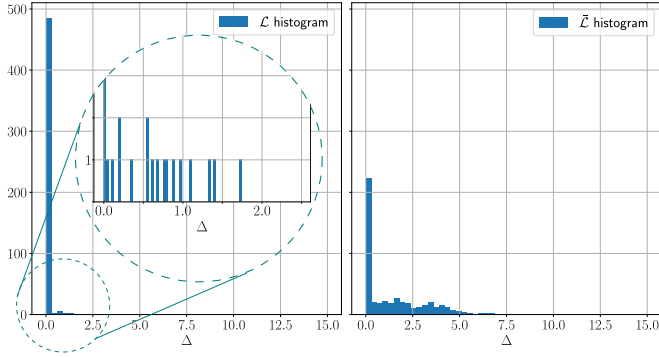
From these samples of the returns and bounds, we build  $\text{P}_{\text{Loss}}$  and  $\text{P}_{\text{bLoss}}$  in Fig. 8. In the right part of Fig. 6 quadrants II and IV, we observe samples that are not action consistent. To assess performance we need to choose some representative  $\Delta$ . Since online we have access exclusively to the simplified problem, let us choose  $\check{\Delta}^* = |\mathbb{E}[\check{g}_k|b_k, \pi, \nu] - \mathbb{E}[\check{g}'_k|b_k, \pi', \nu]|$  and  $\Delta = 0.5\check{\Delta}^*$ . Table I quantifies online characterization against offline  $\text{P}_{\text{Loss}}$  TDF.

In Fig. 9 we focused on  $n = 175$ ; *online* we can conclude that probability that loss incurred by this simplification will be greater than  $\check{\Delta}^*$  is at most 0.11, while actual  $P(\mathcal{L} > \check{\Delta}^*|\cdot)$  is 0.0. Similarly, the probability for loss incurred by this simplification to be greater than  $0.5\check{\Delta}^*$  is at most 0.33, while actual  $P(\mathcal{L} > 0.5\check{\Delta}^*|\cdot)$  is 0.0. In this scenario, the simplification is not absolute action consistent; it means variability described by (46) is sufficient to switch the order of the returns and incur loss  $\Delta$  at some sampled realization.

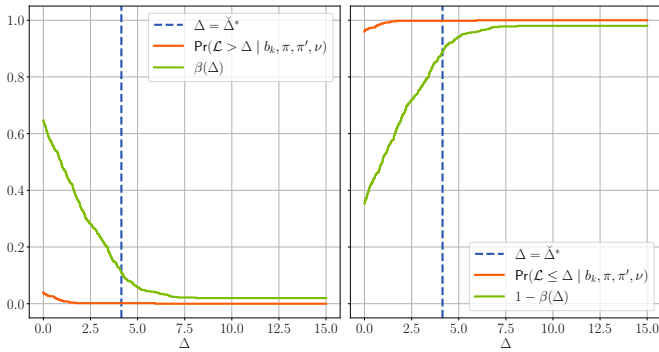
Furthermore, our bounds depend on variance ( $\text{se}^2(n)$ ) of the sample approximation of the reward (55), which, according to (46) does not depend on  $\Delta$ . Hence, as  $\Delta$  decreases towards zero, the contribution of variance versus the difference between simplified returns grows for any realization of  $\check{\mathcal{L}}$ . Therefore,  $\text{P}_{\text{bLoss}}$  departs from  $\text{P}_{\text{Loss}}$  as  $\Delta$  decreases. We observe this behavior in Fig. 9. Moreover, with the diminishing number of samples, this effect is amplified, as demonstrated in Fig. 14, due to growing variance (46). Remarkably, when samples of original returns are more distinct, the effect of variance is nullified. In such a setting, our characterization is incredibly precise, see Fig. 17.

Thus, the behavior of the  $\text{P}_{\text{bLoss}}$  is more conservative in more delicate scenarios, where two candidate policies are close to each other in terms of returns. Importantly, for significantly different policies,  $\text{P}_{\text{bLoss}}$  becomes tighter to  $\text{P}_{\text{Loss}}$ . This brings us to the next section.

b) *Revealing Empirical Absolute Action Consistency:* In this scenario we modified the noise in the observation model as such  $v(x) = w \cdot \|x - x^*\|_2^2$ . In addition we removed one beacon on the way of the second action sequence. We remain with  $w = 0.1$ ,  $m = 50$ ,  $\alpha = 0.01$ ,  $z_{\alpha/2} = 2.56$  and set  $N = 1000$ . In this scenario the returns of two action



**Fig. 8:** Results for scenario 1 - probabilistic action consistency: Histograms of  $P_{\text{Loss}}$  and  $P_{\text{bLoss}}$  for for  $N = 1500$ ,  $n = 175$ ,  $\alpha = 0.01$ ,  $z_{\alpha/2} = 2.56$  (bin width is 0.3, in zoom-in, bin width is 0.03).



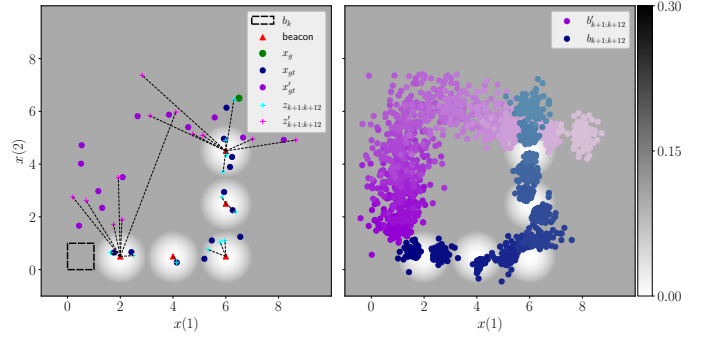
**Fig. 9:** Results for scenario 1 - probabilistic action consistency: Empirical characterization for  $N = 1500$ ,  $n = 175$ ,  $\alpha = 0.01$ ,  $z_{\alpha/2} = 2.56$ , evaluated in a grid with intervals 0.001.

sequences are much more distant. The samples in the right segment of Fig. 11 are more distant from the origin than in Fig. 6. The characterization is shown in Table II. Therefore, the simplification is empirically absolute action consistent. As we see from the Table II, observing  $\beta(\Delta = 0.0)$  we are able to identify online that for  $n = 100$  and  $n = 75$ , probability to receive samples of the returns violating action consistency is at most 0.03, while  $P(\mathcal{L} > 0.0|\cdot)$  is 0.0.

We placed additional results and discussions in Appendix G.

## VI. CONCLUSION

We extended  $\rho$ -POMDP to  $\mathbb{P}\rho$ -POMDP and introduced a framework for quantifying online the effect of simplification, alongside novel stochastic bounds on the return. Our bounds take advantage of the information encoded in the joint distribution of the original and simplified return. The proposed general framework is applicable to any bounds on the return to capture simplification outcomes. We presented experiments that confirmed the benefits of our approach. Our future research will focus on incorporating the contributed framework within simplified online decision making with probabilistic guarantees.



**Fig. 10:** Results for scenario 2 - empirical absolute action consistency: Illustration of one realization of the future in a simulated scenario considering two possible action sequences. We start from  $b_k$  represented by samples uniformly distributed on a unit square. We demonstrated two sequences of observations alongside ground truth state samples, and the closest beacons produced these observations from the left. From the right, we plotted two sequences of the beliefs produced by these two histories. We show 100 most probable samples of each belief.

$n$	$P(\mathcal{L} > 0.0 \cdot)$	$\beta(\Delta = 0.0)$	$\hat{\Delta}^*$	$P(\mathcal{L} > \hat{\Delta}^* \cdot)$	$\beta(\hat{\Delta}^*)$
100	0.0	0.03	17.54	0.0	0.02
75	0.0	0.03	17.14	0.0	0.02
50	0.0	0.06	16.65	0.0	0.02
25	0.0	0.19	15.27	0.0	0.02

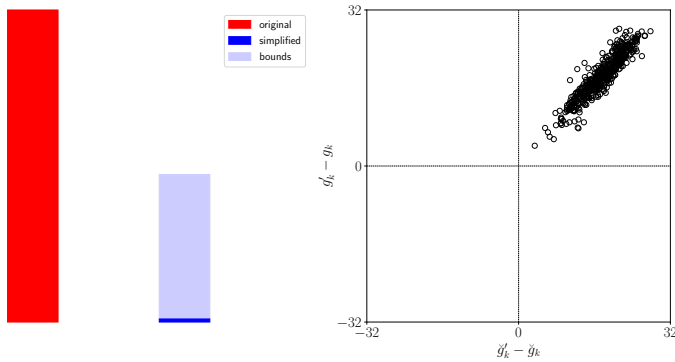
**TABLE II:** Results for scenario 2 - empirical absolute action consistency: Online characterization for  $N = 1000$ ,  $\alpha = 0.01$ ,  $z_{\alpha/2} = 2.56$ .

## ACKNOWLEDGMENTS

This research was supported by the Israel Science Foundation (ISF), by the Israel Ministry of Science & Technology (MOST), and by a donation from the Zuckerman Fund to the Technion Center for Machine Learning and Intelligent Systems (MLIS).

## REFERENCES

- [1] Mauricio Araya, Olivier Buffet, Vincent Thomas, and François Charpillet. A pomdp extension with belief-dependent rewards. *Advances in neural information processing systems*, 23:64–72, 2010.
- [2] Haoyu Bai, David Hsu, and Wee Sun Lee. Integrated perception and planning in the continuous space: A pomdp approach. *The International Journal of Robotics Research*, 33(9):1288–1302, 2014.
- [3] Dimitri P. Bertsekas. *Dynamic programming and optimal control*. 4ed, volume 1. Athena scientific Belmont, MA, 2017.
- [4] Yvo Boers, Hans Driessen, Arunabha Bagchi, and Pranab Mandal. Particle filter based entropy. In *2010 13th International Conference on Information Fusion*, pages 1–8. IEEE, 2010.



**Fig. 11:** Results for scenario 2 - empirical absolute action consistency: (left) Demonstration of runtimes of the total number of the returns for a given extended belief tree where  $N = 1000$  and  $n = 100$ . Note that this illustration agrees with (57); (right) action consistency of the samples of the return.

- [5] Mokrane Bouakiz and Youcef Kebir. Target-level criterion in markov decision processes. *Journal of Optimization Theory and Applications*, 86(1):1–15, 1995.
- [6] Boris Defourny, Damien Ernst, and Louis Wehenkel. Risk-aware decision making and dynamic programming. In *NIPS Workshop on Model Uncertainty and Risk in RL*, 2008.
- [7] Louis Dressel and Mykel J Kochenderfer. Efficient decision-theoretic target localization. In *Twenty-Seventh International Conference on Automated Planning and Scheduling*, 2017.
- [8] Rick Durrett. *Probability: theory and examples*, volume 49. Cambridge university press, 2019.
- [9] Khen Elimelech and Vadim Indelman. Simplified decision making in the belief space using belief sparsification. *arXiv preprint arXiv:1909.00885*, 2019.
- [10] Khen Elimelech and Vadim Indelman. Fast action elimination for efficient decision making and belief space planning using bounded approximations. In *Robotics Research*, pages 843–858. Springer, 2020.
- [11] Elad I Farhi and Vadim Indelman. ix-bsp: Belief space planning through incremental expectation. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 7180–7186. IEEE, 2019.
- [12] Elad I. Farhi and Vadim Indelman. ix-bsp: Incremental belief space planning. *arXiv preprint arXiv:2102.09539*, 2021.
- [13] Mathieu Fehr, Olivier Buffet, Vincent Thomas, and Jilles Dibangoye. rho-pomdps have lipschitz-continuous epsilon-optimal value functions. *Advances in neural information processing systems*, 31:6933–6943, 2018.
- [14] Johannes Fischer and Ömer Sahin Tas. Information particle filter tree: An online algorithm for pomdps with belief-based rewards on continuous domains. In *International Conference on Machine Learning*, pages 3177–3187. PMLR, 2020.
- [15] Neha Priyadarshini Garg, David Hsu, and Wee Sun Lee. Despot-alpha: Online pomdp planning with large state and observation spaces. In *Robotics: Science and Systems*, 2019.
- [16] Vadim Indelman. No correlations involved: Decision making under uncertainty in a conservative sparse information space. *IEEE Robotics and Automation Letters*, 1(1):407–414, 2016.
- [17] Vadim Indelman, Luca Carlone, and Frank Dellaert. Planning in the continuous domain: A generalized belief space approach for autonomous navigation in unknown environments. *The International Journal of Robotics Research*, 34(7):849–882, 2015.
- [18] L. P. Kaelbling, M. L. Littman, and A. R. Cassandra. Planning and acting in partially observable stochastic domains. *Artificial intelligence*, 101(1):99–134, 1998.
- [19] Andrej Kitanov and Vadim Indelman. Topological information-theoretic belief space planning with optimality guarantees. *arXiv preprint arXiv:1903.00927*, 2019.
- [20] Donald E Knuth. *Art of computer programming, volume 2: Seminumerical algorithms*. Addison-Wesley Professional, 2014.
- [21] Mykel J Kochenderfer. *Decision making under uncertainty: theory and application*. MIT press, 2015.
- [22] Levente Kocsis and Csaba Szepesvári. Bandit based monte-carlo planning. In *European conference on machine learning*, pages 282–293. Springer, 2006.
- [23] Hanna Kurniawati and Vinay Yadav. An online pomdp solver for uncertainty planning in dynamic environment. In *Robotics Research*, pages 611–629. Springer, 2016.
- [24] Hanna Kurniawati, David Hsu, and Wee Sun Lee. Sarsop: Efficient point-based pomdp planning by approximating optimally reachable belief spaces. In *Robotics: Science and systems*, volume 2008. Zurich, Switzerland., 2008.
- [25] Christos H Papadimitriou and John N Tsitsiklis. The complexity of markov decision processes. *Mathematics of operations research*, 12(3):441–450, 1987.
- [26] R Tyrrell Rockafellar, Stanislav Uryasev, et al. Optimization of conditional value-at-risk. *Journal of risk*, 2: 21–42, 2000.
- [27] Allison Ryan. Information-theoretic tracking control based on particle filter estimate. In *AIAA Guidance, Navigation and Control Conference and Exhibit*, page 6307, 2008.
- [28] David Silver and Joel Veness. Monte-carlo planning in large pomdps. In *Advances in neural information processing systems*, pages 2164–2172, 2010.
- [29] Per Skoglar, Umut Orguner, and Fredrik Gustafsson. On information measures based on particle mixture for optimal bearings-only tracking. In *2009 IEEE Aerospace conference*, pages 1–14. IEEE, 2009.
- [30] Adhiraj Somani, Nan Ye, David Hsu, and Wee Sun Lee. Despot: Online pomdp planning with regularization. *Advances in neural information processing systems*, 26: 1772–1780, 2013.
- [31] Matthijs TJ Spaan, Tiago S Veiga, and Pedro U Lima. Decision-theoretic planning under uncertainty with in-

- formation rewards for active cooperative perception. *Autonomous Agents and Multi-Agent Systems*, 29(6):1157–1185, 2015.
- [32] Zachary Sunberg and Mykel Kochenderfer. Online algorithms for pomdps with continuous state, action, and observation spaces. *arXiv preprint arXiv:1709.06196*, 2017.
- [33] Richard S Sutton and Andrew G Barto. *Reinforcement learning: An introduction*. MIT press, 2018.
- [34] Sebastian Thrun, Wolfram Burgard, and Dieter Fox. *Probabilistic robotics*. MIT press, 2005.
- [35] Larry Wasserman. *All of statistics: a concise course in statistical inference*. Springer Science & Business Media, 2013.
- [36] Congbin Wu and Yuanlie Lin. Minimizing risk models in markov decision processes with policies depending on target values. *Journal of mathematical analysis and applications*, 231(1):47–67, 1999.
- [37] Nan Ye, Adhiraj Somani, David Hsu, and Wee Sun Lee. Despot: Online pomdp planning with regularization. *JAIR*, 58:231–266, 2017.

APPENDIX A  
DISCUSSION ON BELIEF-MDP

In belief-MDP

$$\mathbb{P}(b_\ell | b_{\ell-1}, \pi_{\ell-1}(b_{\ell-1})) = \int_{z_\ell} \mathbb{P}(b_\ell | b_{\ell-1}, \pi_{\ell-1}(b_{\ell-1}), z_\ell) \mathbb{P}(z_\ell | b_{\ell-1}, \pi_{\ell-1}(b_{\ell-1})), \quad (62)$$

where the usual assumption is that  $\mathbb{P}(b_\ell | b_{\ell-1}, \pi_{\ell-1}(b_{\ell-1}), z_\ell)$  is Dirac's delta function. On the contrary, in our extended setting, this distribution is arbitrary.

APPENDIX B  
PROOFS

**Theorem 4** (Probability that bound bounds). *If*

$$(1 - \alpha) \leq P([l \leq g_k \leq u] = 1 | \mathcal{H}_{k+L}, \nu) \quad (63)$$

and

$$(1 - \alpha) \leq P([l' \leq g'_k \leq u'] = 1 | \mathcal{H}'_{k+L}, \nu), \quad (64)$$

so

$$P([\bar{\mathcal{L}} \geq \mathcal{L}] = 1 | b_k, \pi, \pi', \nu) \geq (1 - \alpha)^2. \quad (65)$$

*Proof:* By definition

$$P([\bar{\mathcal{L}} \geq \mathcal{L}] = 1 | [l \leq g_k \leq u] = 1, [l' \leq g'_k \leq u'] = 1, b_k, \pi, \pi', z_{k+}, z'_{k+}, \nu) = 1. \quad (66)$$

We first apply marginalization over future observations  $z_{k+} \equiv z_{k+1:k+L}$  and  $z'_{k+} \equiv z'_{k+1:k+L}$ , and events  $[l \leq g_k \leq u]$  and  $[l' \leq g'_k \leq u']$ . We then use the fact that given two histories  $\mathcal{H}_{k+L} \triangleq \{b_k, \pi, z_{k+}\}$  and  $\mathcal{H}'_{k+L} \triangleq \{b_k, \pi', z'_{k+}\}$ , the events  $[l \leq g_k \leq u]$  and  $[l' \leq g'_k \leq u']$  are independent of each other. Furthermore, each such event depends exclusively on its own history by design. We have that

$$P([\bar{\mathcal{L}} \geq \mathcal{L}] = 1 | b_k, \pi, \pi', \nu) = \int_{z_{k+}} P([\bar{\mathcal{L}} \geq \mathcal{L} | b_k, \pi, \pi', z_{k+}, z'_{k+}, \nu) \mathbb{P}(z_{k+}, z'_{k+} | b_k, \pi, \pi') \geq \quad (67)$$

$$\int_{z'_{k+}} P([\bar{\mathcal{L}} \geq \mathcal{L}] = 1, [l \leq g_k \leq u] = 1, [l' \leq g'_k \leq u'] = 1 | b_k, \pi, \pi', z_{k+}, z'_{k+}, \nu) \mathbb{P}(z_{k+}, z'_{k+} | b_k, \pi, \pi') = \quad (68)$$

$$\int_{z'_{k+}} P([\bar{\mathcal{L}} \geq \mathcal{L}] = 1 | [l \leq g_k \leq u] = 1, [l' \leq g'_k \leq u'] = 1, b_k, \pi, \pi', z_{k+}, z'_{k+}, \nu) \quad (69)$$

$$P([l \leq g_k \leq u] = 1, [l' \leq g'_k \leq u'] = 1 | b_k, \pi, \pi', z_{k+}, z'_{k+}, \nu) \mathbb{P}(z_{k+}, z'_{k+} | b_k, \pi, \pi') = \quad (70)$$

$$\int_{z_{k+}} P([l \leq g_k \leq u] = 1 | b_k, \pi, z_{k+}, \nu) P([l' \leq g'_k \leq u'] = 1 | b_k, \pi', z'_{k+}, \nu) \mathbb{P}(z_{k+}, z'_{k+} | b_k, \pi, \pi') \geq (1 - \alpha)^2. \quad (70)$$

This completes the proof. ■

**Theorem 5** (Conditional TDF Lower bound).

$$P(\mathcal{L} > \Delta | \bar{\mathcal{L}} \geq \mathcal{L}, b_k, \pi, \pi', \nu) \leq P(\bar{\mathcal{L}} > \Delta | \bar{\mathcal{L}} \geq \mathcal{L}, b_k, \pi, \pi', \nu) \quad \forall \Delta \in \mathbb{R}. \quad (71)$$

*Proof:* Let us recall the definition of the probability space and the random variable [8]. A probability space is a triple  $(\Omega, \mathcal{F}, P)$ , where  $\Omega$  is a set of "outcomes",  $\mathcal{F}$  is a set of "events", and  $P : \mathcal{F} \rightarrow [0, 1]$  is a function that assigns probabilities to events.  $\mathcal{F}$  is a  $\sigma$ -field, i.e., a (nonempty) collection of subsets of  $\Omega$  that satisfy certain conditions.

A real valued function  $X$  defined on  $\Omega$  is said to be a random variable if for every Borel set  $B \subset \mathbb{R}$  we have

$$X^{-1}(B) = \{\omega : X(\omega) \in B\} \in \mathcal{F}. \quad (72)$$

In our case,  $\Omega$  is the same for  $\mathcal{L}$  and  $\bar{\mathcal{L}}$ , since we condition on the same information

$$P(\mathcal{L} > \Delta | \bar{\mathcal{L}} \geq \mathcal{L}, b_k, \pi, \pi', \nu) = P(\{\omega \subseteq \Omega : \mathcal{L}(\omega) > \Delta\}). \quad (73)$$

It follows that

$$\{\omega \subseteq \Omega : \mathcal{L}(\omega) > \Delta\} \subseteq \{\omega \subseteq \Omega : \bar{\mathcal{L}}(\omega) > \Delta\}. \quad (74)$$

This completes the proof. ■

**Theorem 6** (Upper and Lower bounds). Denote  $\beta(\Delta) \triangleq \min \left\{ 1, \frac{P(\bar{\mathcal{L}} > \Delta | b_k, \pi, \pi', \nu)}{(1-\alpha)^2} + 2\alpha - \alpha^2 \right\}$ , so

$$P(\mathcal{L} > \Delta | b_k, \pi, \pi', \nu) \leq \beta(\Delta) \quad (75)$$

and

$$P(\mathcal{L} \leq \Delta | b_k, \pi, \pi', \nu) \geq 1 - \beta(\Delta). \quad (76)$$

*Proof:* To shorten notations let us denote  $|b_k, \pi, \pi', \nu$  by  $|\cdot$  in the proof. Let us express  $P_{\text{Loss}}$  TDF as

$$P(\mathcal{L} > \Delta | \cdot) = P(\mathcal{L} > \Delta | \bar{\mathcal{L}} \geq \mathcal{L}, \cdot) P(\bar{\mathcal{L}} \geq \mathcal{L} | \cdot) + P(\mathcal{L} > \Delta | \bar{\mathcal{L}} < \mathcal{L}, \cdot) P(\bar{\mathcal{L}} < \mathcal{L} | \cdot). \quad (77)$$

Similarly,  $P_{\text{bLoss}}$  TDF reads

$$P(\bar{\mathcal{L}} > \Delta | \cdot) = P(\bar{\mathcal{L}} > \Delta | \bar{\mathcal{L}} \geq \mathcal{L}, \cdot) P(\bar{\mathcal{L}} \geq \mathcal{L} | \cdot) + P(\bar{\mathcal{L}} > \Delta | \bar{\mathcal{L}} < \mathcal{L}, \cdot) P(\bar{\mathcal{L}} < \mathcal{L} | \cdot). \quad (78)$$

Since  $\alpha \in [0, 1)$  it exists  $c \in \mathbb{R}$  such that

$$P(\bar{\mathcal{L}} > \Delta | \cdot) = c(1 - \alpha)^2. \quad (79)$$

This implies

$$P(\bar{\mathcal{L}} > \Delta | \bar{\mathcal{L}} \geq \mathcal{L}, \cdot) P(\bar{\mathcal{L}} \geq \mathcal{L} | \cdot) \leq c(1 - \alpha)^2, \quad (80)$$

$$P(\bar{\mathcal{L}} > \Delta | \bar{\mathcal{L}} \geq \mathcal{L}, \cdot) \leq c \underbrace{\frac{(1 - \alpha)^2}{P(\bar{\mathcal{L}} \geq \mathcal{L} | \cdot)}}_{\leq 1} \leq c. \quad (81)$$

Moreover, using that  $P(\bar{\mathcal{L}} \geq \mathcal{L} | \cdot) + P(\bar{\mathcal{L}} < \mathcal{L} | \cdot) = 1$ , we obtain

$$\begin{aligned} P(\mathcal{L} > \Delta | \cdot) &= P(\mathcal{L} > \Delta | \bar{\mathcal{L}} \geq \mathcal{L}, \cdot) P(\bar{\mathcal{L}} \geq \mathcal{L} | \cdot) + P(\mathcal{L} > \Delta | \bar{\mathcal{L}} < \mathcal{L}, \cdot) (1 - P(\bar{\mathcal{L}} \geq \mathcal{L} | \cdot)) \leq \\ &P(\mathcal{L} > \Delta | \bar{\mathcal{L}} \geq \mathcal{L}, \cdot) + 1 - (1 - \alpha)^2 \leq c + 2\alpha - \alpha^2. \end{aligned} \quad (82)$$

but  $c = \frac{P(\bar{\mathcal{L}} > \Delta | b_k, \pi, \pi', \nu)}{(1-\alpha)^2}$ . We showed that

$$P(\mathcal{L} > \Delta | \cdot) \leq \frac{P(\bar{\mathcal{L}} > \Delta | b_k, \pi, \pi', \nu)}{(1 - \alpha)^2} + 2\alpha - \alpha^2. \quad (83)$$

Furthermore, by definition of TDF

$$P(\mathcal{L} > \Delta | \cdot) \leq 1. \quad (84)$$

We write the above two relations compactly as

$$P(\mathcal{L} > \Delta | \cdot) \leq \beta(\Delta), \quad (85)$$

where  $\beta(\Delta) = \min \left\{ 1, \frac{P(\bar{\mathcal{L}} > \Delta | b_k, \pi, \pi', \nu)}{(1-\alpha)^2} + 2\alpha - \alpha^2 \right\}$ . Clearly

$$P(\mathcal{L} \leq \Delta | b_k, \pi, \pi', \nu) = 1 - P(\mathcal{L} > \Delta | b_k, \pi, \pi', \nu) \geq 1 - \beta(\Delta). \quad (86)$$

This completes the proof. ■

APPENDIX C  
CALCULATION OF PLOSS OFFLINE

Similar to  $\text{PbLoss}$ , one approach to obtain  $\text{PLOSS}$  *offline* is to sample  $(g_k, g'_k, \check{g}_k, \check{g}'_k) \sim \mathbb{P}(g_k, g'_k, \check{g}_k, \check{g}'_k | b_k, \pi, \pi', \nu)$ .  $\text{PLOSS}$  is represented by  $\{f_{\mathcal{L}}(g_k, g'_k, \check{g}_k, \check{g}'_k)\}$ .

To generate samples  $(g_k, g'_k, \check{g}_k, \check{g}'_k)$  we marginalize over future measurements  $z_{k+} \equiv z_{k+1:k+L}$  and  $z'_{k+} \equiv z'_{k+1:k+L}$ . As we mentioned in the main manuscript

$$\begin{aligned} & \int_{z_{k+}}^{z'_{k+}} \mathbb{P}(g_k, g'_k, \check{g}_k, \check{g}'_k | b_k, \pi, \pi', \nu, z_{k+}, z'_{k+}) \mathbb{P}(z_{k+}, z'_{k+} | b_k, \pi, \pi') dz_{k+} dz'_{k+} = \\ & \int_{z_{k+}}^{z'_{k+}} \mathbb{P}(g_k, \check{g}_k | \mathcal{H}_{k+L}, \nu) \mathbb{P}(g'_k, \check{g}'_k | \mathcal{H}'_{k+L}, \nu) \cdot \mathbb{P}(z_{k+}, z'_{k+} | b_k, \pi, \pi') dz_{k+} dz'_{k+}. \end{aligned} \quad (87)$$

We take samples of  $\mathbb{P}(z_{k+}, z'_{k+} | b_k, \pi, \pi')$  from the corresponding extended belief trees built for  $\text{PbLoss}$ . To sample

$$\mathbb{P}(g_k, \check{g}_k | \mathcal{H}_{k+L}, \nu) \quad , \quad \mathbb{P}(g'_k, \check{g}'_k | \mathcal{H}'_{k+L}, \nu), \quad (88)$$

we use the original (not simplified) rewards calculated from the beliefs present at the belief tree (belief tree does not undergo simplification) and their simplified counterparts.

APPENDIX D  
DERIVATION OF THE DISTRIBUTION OF A LINEAR FUNCTION OF GAUSSIAN RANDOM VECTOR

We are interested in the following distribution

$$A \cdot \begin{pmatrix} g_k \\ \check{g}_k \end{pmatrix}, \quad (89)$$

where  $A = \begin{pmatrix} 1 & -1 \end{pmatrix}$ . Denote by  $\phi_x(t) = \exp(it^T \mu - \frac{1}{2}t^T \Sigma t)$  the characteristic function of  $x = \begin{pmatrix} g_k \\ \check{g}_k \end{pmatrix}$  and by  $\phi_y(t)$  the desired characteristic function (of  $y = Ax$ ), we have that

$$\phi_y(t) = \mathbb{E} [\exp(it^T Ax)] = \quad (90)$$

$$\mathbb{E} [\exp(i(A^T t)^T x)] = \quad (91)$$

$$\phi_x(A^T t) = \exp\left(i(A^T t)^T \mu - \frac{1}{2}(A^T t)^T \Sigma (A^T t)\right) = \quad (92)$$

$$\exp\left(it^T A \mu - \frac{1}{2}t^T A \Sigma A^T t\right). \quad (93)$$

In other words,  $y = g_k - \check{g}_k$  is zero mean Gaussian with the following variance

$$\text{var}(y) = \text{se}^2(N) + \text{se}^2(n) - 2\text{cov}. \quad (94)$$

APPENDIX E  
DERIVATION AND COMPLEXITY OF TWO DEPENDENT EXTENDED BELIEF POLICY TREES

In a theoretical form,

$$\begin{aligned} & \mathbb{P}(z_{k+1:k+L}, z'_{k+1:k+L} | b_k, \pi, \pi') = \mathbb{P}(z_{k+1}, z'_{k+1} | b_k, \pi_k, \pi'_k) \\ & \prod_{\ell=k+2}^{k+L-1} \mathbb{P}(z_\ell | b_{\ell-1}, \pi_{\ell-1}) \mathbb{P}(z'_\ell | b'_{\ell-1}, \pi'_{\ell-1}) \int_{b_\ell} \mathbb{P}(b_\ell | b_{\ell-1}, \pi_{\ell-1}, z_\ell) \int_{b'_\ell} \mathbb{P}(b'_\ell | b'_{\ell-1}, \pi'_{\ell-1}, z'_\ell) \\ & \mathbb{P}(z_{k+L} | b_{k+L-1}, \pi_{k+L-1}) \mathbb{P}(z'_{k+L} | b'_{k+L-1}, \pi'_{k+L-1}). \end{aligned} \quad (95)$$

However, realizations of the future are correlated through the belief from planning time (present),

$$\begin{aligned} & \mathbb{P}(z_{k+1}, z'_{k+1} | b_k, \pi_k, \pi'_k) = \int_{x_{k+1}, x'_{k+1}} \mathbb{P}_Z(z_{k+1} | x_{k+1}) \mathbb{P}_Z(z'_{k+1} | x'_{k+1}) \int_{x_k} \mathbb{P}(x_{k+1}, x'_{k+1} | x_k, b_k, \pi_k, \pi'_k) b_k(x_k) = \\ & \int_{x_{k+1}, x'_{k+1}} \mathbb{P}_Z(z_{k+1} | x_{k+1}) \mathbb{P}_Z(z'_{k+1} | x'_{k+1}) \int_{x_k} \mathbb{P}_T(x_{k+1} | \pi_k(b_k), x_k) \mathbb{P}_T(x'_{k+1} | \pi'_k(b_k), x_k) b_k(x_k). \end{aligned} \quad (96)$$

In practice, the mutual likelihood of the observations  $\mathbb{P}(z_{k+}, z'_{k+} | b_k, \pi, \pi')$  (95) corresponds to two extended belief policy trees, starting from the same root ( $b_k$ ) and having the same rule for choosing rollouts.



Below we discuss the construction of the extended belief tree. Let  $N$  be a number of samples of the posterior belief. From now let us assume that  $\psi_{st}$  is an off-the-shelf particle filter with low-variance re-sampling [34]. The entire belief update process complexity is  $\mathcal{O}(N)$ . We choose the samples of the belief for creating the observations according to the following scheme, which is inspired by progressive widening [32]. Let  $n_z^{(\ell)}$  be number of observations generated by each belief at level  $\ell$  of the tree. We specify  $n_z^{(1)}$  (the number of observations generated by  $b_k$ ) and the dwindle factor  $c$ . Starting from  $\ell = 2$  the number of observations generated by each belief on level  $\ell$  in the tree is calculated as  $n_z^{(\ell)} = \max\{1, \lfloor \frac{n_z^{(1)}}{(\ell-1) \cdot c} \rfloor\}$ . We sample observations from resampled posterior with Fisher-Yates shuffling (with early termination) [20]. This algorithm is  $\mathcal{O}(N)$  for initialization, plus  $\mathcal{O}(n_z^{(\ell)})$  for random shuffling.

In our extended belief policy tree, there may be many beliefs stemming from an observation. Denote this number by  $n_b$ . The complexity of constructing the tree is

$$\mathcal{O}(N) \sum_{\ell=1}^{L-1} \prod_{i=1}^{\ell} n_b n_z^{(i)}. \quad (97)$$

At each level of the tree beside the bottom, we must apply a particle filter number of times equal to the total number of the beliefs at the next level, which is  $\prod_{i=1}^{\ell} n_b n_z^{(i)}$  at level  $\ell$ . Also, we need to subsample observations at the current level. Since the number of beliefs at the next level is not smaller than at the current level, and the subsampler and particle filter complexity is linear in  $N$ , we are left with (97). Let us mention that sampling from the belief and application of particle filters on each level can be done in parallel.

#### APPENDIX F DIFFERENTIAL ENTROPY

We adopt the differential entropy estimator from [29, 4]. Suppose  $b_k$  is given by a weighted set of samples

$$b_k = \sum_{i=1}^n w_k^i \delta(x_k - x_k^i). \quad (98)$$

Propagated belief is

$$b_{k+1}^- = \sum_{i=1}^n w_k^i \mathbb{P}_T(x_{k+1} | a_k, x_k^i) \approx \sum_{i=1}^n w_k^i \delta(x_{k+1} - x_{k+1|k}^i). \quad (99)$$

The posterior is

$$b_{k+1} \approx \sum_{i=1}^n w_{k+1}^i \delta(x_{k+1} - x_{k+1|k}^i). \quad (100)$$

We plug the particle approximation of the propagated belief (99) into the likelihood of the observation in subsequent time instant and obtain

$$\mathbb{P}(z_{k+1} | b_k, a_k) = \int_{x_{k+1}} \mathbb{P}_Z(z_{k+1} | x_{k+1}) b_{k+1}^-(x_{k+1}) \approx \sum_{i=1}^n w_k^i \mathbb{P}_Z(z_{k+1} | x_{k+1|k}^i). \quad (101)$$

We assumed here that weights are normalized. Using Bayes rule, we have

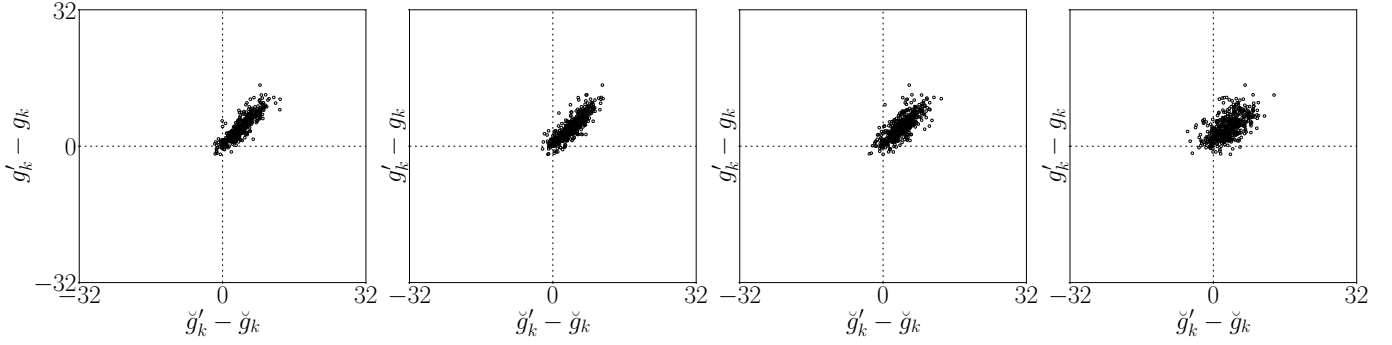
$$\begin{aligned} \mathcal{H}(b_{k+1}) &= - \int_{x_{k+1}} b_{k+1}(x_{k+1}) \ln \mathbb{P}_Z(z_{k+1} | x_{k+1}) \\ &\quad - \int_{x_{k+1}} b_{k+1}(x_{k+1}) \ln b_{k+1}^-(x_{k+1}) + \ln \mathbb{P}(z_{k+1} | b_k, a_k), \end{aligned} \quad (102)$$

and, therefore,

$$\begin{aligned} \mathcal{H}(b_{k+1}) &= - \sum_{j=1}^n w_{k+1}^j \ln \mathbb{P}_Z(z_{k+1} | x_{k+1|k}^j) \\ &\quad - \sum_{j=1}^n w_{k+1}^j \ln \sum_{i=1}^n w_k^i \mathbb{P}_T(x_{k+1|k}^j | a_k, x_k^i) + \ln \sum_{i=1}^n w_k^i \mathbb{P}_Z(z_{k+1} | x_{k+1|k}^i). \end{aligned} \quad (103)$$

#### APPENDIX G ADDITIONAL RESULTS AND DISCUSSIONS

This section shows additional results for two of the scenarios we presented in the main manuscript. In all our simulations we set  $w = 0.1$ ,  $m = 50$ ,  $\alpha = 0.01$ ,  $z_{\alpha/2} = 2.56$ , and the total number of observations is 500.



**Fig. 12:** Results for scenario 1 - probabilistic action consistency: We demonstrate from the left to the right action consistency of the samples of the returns for  $n = 175, n = 125, n = 75, n = 25$ , whereas  $N = 1500$ . As we see, samples violating action consistency are present at all graphs.

	$n = 175$	$n = 150$	$n = 125$	$n = 100$	$n = 75$	$n = 50$	$n = 25$
$g_k$ time [sec]	104957	69658	95651	69713	68584	96354	66513
$\check{g}_k$ and $l, u$ time [sec]	72694	34842	33759	15498	8293	5589	969
$\check{g}_k$ time [sec]	1454	661	669	298	172	119	14
$l, u$ time [sec]	71240	34181	33090	15200	8121	5469	955

**TABLE III:** Results for scenario 1 - probabilistic action consistency: run times for  $N = 1500$ .

a) *Characterizing Probabilistic Action Consistency:* In this scenario

$$\Sigma_v(x) = v(x) \cdot I, \quad v(x) = w \cdot \min\{1, \|x - x^*\|_2^2\}, \quad (104)$$

where  $x^*$  is the location of the light beacon closest to  $x$ .

We showed an illustration of this scenario in Fig. 7. In Fig. 12, we demonstrated scatter plots that show samples of the simplified and original returns' differences. We identify that with decreasing  $n$ , more samples are not action consistent. This phenomenon is corroborated by the histograms of  $\mathcal{L}$  in Fig. 13. We recite the model

$$\begin{pmatrix} g_k \\ \check{g}_k \end{pmatrix} | \mathcal{H}_{k+L}, \nu \sim \mathcal{N} \left( \begin{pmatrix} \mu \\ \mu \end{pmatrix}, \begin{pmatrix} \text{se}^2(N) & \text{cov} \\ \text{cov} & \text{se}^2(n) \end{pmatrix} \right). \quad (105)$$

With decreasing  $n$ , the variance ( $\text{se}^2(n)$ ) in (105) grows, making more samples of the returns are not action consistent. Moreover, since our probabilistic bounds on the return are based on ( $\text{se}(n)$ ), it is harder to characterize  $\text{PLoss}$  with  $\text{PbLoss}$ . This is observable in histograms of  $\bar{\mathcal{L}}$  in Fig. 13 and empirical characterization in Fig. 14, where

$$\text{empirical PLoss TDF: } \frac{\text{number of samples of } \mathcal{L}, \text{ satisfying } \mathcal{L} > \Delta}{\text{number of all samples of } \mathcal{L}}, \quad (106)$$

$$\text{empirical PbLoss TDF: } \frac{\text{number of samples of } \bar{\mathcal{L}}, \text{ satisfying } \bar{\mathcal{L}} > \Delta}{\text{number of all samples of } \bar{\mathcal{L}}}. \quad (107)$$

b) *Revealing Empirical Absolute Action Consistency:* Here the covariance matrix of the observation model is

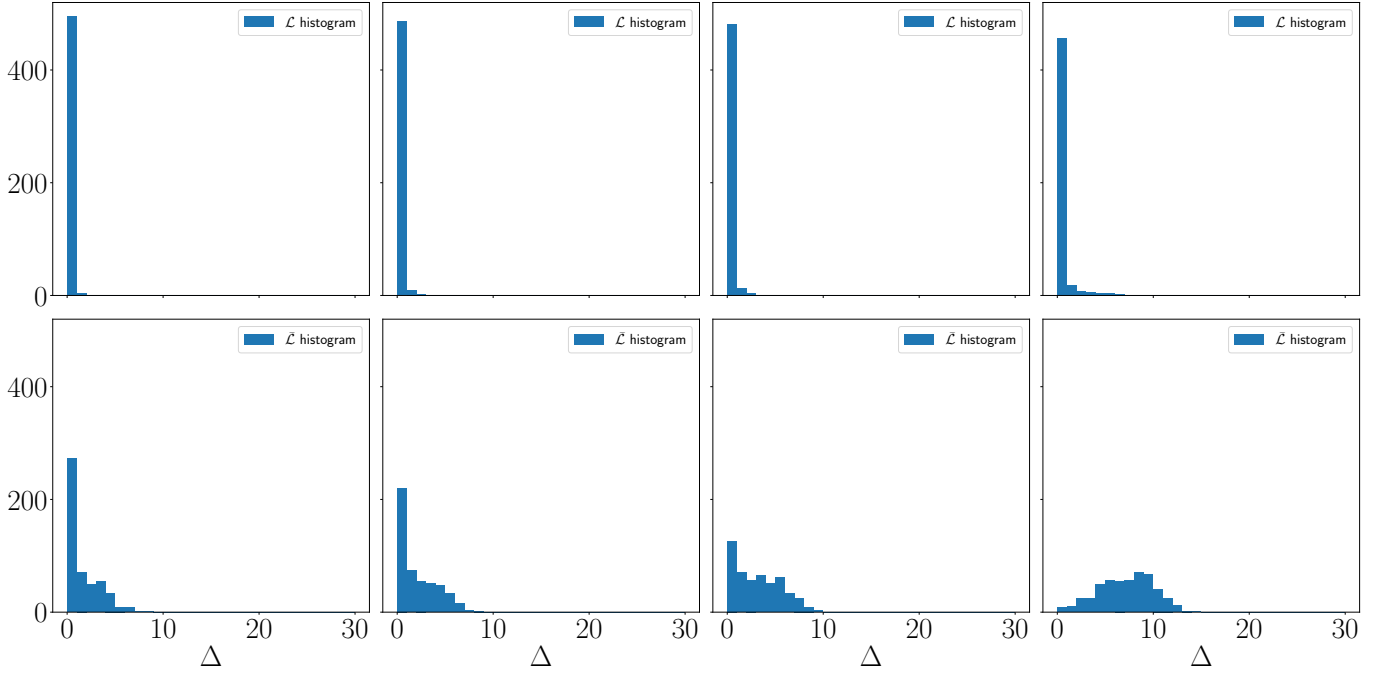
$$\Sigma_v(x) = v(x) \cdot I, \quad v(x) = w \cdot \|x - x^*\|_2^2. \quad (108)$$

We demonstrated this scenario in Fig. 10. As we can see in Fig. 15, the clouds of samples are farther from the origin than in the previous scenario. Therefore, two action sequences are more distant. In this case, the simplification is empirically absolute action consistent, as we observe in the histograms of  $\mathcal{L}$  in Fig. 16 and empirical characterization shown in Fig. 17. We report run times for two scenarios in Table III and Table IV, respectively.

## APPENDIX H TECHNICAL CHARACTERISTICS OF COMPUTERS USED IN SIMULATIONS

Our simulations are written in Julia language with a multi-threaded calculation of immediate reward. We used 4 computers with the following characteristics:

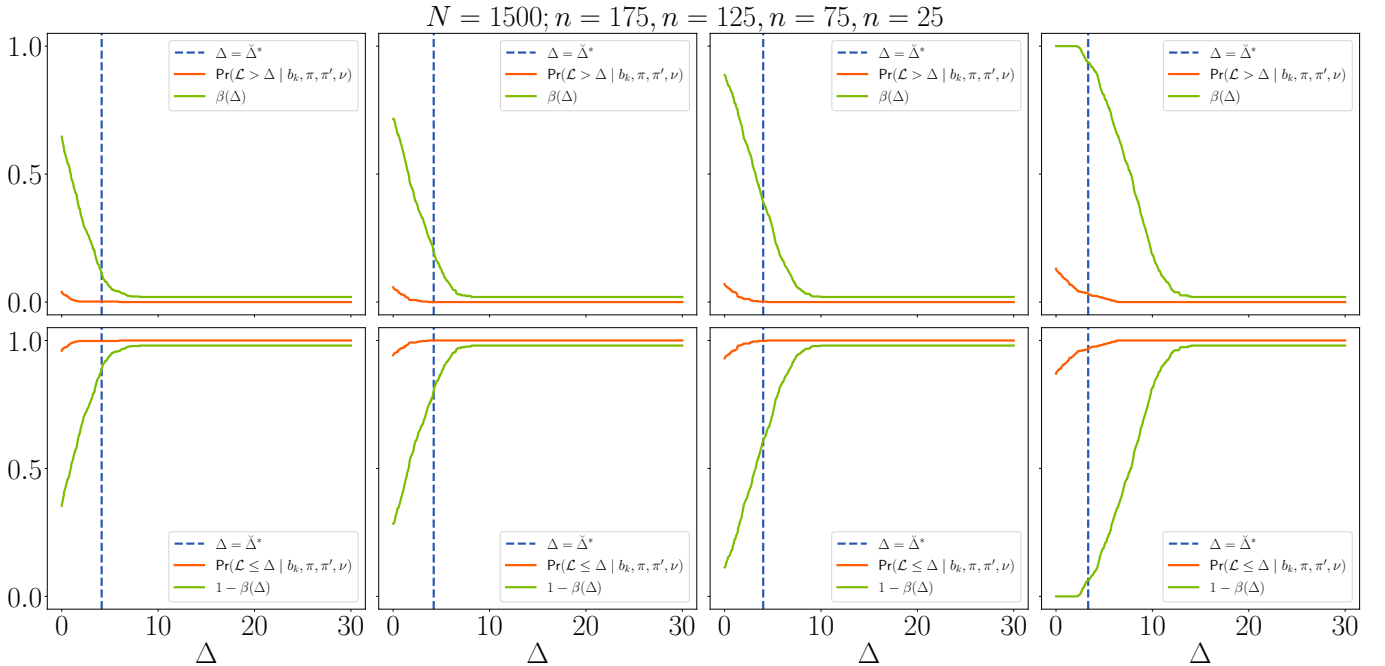
- 1) 40 cores Intel(R) Xeon(R) E5-2670 v2 with 256 GB of RAM working at 2.50GHz;
- 2) 24 cores Intel(R) Core(TM) i9-7920X with 64 GB of RAM working at 2.90GHz;
- 3) 20 cores Intel(R) Xeon(R) E5-2630 v4 with 64 GB of RAM working at 2.20GHz;
- 4) 20 cores Intel(R) Core(TM) i9-9820X with 64 GB of RAM working at 3.30GHz.



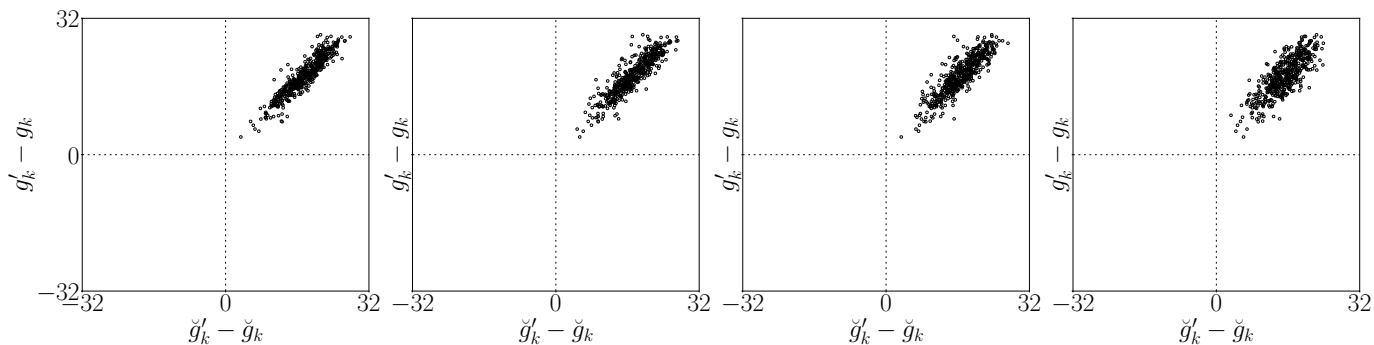
**Fig. 13:** Results for scenario 1 - probabilistic action consistency: Histograms of  $P_{\text{Loss}}$  and  $Pb_{\text{Loss}}$  for  $N = 1500$ ,  $\alpha = 0.01$ ,  $z_{\alpha/2} = 2.56$ , bin width is 1.0; from the left to the right  $n = 175$ ,  $n = 125$ ,  $n = 75$ ,  $n = 25$ .

	$n = 100$	$n = 75$	$n = 50$	$n = 25$
$g_k$ time [sec]	36745	45187	44899	30889
$\check{g}_k$ and $l, u$ time [sec]	17361	12546	4388	844
$\check{g}_k$ time [sec]	363	247	65	14
$l, u$ time [sec]	16998	12299	4323	830

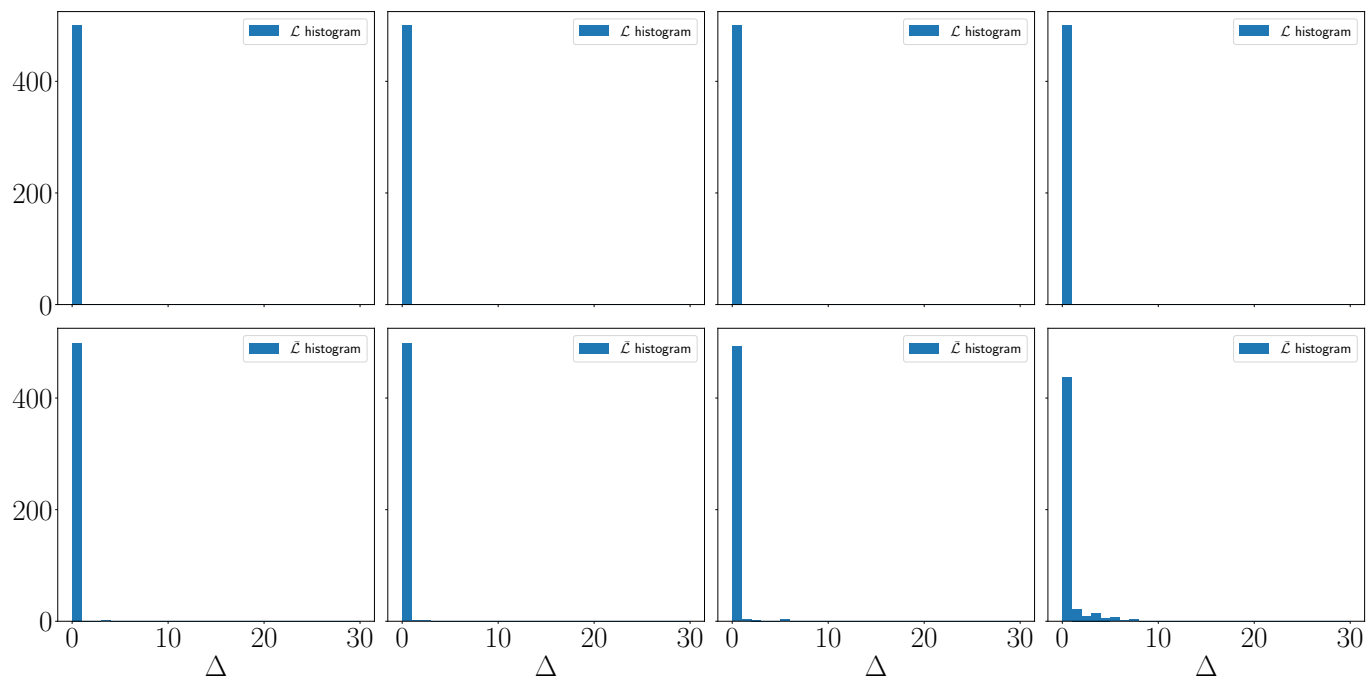
**TABLE IV:** Results for scenario 2 - empirical absolute action consistency: run times for  $N = 1000$ .



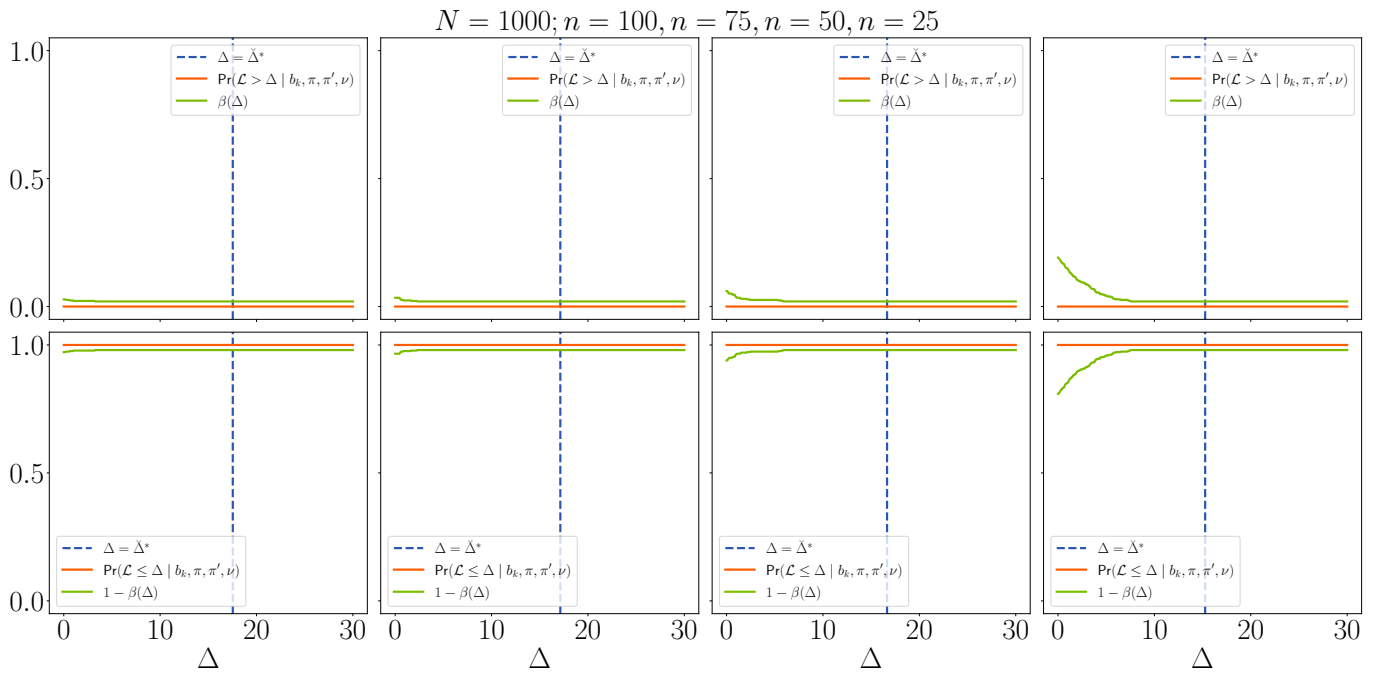
**Fig. 14:** Results for scenario 1 - probabilistic action consistency: Empirical characterization for  $N = 1500$ ,  $\alpha = 0.01$ ,  $z_{\alpha/2} = 2.56$ , evaluated in a grid with intervals 0.001; from the left to the right  $n = 175$ ,  $n = 125$ ,  $n = 75$ ,  $n = 25$ .



**Fig. 15:** Results for scenario 2 - empirical absolute action consistency: We demonstrate from the left to the right action consistency of the samples of the returns for  $n = 100, n = 75, n = 50, n = 25$ , whereas  $N = 1000$ . As we see, all the samples are action consistent.



**Fig. 16:** Results for scenario 2 - empirical absolute action consistency: Histograms of  $\text{PLoss}$  and  $\text{PbLoss}$  for  $N = 1000$ ,  $\alpha = 0.01$ ,  $z_{\alpha/2} = 2.56$ , bin width is 1.0; from the left to the right  $n = 100, n = 75, n = 50, n = 25$ .



**Fig. 17:** Results for scenario 2 - empirical absolute action consistency: Empirical characterization for  $N = 1000$ ,  $\alpha = 0.01$ ,  $z_{\alpha/2} = 2.56$ , evaluated in a grid with intervals 0.001; from the left to the right  $n = 100, n = 75, n = 50, n = 25$ .