# Voronoi Progressive Widening: Efficient Online Solvers for Continuous State, Action, and Observation POMDPs

Michael H. Lim, Claire J. Tomlin and Zachary N. Sunberg

*Abstract*— This paper introduces Voronoi Progressive Widening (VPW), a generalization of Voronoi optimistic optimization (VOO) and action progressive widening to partially observable Markov decision processes (POMDPs). Tree search algorithms can use VPW to effectively handle continuous or hybrid action spaces by efficiently balancing local and global action searching. This paper proposes two VPW-based algorithms and analyzes them from theoretical and simulation perspectives. Voronoi Optimistic Weighted Sparse Sampling (VOWSS) is a theoretical tool that justifies VPW-based online solvers, and it is the first algorithm with global convergence guarantees for continuous state, action, and observation POMDPs. Voronoi Optimistic Monte Carlo Planning with Observation Weighting (VOMCPOW) is a versatile and efficient algorithm that consistently outperforms state-of-the-art POMDP algorithms in several simulation experiments.

## I. INTRODUCTION

The partially observable Markov decision process (POMDP) is a flexible mathematical framework for expressing stochastic sequential decision problems. POMDPs can represent a wide range of real world problems such as autonomous driving [1], [2], cancer screening [3], spoken dialog systems [4], and aircraft collision avoidance [5]. A POMDP is an optimization problem in which we aim to find a policy that maps states to actions which will control the state to maximize the expected sum of rewards. Finding an optimal POMDP policy is computationally demanding because of the uncertainty introduced by imperfect observations [6]. One of the most popular approaches to deal with this computational challenge is to use *online* algorithms that look for local approximate policies as the agent interacts with the environment rather than computing a global policy that maps every possible outcome to an action. Many online POMDP algorithms are variants of Monte Carlo tree search (MCTS) [7], [8], [9] or other tree search variants [10], [11].

The research presented here concerns *continuous space POMDPs*, defined as POMDPs with continuous state, action and observation spaces. If tree search is naively applied to continuous space POMDPs, an immediate problem is that the branching factor of the tree will be infinite, preventing evaluation of future consequences of actions deep in the tree. This problem has been thoroughly studied, and the most popular solutions are sparse sampling [15], [13], [16] and progressive widening (PW) [17], [9], which limit the branching factor by only considering a randomly-selected subset of the states, actions and observations. For continuous

Michael H. Lim and Claire J. Tomlin are with the Electrical Engineering and Computer Sciences department, University of California, Berkeley. Zachary N. Sunberg is with the Aerospace Engineering Sciences department, University of Colorado Boulder.
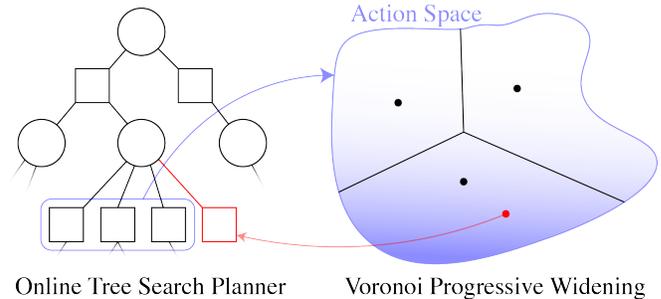
Fig. 1. Voronoi progressive widening (VPW) guides the selection of actions to widen the search tree. Circles denote state/belief nodes, squares denote action nodes.

states and observations, this random sampling is used to approximate the expectation of the next-step value in Bellman's equation. Since the Monte Carlo integration needed to approximate these expectations has straightforward and robust convergence guarantees, sparse sampling and PW are sufficient.

Continuous action spaces are, however, much more difficult to accommodate. Instead of simply estimating expectations, planning in a problem with a continuous action space requires solving a nonconvex optimization problem at each node in the tree. PW has been applied to these problems [17], [9], but, since PW considers a randomly sampled set of actions, it wastes a large amount of computation on suboptimal parts of the action space and has only been demonstrated on small action spaces such as one-dimensional intervals of real numbers. Voronoi optimistic optimization (VOO) [12] is an approach for sequential decision problems with deterministic and fully-observable dynamics that attempts to focus computation on more promising parts of the action space by partitioning it into Voronoi cells and sampling from cells corresponding to previously successful actions.

In this work, we propose Voronoi Progressive Widening (VPW), a versatile technique to modify tree search algorithms to effectively handle continuous or hybrid action spaces. VPW generalizes VOO and PW to problems with both transition and observation uncertainties. Like PW, it balances exploring previously proposed actions and searching for new action candidates. However, VPW searches for new candidates in a much more efficient way via VOO, and balances local and global searching without relying on any additional prior information or expert knowledge. Furthermore, VPW does not require significant additional computation time and can handle both continuous and hybrid

| Solver | Problem | Continuous Space | Source | Brief Description |
|---|---|---|---|---|
| VOOT | MDP | $S,A$ | [12] | VOO for deterministic MDPs; optimality guarantees; practical. |
| POWSS | POMDP | $S,O$ | [13] | Sparse sampling with observation weights; optimality guarantees. |
| **VOWSS** | POMDP | $S,A,O$ | **New** | Extends POWSS with VPW; optimality guarantees. |
| POMCPOW | POMDP | $S,A,O$ | [9] | Combines POMCP and DPW; practical. |
| **VOMCPOW** | POMDP | $S,A,O$ | **New** | Extends POMCPOW with VPW; practical. |
| BOMCP | POMDP | $S,A,O$ | [14] | Extends POMCPOW with Gaussian processes; practical. |
| VOSS | MDP | $S,A$ | Appendix | Extends sparse sampling with VPW; optimality guarantees. |

action spaces with relative ease.

There are two main contributions in this work. First, we prove theoretical guarantees about certain parts of the VPW approach via Voronoi Optimistic Weighted Sparse Sampling (VOWSS). VOWSS is special case of VPW applied to Partially Observable Weighted Sparse Sampling (POWSS) [13] with fixed number of action samples, which integrates the two convergence guarantees from VOO and POWSS into a single guarantee. Specifically, this is the first known solver to have global convergence guarantees for continuous space POMDPs. This shows that VOO and action PW were extended in a way that VPW not only can handle transition and observation uncertainties, but also combines the two procedures in a theoretically sound way.

Second, we propose an efficient VPW-based algorithm: Voronoi Optimistic Monte Carlo Planning with Observation Weighting (VOMCPOW). VOMCPOW can also be thought of as a practical extension of VOWSS. The experiments over different domains show the practical effectiveness of VPW for handling continuous and hybrid action space problems. VOMCPOW consistently outperforms state-of-the-art continuous space POMDP solvers based on PW, such as Partially Observable Monte Carlo Planning with Observation Widening (POMCPOW) [9] and Bayesian Optimized Monte Carlo Planning (BOMCP) [14] in several simulation experiments. Table I summarizes the solvers that are studied in this paper.

The remainder of this paper proceeds as follows: First, Sections II and III review preliminary definitions and previous work. Then, Section IV introduces the VPW algorithm, and describes its strengths as well as how to efficiently implement it in practice. Section V presents the VOWSS algorithm and theoretical analysis justifying VPW. Finally, Section VI empirically shows the optimality of VOWSS action selection and the efficiency and robustness of VOMCPOW over three different simulation experiments.

## II. PRELIMINARIES

*1) POMDP Formulation:* A POMDP is defined by a 7-tuple $(S,A,O,\mathcal{T},\mathcal{Z},R,\gamma)$: $S$ is the state space, $A$ is the action space, $O$ is the observation space, $\mathcal{T}$ is the transition density $\mathcal{T}(s'|s,a)$, $\mathcal{Z}$ is the observation density $\mathcal{Z}(o|a,s')$, $R$ is the reward function, and $\gamma \in [0,1)$ is the discount factor [18], [19]. For POMDPs, since the agent receives only observations, the agent can infer the state by maintaining a belief $b_t$ at each step $t$ and updating it with new action

and observation pair $(a_{t+1},o_{t+1})$ via Bayesian filtering [20]. A policy, denoted with $\pi$, maps beliefs $b_t$ to actions $a_t$. The agent seeks an optimal policy, $\pi^*$ that maximizes the expected cumulative reward.

For simplicity, in the theoretical portion of this paper, we focus on problems with a finite horizon length $D$, though the basic concepts empirically work well in the infinite horizon case. The state value function $V$ and action value function $Q$ for a given belief $b$ and policy $\pi$ at step $t$ by Bellman updates for $t \in [0,D-1]$ are defined as follows:

$$V_t^\pi(b) = \mathbb{E}[\sum_{i=t}^{D-1} \gamma^{i-t} R(s_i, \pi(s_i))|b], \ V_D^\pi(b) = 0, \quad (1)$$

$$Q_t^\pi(b,a) = \mathbb{E}\left[R(s,a) + \gamma V_{t+1}^\pi(bao)|b\right], \quad (2)$$

where $bao$ indicates the belief $b$ updated with $(a,o)$. The optimal value functions at each depth $t$ should satisfy

$$V_t^\star(b) = \max_{a \in A} Q_t^\star(b,a), \ \pi_t^\star(b) = \arg\max_{a \in A} Q_t^\star(b,a), \quad (3)$$

$$Q_t^\star(b,a) = \mathbb{E}\left[R(s,a) + \gamma V_{t+1}^\star(bao)|b\right]. \quad (4)$$

*2) Generative models:* In some cases, it is not necessary to explicitly evaluate the probability density of the transition or observation distributions, and merely generating samples is sufficient. In this work, we use a *generative model G* that generates state, reward, and observation samples.

## III. ADDITIONAL RELATED WORK

This section expands on the introduction to cover more previous work in solving continuous action MDPs and POMDPs with online tree search. Several techniques use double progressive widening (DPW) [17], originally designed to solve continuous space MDPs. Most notably, POMCPOW and PFT-DPW [9] extend POMCP and DPW to handle continuous space POMDPs. However, these algorithms use DPW as it is proposed with the inefficient action sampling that is not suitable for large problems.

One effective direction to handle continuous action spaces has been to use continuous bandits to sample new actions. Particularly, hierarchical optimistic optimization (HOO) and the corresponding HOOT algorithm [21] began work on continuous bandit algorithms applied to MCTS. This is followed by works such as HOLOP [22] that plans with open-loop trajectories instead of individual actions using HOO, POLY-HOOT [23] with polynomial guarantees, and VOO and VOOT [12] that we build on in this work.

**Algorithm 1** VOO Algorithms [12]

**Global Variables:** $A, D(\cdot,\cdot), \omega, \Sigma$.
**Algorithm:** VOO($\mathbf{a}, \mathbf{Q}$)
**Input:** Array of Voronoi centers $\mathbf{a} = [a_i]$ and their function values $\mathbf{Q} = [Q(a_i)]$.
**Output:** A sample $a$.

1: $u \leftarrow \text{Unif}[0,1]$
2: **if** $u \leq \omega$ or $|\mathbf{Q}| = 0$ **then**
3:     $a \leftarrow \text{Unif}(A)$
4: **else**
5:     $a \leftarrow \text{BESTVORONOICELL}(\mathbf{a}, \mathbf{Q})$
6: **return** $a$

**Algorithm:** BESTVORONOICELL($\mathbf{a}, \mathbf{Q}$)
**Input:** Array of Voronoi centers $\mathbf{a} = [a_i]$ and their function values $\mathbf{Q} = [Q(a_i)]$.
**Output:** A sample $a$.

1: $a^* \leftarrow \arg\max_{\mathbf{a}} \mathbf{Q}$
2: **while** $a$ is not closest **do**
3:     $a \leftarrow N(a^*, \Sigma)$
4:     Check if $D(a, a^*) \leq D(a, a_i) \, \forall a_i \in \mathbf{a}, \, a_i \neq a^*$
5: **return** $a$

---

**Algorithm 2** Action Progressive Widening [17]

**Global Variables:** $k_a, \alpha_a, A, c$.
**Algorithm:** PW($h$)
**Input:** Belief/history node $h$ in the MCTS tree.
**Output:** An action $a$.

1: **if** $|C(h)| = 0$ or $|C(h)| \leq k_a N(h)^{\alpha_a}$ **then**
2:     $a \leftarrow \text{Unif}(A)$
3:     $C(h) \leftarrow C(h) \cup \{a\}$
4: **else**
5:     $a \leftarrow \arg\max_{a \in C(h)} Q(h,a) + c\sqrt{\frac{\log N(h)}{N(h,a)}}$
6: **return** $a$

---

Another direction is to use Bayesian optimization to efficiently sample new actions. CBTS [24] uses Gaussian processes (GP) to tackle the random action sampling problem. However, CBTS does not use progressive widening approach and uses a separate GP at each belief node, which limits its optimization scope and branching capabilities. Most recently, BOMCP [14] extends POMCPOW by posing the random action sampling step of DPW as a Bayesian optimization problem over all the belief and action nodes. Despite the effectiveness of BOMCP, GPs are very computationally expensive to fit especially over the joint belief and action space, and BOMCP trades off computation time for sample efficiency compared to POMCPOW.

Other directions include GPS-ABT [25] that uses generalized pattern search to find local optima, PA-POMCPOW that additionally incorporates score functions [26], and VG-UCT that calculates gradient of the value function [27].

## IV. VORONOI PROGRESSIVE WIDENING

In this section, we first introduce VOO and PW to motivate the formulation of VPW and VOMCPOW.

### A. Voronoi Optimistic Optimization

Voronoi optimistic optimization (VOO) [12] is a continuous multi-armed bandit algorithm that adaptively explores the sampling space by partitioning the space into Voronoi cells. We define a Voronoi cell by the set of points closest to the corresponding Voronoi center compared to the other centers, and the best Voronoi cell is the Voronoi cell with the center that achieves the highest function value estimate. In our setting, the sampling space is the action space, the Voronoi centers are the sampled actions, and the function values are the $Q$-value estimates at the actions. Since Voronoi cells are solely defined by distances to Voronoi centers, VOO

additionally takes in a distance metric $D(\cdot,\cdot)$ as an input, and scales well to higher dimensions unlike HOO.

In Algorithm 1, we define the adapted functions VOO and BESTVORONOICELL. VOO searches the action space either uniformly with probability $\omega$ or from the best Voronoi cell with probability $1 - \omega$. BESTVORONOICELL samples an action from the best Voronoi cell via rejection sampling. In practice, VOO sampling is best implemented by using Gaussian rejection sampling centered around the current best action [12]. While we reflected this in our algorithm definition, we note that the theoretical guarantees only hold for uniform rejection sampling. Furthermore, to improve computation time for the experiments in Section VI, we limit the maximum number of rejected samples via methods described in Appendix B. VOO has previously been extended to Voronoi Optimistic Optimization Tree (VOOT) [12], but only for deterministic MDPs.

### B. Progressive Widening

Progressive widening (PW) [17] tackles continuous state and action spaces by gradually expanding the state and action sample sets. PW limits the number of sampled children to $k \cdot N^\alpha$, where $N$ corresponds to the number of visits to the parent node in the search tree, and $k, \alpha$ are the widening parameters. Algorithm 2 describes action progressive widening. Here, $h$ is the belief/history node in the MCTS tree, $C(h)$ the list of children action nodes, $N$ the number of visits, and $c$ the Upper Confidence Bound exploration parameter. In the action space, PW balances two modes of exploration: (1) exploring branches with previously proposed actions (line 2), and (2) searching for new suitable action candidates (line 5). However, since PW samples new actions uniformly from the action space, it is rather sample inefficient and does not take into account any information gained from previous samples.

### C. Voronoi Progressive Widening

We now introduce Voronoi Progressive Widening (VPW) in Algorithm 3, which generalizes VOO and PW. Essentially, VPW progressively widens with a VOO sample instead of a uniform sample. With this formulation, PW is simply VPW with $\omega = 1$, and VOO is VPW with $k_a \cdot N^{\alpha_a} = +\infty$.

Compared to PW, VPW only additionally relies on having a collection of $Q$-value estimates, which are typically calculated and stored for MCTS solvers. While this means that we require the $Q$-value estimates to be relatively faithful to the

**Algorithm 3** Voronoi Progressive Widening

**Global Variables:** $k_a, \alpha_a, \omega, A, c$.
**Algorithm:** VPW($h$)
**Input:** Belief/history node $h$ in the MCTS tree.
**Output:** An action $a$.

1: **if** $|C(h)| = 0$ or $|C(h)| \leq k_a N(h)^{\alpha_a}$ **then**
2:     $a \leftarrow \text{VOO}(C(h), [Q(h, \cdot)])$
3:     $C(h) \leftarrow C(h) \cup \{a\}$
4: **else**
5:     $a \leftarrow \arg\max_{a \in C(h)} Q(h, a) + c\sqrt{\frac{\log N(h)}{N(h, a)}}$
6: **return** $a$

---

actual $Q$-values and also requires an informative rollout policy that can guide the solver to more optimal actions, this is not a requirement specific to VPW as a typical MCTS solver should aim to have both of these components. Furthermore, since the VOO step does not require a significant additional time, VPW operates in a similar time scale as PW/DPW while being able to efficiently sample action candidates that are closer to the optimal action. The only other requirement of VPW is the distance metric on the action space, but this is often simple to define, and in fact lets VPW to tackle hybrid action spaces such as the one in Section VI-B.

### D. VOMCPOW

The VPW technique is versatile as it can be applied to any MCTS solver for MDPs and POMDPs without requiring any additional expert or domain knowledge. For instance, one could consider Sparse-UCT-VPW algorithm for continuous space MDPs that uses VPW criterion with the Sparse-UCT algorithm [28] to handle continuous action spaces. Specifically, we demonstrate the versatility of VPW through modifying POMCPOW into VOMCPOW (Voronoi Optimistic Monte Carlo Planning with Observation Weighting) by simply swapping out action PW with VPW. VOMCPOW consistently outperforms POMCPOW and BOMCP by a statistically significant margin across all of our experiments.

### V. CONVERGENCE ANALYSIS

For convergence analysis, we introduce and Voronoi Optimistic Weighted Sparse Sampling (VOWSS), a VPW-based continuous space online POMDP solver that guarantees convergence to an optimal policy. VOWSS justifies the usage of VPW-based solvers that build upon sparse sampling [15], particle weighting and VOO, such as VOMCPOW.

### A. VOWSS

We define ESTIMATEV and ESTIMATEQ functions in Algorithm 4. Here, $C_s, C_a$ are the state and action widths, respectively. ESTIMATEV is a subroutine that returns the value function, $V$, for an estimated state or belief, by calling ESTIMATEQ for each action and returning the maximum. Similarly, ESTIMATEQ performs sampling and recursive calls to ESTIMATEV to estimate the $Q$-function at a given step with a weighted average in the style of POWSS [13]. The sampled states $s_i'$ are inserted into each next-step belief particle set $\overline{bao}_j$ with the new weights $w_i' = w_i \cdot \mathcal{Z}(o_j | a, s_i')$,

---

**Algorithm 4** Value estimation algorithms for VOWSS

**Global Variables:** $\gamma, G, C_s, C_a, D$.
**Algorithm:** ESTIMATEV($\bar{b}, d$)
**Input:** Belief particle set $\bar{b} = \{(s_i, w_i)\}$, depth $d$.
**Output:** A scalar $\hat{V}_d^\star(\bar{b})$ that is an estimate of $V_d^\star(\bar{b})$.

1: If $d \geq D$ the max depth, then return 0.
2: For $i = 1, \cdots, C_a$, sequentially run VPW($\bar{b}$) to sample actions $a_i$ and estimate the corresponding $Q$-values with $\hat{Q}_d^\star(\bar{b}, a_i) = \text{ESTIMATEQ}(\bar{b}, a, d)$.
3: Return the value function:

$$\hat{V}_d^\star(\bar{b}) = \max_{i=1,\cdots,C_a} \hat{Q}_d^\star(\bar{b}, a_i).$$

**Algorithm:** ESTIMATEQ($\bar{b}, a, d$)
**Input:** Belief particle set $\bar{b} = \{(s_i, w_i)\}$, action $a$, depth $d$.
**Output:** A scalar $\hat{Q}_d^\star(\bar{b}, a)$ that is an estimate of $Q_d^\star(b, a)$.

1: For each particle-weight pair $(s_i, w_i)$ in $\bar{b}$, generate $s_i', o_i, r_i$ from $G(s_i, a)$.
2: For each observation $o_j$ from previous step, iterate over $i = 1, \cdots, C_s$ to insert $(s_i', w_i \cdot \mathcal{Z}(o_j | a, s_i'))$ to a new belief particle set $\overline{bao}_j$.
3: Return the $Q$-value estimate:

$$\hat{Q}_d^\star(\bar{b}, a) = \frac{\sum_{i=1}^{C_s} w_i(r_i + \gamma \cdot \text{ESTIMATEV}(\overline{bao}_i, d+1))}{\sum_{i=1}^{C_s} w_i}.$$

---

which are the adjusted probability of hypothetically sampling observation $o_j$ from state $s_i'$.

Thus, the VOWSS policy action is obtained by calling ESTIMATEV($\bar{b}_0, 0$) at the root node and taking the action that corresponds to the maximizing $Q$-value. Note that the particle belief set is initialized by drawing samples from $b_0$ and normalizing weights to $1/C_s$. Like POWSS, VOWSS is not very computationally efficient as it fully expands the sparsely sampled tree. Thus, VOWSS mainly demonstrates theoretical convergence and are only practically applicable to very small problems. The algorithms in this section are written in a style closer to a plain English format to enhance the readability of mathematical machinery used for these algorithms, since VOWSS is more theoretical in nature.

It is worth noting that with our recursive formulation, VOOT [12] can be recast as a $Q$-function estimation algorithm that simply returns the one sample $Q$-function estimate. Particularly, this means that ESTIMATEQ for VOOT would look identical to VOWSS's ESTIMATEQ function except with $C_s$ set to 1 as VOOT assumes no transition and observation uncertainties, with no action width decay. Consequently, VOOT is notably faster than VOWSS and thus practical for deterministic MDP problems, since it only needs to sample the deterministic next step state once. Nevertheless, due to the hierarchical structure of VOOT, the actual VOO algorithm as well as the theoretical guarantees of VOOT should remain identical when recast into this recursive form. We also discuss the stochastic MDP case in Appendix C.

## B. Convergence Guarantees

We obtain global convergence to the optimal value functions for continuous space POMDPs by combining VOO that globally optimizes over the entire action space and POWSS that estimates value functions with arbitrary precision. To our knowledge, VOWSS is the first continuous space POMDP algorithm to have global convergence guarantees without relying on any discretization schemes for any of the spaces. In the subsequent sections, we prove that VOWSS policy can be made to perform arbitrarily close to the optimal policy by increasing both the state and action widths. Our results are derived with $k_a \cdot N^{\alpha_a} = +\infty$ (VPW = VOO), and uniform rejection sampling for BESTVORONOICELL.

The convergence proof of VOWSS relies on the proofs for the ancestor algorithms: POWSS and VOO. Since the proof of VOWSS requires the union of all the regularity conditions for POWSS and VOO, which deal with local smoothness of rewards and observation density requirements, we will not repeat them here and detail them in Appendix A.2. One of the conditions requires that the reward function $R$ be bounded and measurable, $||R||_\infty \leq R_{\max}$, and consequently the value is bounded by $V_{\max} \equiv R_{\max}/(1-\gamma)$ for $\gamma < 1$.

*Theorem 1 (VOWSS Inequality):* Suppose we choose the action sampling width $C_a$ and state sampling width $C_s$ such that under the union of regularity conditions specified by [13] and [12], the intermediate POWSS bounds and VOO bounds in Lemma 1 are satisfied at every depth of the tree. Then, the following bounds for the VOWSS estimator $\hat{V}_{\text{VOWSS},d}(\bar{b})$ hold for all $d \in [0, D-1]$ in expectation:

$$|V_d^\star(b) - \hat{V}_{\text{VOWSS},d}(\bar{b})| \leq \eta_{C_a} + \alpha_{C_s}. \tag{5}$$

Here, $\eta_{C_a}$ and $\alpha_{C_s}$ are non-negative bounds that tend to 0 as we increase $C_a$ and $C_s$, respectively. We show a brief outline of the proof of Theorem 1. In order to prove Theorem 1, we first need to prove the intermediate bounds.

*Lemma 1 (VOWSS Intermediate Inequality):* Suppose with our notation, the POWSS estimators at all depths $d$ are within $\varepsilon$ of their mean values with probability $1 - p$, and the VOO agents have regret bounds of $\mathscr{R}_{C_a}$. The following inequalities hold for all $d \in [0, D-1]$ in expectation:

$$|V_d^\star(b) - \hat{V}_d^{C_a}(b)| \leq \eta_{C_a}(d), \tag{6}$$

$$|\hat{V}_d^{C_a}(b) - \hat{V}_{\text{VOWSS},d}(\bar{b})| \leq \alpha_{C_s}(d). \tag{7}$$

In Lemma 1, we aim to bound the inequality in (5) by applying triangle inequality to the two split terms: the VOO-like regret bound with $\eta_{C_a}(d)$ and the POWSS-like concentration bound with $\alpha_{C_s}(d)$. The exact definitions of these intermediate bound terms are defined and explained further in Appendix A.2. We define each value function used in Lemma 1 as the following:

$$V_d^\star(b) \equiv \max_{a \in A} R(b,a) + \gamma \mathbb{E}[V_{d+1}^\star(bao)|b],$$

$$\hat{V}_d^{C_a}(b) \equiv \max_{a \in VOO(A)} R(b,a) + \gamma \mathbb{E}[\hat{V}_{d+1}^{C_a}(bao)|b], \tag{8}$$

$$\hat{V}_{\text{VOWSS},d}(\bar{b}) \equiv \max_{a \in VOO(A)} \frac{\sum_{i=1}^{C_s} w_i(r_i + \gamma \cdot \hat{V}_{\text{VOWSS},d+1}(\overline{bao_i}))}{\sum_{i=1}^{C_s} w_i}.$$

Here, $V_d^\star(b)$ is the optimal value function as per the conventional definition, and $\hat{V}_{\text{VOWSS},d}(\bar{b})$ is the VOWSS estimator for the optimal value function. In addition, we introduce the theoretical intermediate term $\hat{V}_d^{C_a}(b)$ that bridges the gap between the regret bound and the concentration bound by only performing the VOO step and not the sparse sampling step. Using this intermediate term, we can further subdivide the intermediate inequalities to obtain the appropriate regret and concentration bounds using results from [12] and [13], respectively. Proving Lemma 1 proves Theorem 1.

## VI. EXPERIMENTS

The numerical experiments in this section confirm the theoretical results of Section V-B, as well as showcase the effectiveness of VPW. We demonstrate the performances of our algorithms in three different control experiments: Linear-Quadratic-Gaussian (LQG) control, Van Der Pol Tag, and modified lunar lander problems. When running experiments for POMCPOW, VOMCPOW, and BOMCP, we use the rollout policy heuristic to estimate the value function as well as take the first action to be the rollout policy action at each newly generated belief node. To determine the hyperparameters of the solvers, we used cross-entropy method (CEM) [29] to maximize the mean expected reward of a hyperparameter set when the optimal hyperparameters are not already available from previous experiments. The VOMCPOW and BOMCP hyperparameters were trained by first initializing them with POMCPOW hyperparameters and then using CEM again including all the additional hyperparameters specific to the solvers. The only hyperparameter that was manually chosen is the Gaussian covariance matrix for VOO rejection sampling. The code for the experiments is built on the POMDPs.jl framework [30]. The exact hyperparameters used for each experiment are given in Appendix B.

## A. Linear-Quadratic-Gaussian Control

We first test our algorithms on a simple 2D-action space LQG control system. Here, we choose a relatively simple problem setup such that we can easily understand and visualize the results while allowing VOWSS to still plan within a reasonable time. Here, VOWSS uses VPW like VOO by setting $k_a \cdot N^{\alpha_a} = +\infty$, but with Gaussian rejection sampling. The dynamics and observation models are

$$x_{t+1} = x_t + u_t + v_t; \quad y_t = x_t + w_t; \quad v_t, w_t \overset{i.i.d.}{\sim} N(0, \sigma^2 I). \tag{9}$$

The initial state $x_0$ is distributed as $N([-10, 10], \sigma^2 I)$, and $\sigma = 0.1$ for all $x_0, v_t, w_t$. We aim to minimize the cost function $J(x_0)$, while planning for two steps ($N = 2$):

$$J(x_0) = \mathbb{E}[x_N^T x_N + \sum_{t=0}^{N-1} (x_i^T x_i + u_i^T u_i)]. \tag{10}$$

The analytical answer can be obtained by solving the LQG backup equation, which is identical to the LQR solution,

$$u_0^\star = -K_0 \hat{x}_0 = -0.6 \cdot \hat{x}_0 = [6.0, -6.0]. \tag{11}$$

Since the problem has an analytical solution, we can directly compare the actions each solver chooses to the
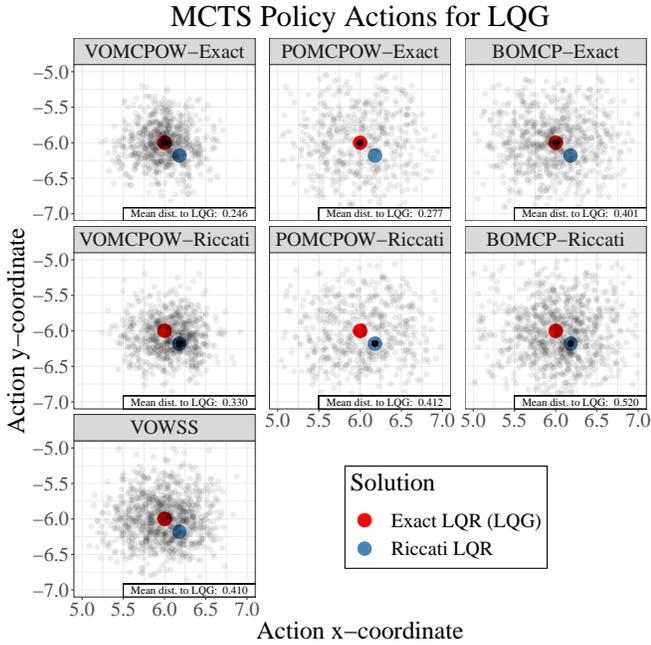
Fig. 2. Scatter plots of the first actions chosen by policies from different solvers and rollouts for the two step LQG problem. The exact LQR solution, which is the same as the LQG solution, is shown as a red dot at $[6.0, -6.0]$, and the Riccati solution as a blue dot at approximately $[6.18, -6.18]$.
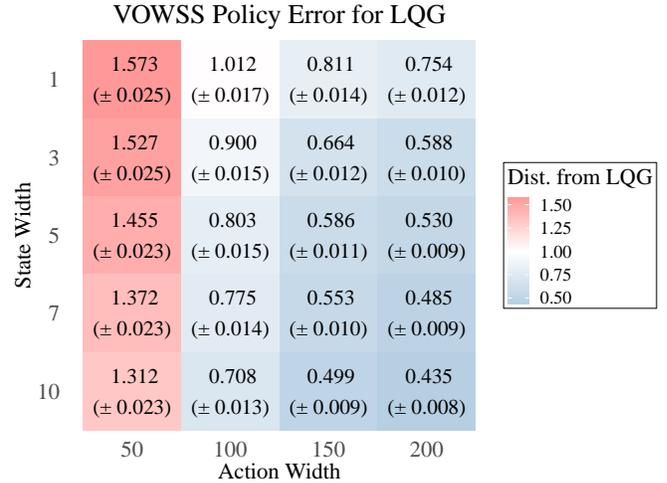


Fig. 3. Tabular summary of the first actions chosen by VOWSS policy, where we show the mean Euclidean distance from the LQG solution, and the corresponding standard error in parentheses.

analytical solution, for VOWSS, POMCPOW, VOMCPOW and BOMCP. For the solvers that require a rollout policy, we test both the finite horizon LQR solution (referred to as "exact policy") and the steady-state solution to the discrete time algebraic Riccati equation (referred to as "Riccati policy") [31] as the rollout policies to observe effects on each solver. For this particular scenario, the Riccati solution is $u_t^* \approx -0.618 \cdot \hat{x}_t$, and $u_0^* \approx [6.18, -6.18]$, slightly different from the exact policy.

When choosing hyperparameters through CEM, we use the Riccati policy as the rollout policy. The number of queries for POMCPOW and VOMCPOW are set to 1000 and BOMCP to 100, which correspond to approximately 0.1 seconds of planning time. For VOWSS, the state width is set to 10 and action width to 200. Furthermore, we introduce the action width decay term $\gamma_a$ such that the number of VPW iterations for a given height is $\gamma_a$ times the VPW iterations for the previous height, similar to the way VOOT is modified for long horizons [12]. The action width decay is manually set to $\gamma_a = 0.4$ after a careful inspection to balance performance and run time, which means that the action width for first step is 200, and the action width for second step is 80. While this results in number of evaluations on the order of $10^6$, we include the performance of VOWSS as a reference.

We first show the scatter plots of the actions chosen by each solver in Figure 2. Each scatter plot shows the result of 1000 simulations. While POMCPOW and BOMCP find solutions biased towards the rollout actions, we see that the effect is much less pronounced in VOMCPOW. Specifically, when the rollout policy is set to Riccati policy, both POMCPOW and BOMCP heavily resort to picking the Riccati solution as

shown by the large point mass corresponding to the Riccati solution in the scatter plots. Even though VOMCPOW still produces a sizable mass of points on and surrounding the Riccati solution, it otherwise picks solutions that are close to the LQG solution. VOMCPOW scatter plots show that the bias and the variance of the optimal action estimates are noticeably smaller than those of POMCPOW and BOMCP. Although VOWSS cannot be directly compared to the other solvers since the number of iterations is not on the same order and it doesn't rely on a rollout policy, the overall shape of the scatter plot looks similar to those of VOMCPOW.

In addition, we study the effects of varying the state and action widths for VOWSS. Figure 3 shows the mean Euclidean distance to the LQG solution for different combinations of state and action widths. Each cell in the table shows the result of 1000 simulations, and the other hyperparameters are left intact. As we increase either the state or action width, the mean distance and the standard error decrease. This indicates that VOWSS chooses better actions by increasing the state and action widths, which confirms our theoretical results: larger state widths should increase value estimation accuracy and larger action widths should improve action optimization.

### B. Van Der Pol Tag

Next, we test POMCPOW and VOMCPOW on the Van Der Pol Tag (VDP Tag) problem introduced by [9], which is the only problem in their work with continuous state, observation and (hybrid) action spaces. In VDP Tag, an agent navigates through 2D box to tag a target with randomized initial state that follows a dynamics model defined by the Van Der Pol differential equation. The agent moves at a fixed speed, but can choose the direction of travel and can decide to make an accurate observation of where the target is at a higher cost than making the default noisy observation. While the target can freely move within the 2D box, the agent is blocked when it comes into contact with a barrier.
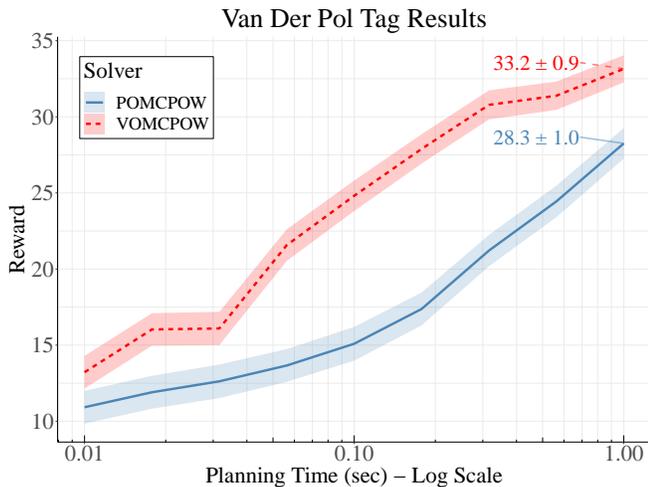
Fig. 4. Mean rewards for POMCPOW and VOMCPOW for VDP Tag. Ribbons indicate one standard error.



Fig. 5. Mean rewards for POMCPOW, VOMCPOW and BOMCP for the modified lunar lander problem. Ribbons indicate one standard error.

We show the mean reward for 1000 simulations for each solver and each planning time plotted in log scale of seconds in Figure 4. We observe that VOMCPOW outperforms POM-CPOW at every planning time by a statistically significant margin. It is also worth noting that VOMCPOW takes almost an order of magnitude less planning time to reach the mean reward of 25 compared to POMCPOW. While VDP Tag is a continuous space POMDP, the rewards are still discrete in both state and action spaces, which suggests that even with discrete jumps in the reward function, VPW can still optimize to find better actions. In addition, while it is possible to adapt BOMCP to solve VDP Tag, it requires nontrivial modifications to the action space due to the action space being hybrid ($A = [0, 2\pi] \times \{0, 1\}$) and the angle space being a modular space. This further illustrates the ease of implementing VPW since we only need to supply VPW with a distance metric on the action space, which makes VPW suitable for hybrid action spaces as well.

### C. Lunar Lander

Lastly, we test POMCPOW, VOMCPOW and BOMCP on the modified lunar lander problem proposed in [14], where the main objective is to guide a vehicle to land in a target zone safely. In this version of lunar lander, the vehicle state is defined as $(x, y, \theta, \dot{x}, \dot{y}, \omega)$, and we only obtain noisy observations of the angular rate, horizontal speed, and above-ground level. The action space is defined by the tuple $(T, F_x, \delta)$ where $T$ is the main vertical thrust, $F_x$ the corrective horizontal thrust, and $\delta$ the offset. The default rollout policy is proportional control based on the observation.

Once again, we show the mean reward for 1000 simulations for each solver and each planning time plotted in log scale of seconds in Figure 5. VOMCPOW outperforms both POMCPOW and BOMCP by a statistically significant margin once the planning time exceeds 0.1 seconds. Below 0.1 seconds, the performances of POMCPOW and VOMCPOW are similar due to the problem having a long horizon, in which the effects of VPW will not be as apparent since best
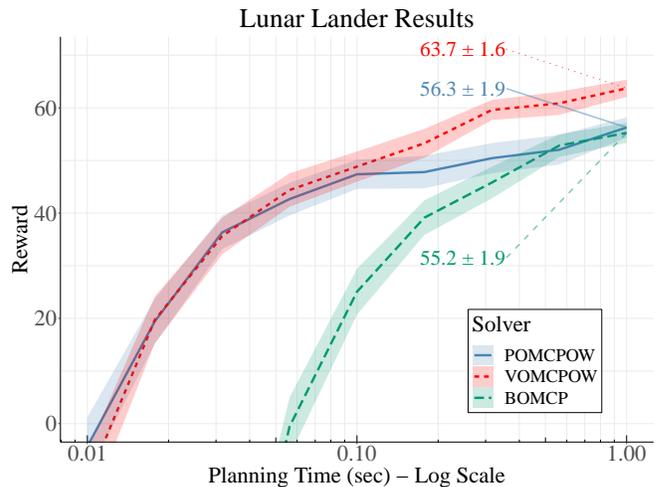
Voronoi cell actions will not be sampled very often.

We note that the performance of POMCPOW here seemingly differs from the results in [14], which might be due to a several factors. The lander problem has a high variance in the returns, since failure to land results in a steep penalty of -1000. This failure mode does not happen very often even within 1000 iterations, so the number of failures significantly impacts the mean rewards. Furthermore, we are comparing the algorithms based on planning time rather than number of queries. POMCPOW queries much faster than BOMCP, and average planning time may not sufficiently capture the relationship between total planning time and number of queries for a progressively built tree. It is also possible that authors of [14] did not set POMCPOW first action to be the rollout action.

## VII. CONCLUSION

In this paper, we have introduced Voronoi Progressive Widening (VPW), a versatile technique to effectively handle continuous or hybrid action spaces in MDPs and POMDPs. Consequently, we proposed two VPW-based tree search algorithms to demonstrate convergence guarantees and efficiency, justifying the theoretical soundness, versatility and practicality of VPW for many continuous or hybrid action space POMDPs.

This study has yielded a several key insights, which suggest some promising future directions to explore. While each of the VPW-based algorithms either enjoy theoretical guarantees or computational efficiency, the gap between theoretical soundness and practicality still remains. In particular, the state and action selection heuristics as it is utilized in POMCPOW still need to be integrated into the theoretical algorithms like POWSS to bridge the gap between theory and practice. On the other hand, we believe that the VPW technique itself should also have some theoretical guarantees similar to the DPW technique [17].

Additionally, since the convergence guarantees of VOO hold as long as the reward function is locally smooth around

at least one global optimum [12], more extensive study could be done for effectiveness of VPW on MDPs and POMDPs with discrete reward structures. Similar to how the Rényi divergence requirement for POWSS [13] could be a difficulty indicator for likelihood-weighted sparse tree solvers, the proof techniques utilized for VOO and DPW could offer insights on what continuous or hybrid action space problems are harder to solve than others.

## ACKNOWLEDGEMENTS

## REFERENCES

[1] H. Bai, S. Cai, N. Ye, D. Hsu, and W. S. Lee, "Intention-aware online POMDP planning for autonomous driving in a crowd," in *IEEE International Conference on Robotics and Automation*. Seattle, WA, USA: IEEE, 2015, pp. 454–460.

[2] Z. N. Sunberg, C. J. Ho, and M. J. Kochenderfer, "The value of inferring the internal state of traffic participants for autonomous freeway driving," in *American Control Conference*. Seattle, WA, USA: IEEE, 2017, pp. 3004–3010.

[3] T. Ayer, O. Alagoz, and N. K. Stout, "A POMDP approach to personalize mammography screening decisions," *Operations Research*, vol. 60, no. 5, pp. 1019–1034, 2012.

[4] S. Young, M. Gašić, B. Thomson, and J. D. Williams, "POMDP-based statistical spoken dialog systems: A review," *IEEE*, vol. 101, no. 5, pp. 1160–1179, 2013.

[5] J. E. Holland, M. J. Kochenderfer, and W. A. Olson, "Optimizing the next generation collision avoidance system for safe, suitable, and acceptable operational performance," *Air Traffic Control Quarterly*, vol. 21, no. 3, pp. 275–297, 2013.

[6] C. H. Papadimitriou and J. N. Tsitsiklis, "The complexity of Markov decision processes," *Mathematics of Operations Research*, vol. 12, no. 3, pp. 441–450, 1987.

[7] C. B. Browne, E. Powley, D. Whitehouse, S. M. Lucas, P. I. Cowling, P. Rohlfshagen, S. Tavener, D. Perez, S. Samothrakis, and S. Colton, "A survey of monte carlo tree search methods," *IEEE Transactions on Computational Intelligence and AI in Games*, vol. 4, no. 1, pp. 1–43, 2012.

[8] D. Silver and J. Veness, "Monte-Carlo planning in large POMDPs," in *Advances in Neural Information Processing Systems*. Curran Associates, Inc., 2010, pp. 2164–2172.

[9] Z. Sunberg and M. J. Kochenderfer, "Online algorithms for POMDPs with continuous state, action, and observation spaces," in *International Conference on Automated Planning and Scheduling*. Delft, Netherlands: AAAI Press, 2018.

[10] N. Ye, A. Somani, D. Hsu, and W. S. Lee, "DESPOT: Online POMDP planning with regularization," *Journal of Artificial Intelligence Research*, vol. 58, pp. 231–266, 2017.

[11] H. Kurniawati and V. Yadav, "An online POMDP solver for uncertainty planning in dynamic environment," in *Robotics Research*. Springer, 2016, pp. 611–629.

[12] B. Kim, K. Lee, S. Lim, L. Kaelbling, and T. Lozano-Perez, "Monte carlo tree search in continuous spaces using voronoi optimistic optimization with regret bounds," in *AAAI Conference on Artificial Intelligence*, vol. 34. AAAI Press, 2020, pp. 9916–9924.

[13] M. H. Lim, C. Tomlin, and Z. N. Sunberg, "Sparse tree search optimality guarantees in pomdps with continuous observation spaces," in *International Joint Conference on Artificial Intelligence, IJCAI-20*. International Joint Conferences on Artificial Intelligence, Inc., 7 2020, pp. 4135–4142.

[14] J. Mern, A. Yildiz, Z. Sunberg, T. Mukerji, and M. J. Kochenderfer, "Bayesian optimized Monte Carlo planning," in *AAAI Conference on Artificial Intelligence*. AAAI Press, 2021.

[15] M. Kearns, Y. Mansour, and A. Y. Ng, "A sparse sampling algorithm for near-optimal planning in large Markov decision processes," *Machine Learning*, vol. 49, no. 2, pp. 193–208, 2002.

[16] N. P. Garg, D. Hsu, and W. S. Lee, "DESPOT-$\alpha$: Online POMDP planning with large state and observation spaces," in *RSS*, 2019.

[17] A. Couëtoux, J.-B. Hoock, N. Sokolovska, O. Teytaud, and N. Bonnard, "Continuous upper confidence trees," in *Learning and Intelligent Optimization*. Rome, Italy: Springer, 2011.

[18] M. J. Kochenderfer, *Decision Making Under Uncertainty: Theory and Application*. Massachusetts: MIT Press, 2015.

[19] D. Bertsekas, *Dynamic Programming and Optimal Control*. Massachusetts: Athena Scientific, 2005.

[20] L. P. Kaelbling, M. L. Littman, and A. R. Cassandra, "Planning and acting in partially observable stochastic domains," *Artificial Intelligence*, vol. 101, no. 1, pp. 99 – 134, 1998.

[21] C. Mansley, A. Weinstein, and M. L. Littman, "Sample-based planning for continuous action markov decision processes," in *International Conference on Automated Planning and Scheduling*. AAAI Press, 2011, p. 335–338.

[22] A. Weinstein and M. L. Littman, "Bandit-based planning and learning in continuous-action markov decision processes," in *International Conference on Automated Planning and Scheduling*. AAAI Press, 2012, p. 306–314.

[23] W. Mao, K. Zhang, Q. Xie, and T. Başar, "Poly-hoot: Monte-carlo planning in continuous space mdps with non-asymptotic analysis," in *Advances in Neural Information Processing Systems*, vol. 33. Curran Associates, Inc., 2020.

[24] P. Morere, R. Marchant, and F. Ramos, "Bayesian optimisation for solving continuous state-action-observation pomdps," in *Advances in Neural Information Processing Systems*, vol. 30. Curran Associates, Inc., 2016.

[25] K. M. Seiler, H. Kurniawati, and S. P. N. Singh, "An online and approximate solver for pomdps with continuous action space," in *IEEE International Conference on Robotics and Automation*. IEEE, 2015, pp. 2290–2297.

[26] J. Mern, A. Yildiz, L. Bush, T. Mukerji, and M. J. Kochenderfer, "Improved POMDP tree search planning with prioritized action branching," in *AAAI Conference on Artificial Intelligence*. AAAI Press, 2021.

[27] J. Lee, W. Jeon, G.-H. Kim, and K.-E. Kim, "Monte-carlo tree search in continuous action spaces with value gradients," in *AAAI Conference on Artificial Intelligence*, vol. 34. AAAI Press, 2020, pp. 4561–4568.

[28] R. Bjarnason, A. Fern, and P. Tadepalli, "Lower bounding Klondike solitaire with Monte-Carlo planning," in *International Conference on Automated Planning and Scheduling*. Thessaloniki, Greece: AAAI Press, 2009.

[29] S. Mannor, R. Rubinstein, and Y. Gat, "The cross entropy method for fast policy search," in *International Conference on Machine Learning*. AAAI Press, 2003, p. 512–519.

[30] M. Egorov, Z. N. Sunberg, E. Balaban, T. A. Wheeler, J. K. Gupta, and M. J. Kochenderfer, "POMDPs.jl: A framework for sequential decision making under uncertainty," *Journal of Machine Learning Research*, vol. 18, no. 26, pp. 1–5, 2017.

[31] G. Chow, *Analysis and Control of Dynamic Economic Systems*. New York: John Wiley & Sons, 1975.

**Contents**

*A. Mathematical Proofs*

*1) Analogous Notations and Concepts to Previous Works:* In our work, we build upon a lot of previous works on sparse sampling, POWSS, and VOO, while adapting notations and concepts to be consistent within our own work. Table I explicitly connects the notations we use in our analyses that are slightly different from the notations used by [15], [12], [13].

TABLE I

SUMMARY OF CORRESPONDING NOTATIONS AND QUANTITIES.

| Proof | Notation | Original Notation | Previous Works | Short Description |
|---|---|---|---|---|
| General | | | | |
| | $C_a$ | $n$ | VOO [12] | Number of action samples |
| | $C_s$ | $C$ | Sparse Sampling [15] | Number of state samples |
| | | | POWSS [13] | Number of state/observation samples |
| | $D$ | $H$ | VOO [12] | Horizon of the problem |
| | | | Sparse Sampling [15] | |
| | $d$ | $h$ | VOO [12] | Given depth of the tree |
| | $d$ | $h, n$ | Sparse Sampling [15] | Given depth of the tree/estimator |
| VOWSS | | | | |
| | $\mathscr{R}_{C_a}$ | $\mathscr{R}_n$ | VOO [12, Theorem 1] | VOO agent regret bound |
| | $\varepsilon$ | $\frac{\lambda}{1-\gamma}$ | POWSS [13, Theorem 2] | POWSS concentration bound |
| | $p$ | $\delta$ | POWSS [13, Theorem 2] | POWSS bound tail probability |
| VOSS | | | | |
| | $\mathscr{R}_{C_a}$ | $\mathscr{R}_n$ | VOO [12, Theorem 1] | VOO agent regret bound |
| | $\varepsilon$ | $\lambda$ | Sparse Sampling [15, Lemma 3] | Chernoff concentration bound |
| | $p$ | $e^{-\lambda^2 C/V_{\max}^2}$ | Sparse Sampling [15, Lemma 3] | Chernoff bound tail probability |

Similarly, Table II connects the concepts we use to the analogous concepts used by [15], [12], [13].

TABLE II

SUMMARY OF ANALOGOUS CONCEPTS.

| Proof | Notation | Original Notation | Previous Works | Short Description |
|---|---|---|---|---|
| VOWSS | | | | |
| | $\eta_{C_a}(d)$ | $\eta(h)$ | VOO [12, Theorem 2] | Value function regret bound |
| | $\alpha_{C_s}(d)$ | $\alpha_d$ | POWSS [13, Lemma 2] | POWSS recursive bound |
| VOSS | | | | |
| | $\eta_{C_a}(d)$ | $\eta(h)$ | VOO [12, Theorem 2] | Value function regret bound |
| | $\alpha_{C_s}(d)$ | $\alpha_n$ | Sparse Sampling [15, Lemma 4] | Sparse sampling recursive bound |

*2) VOWSS Conditions:* Before we prove the VOWSS inequality theorem, we present the necessary conditions for this analysis. The conditions required for the proof is the union of regularity conditions required for POWSS and VOO. The POWSS conditions are the following, with appropriate adaptations to our formalism:

(i) (Continuous spaces) $S, A, O$ are continuous spaces.
(ii) (Bounded Rényi divergence) For any observation sequence $\{o_n\}_j$, the densities $\mathscr{Z}, \mathscr{T}, b_0$ are chosen such that the Rényi divergence of the target distribution $\mathscr{P}^d$ and sampling distribution $\mathscr{Q}^d$ is bounded above by $d_\infty^{\max} < +\infty$ a.s. for all $d = 0, \cdots, D-1$:

$$d_\infty(\mathscr{P}^d || \mathscr{Q}^d) = \text{ess sup}_{x \sim \mathscr{Q}^d} w_{\mathscr{P}^d/\mathscr{Q}^d}(x) \leq d_\infty^{\max}$$

(iii) (Bounded reward function) The reward function $R$ is Borel and bounded by a finite constant $||R||_\infty \leq R_{\max} < +\infty$ a.s., and $V_{\max} \equiv \frac{R_{\max}}{1-\gamma} < +\infty$.
(iv) (Generative model) We can sample from the generating function $G$ and evaluate the observation probability density $\mathscr{Z}$.
(v) (Finite horizon) The POMDP terminates after $D < \infty$ steps.

In our context, the target distribution $\mathscr{P}^d$ corresponds to the conditional observation density of the state trajectory, $\mathscr{Q}^d$ the marginal state trajectory density, and $w_{\mathscr{P}^d/\mathscr{Q}^d}$ the importance weights:

$$\mathscr{P}^d = \mathscr{P}^d_{\{o_n\}_j}(\{s_n\}_i) = \frac{(\mathscr{Z}^{i,j}_{1:d})(\mathscr{T}^i_{1:d})b_0^i}{\int_{S^{d+1}}(\mathscr{Z}^j_{1:d})(\mathscr{T}_{1:d})b_0 ds_{0:d}}, \tag{1}$$

$$\mathscr{Q}^d = \mathscr{Q}^d(\{s_n\}_i) = (\mathscr{T}^i_{1:d})b_0^i, \tag{2}$$

$$w_{\mathscr{P}^d/\mathscr{Q}^d}(\{s_n\}_i) = \frac{(\mathscr{Z}^{i,j}_{1:d})}{\int_{S^{d+1}}(\mathscr{Z}^j_{1:d})(\mathscr{T}_{1:d})b_0 ds_{0:d}}. \tag{3}$$

Essentially, condition (ii) means that the conditional observation density cannot be much larger than the marginal density for any given state trajectory. With the above conditions, the state sampling width $C_s$ is chosen such that the following holds:

$$p \geq 3C_a(3C_a \cdot C_s)^D \exp(-C_s \cdot t_{\max}^2), \tag{4}$$

$$p = \frac{\varepsilon}{(1-\gamma)V_{\max}D}, \quad t_{\max}(\varepsilon, C_s) = \frac{\varepsilon(1-\gamma)}{3V_{\max}d_\infty^{\max}} - \frac{1}{\sqrt{C_s}}. \tag{5}$$

This ensures the POWSS type bound $|Q_d^\star(b,a) - \hat{Q}_d^\star(\bar{b},a)| \leq \varepsilon$ holds with probability at least $1-p$ for each $b, a$ at depth $d$.

The VOO conditions are the following, with appropriate adaptations to our formalism. Note that with our notation, VOO aims to optimize a function $Q(a)$ over some space $A$:

(i) (Translation-invariant semi-metric) $D : A \times A \to \mathbb{R}^+$ is such that $\forall x, y, z \in A$, $D(x,y) = D(y,x)$, $D(x,y) = 0$ iff $x = y$, and $D(x+z, y+z) = D(x,y)$.
(ii) (Local smoothness of $Q$) There exists at least one global optimum $a^* \in A$ of $Q$ such that $\forall a \in A, Q(a^*) - Q(a) \leq L \cdot D(a, a^*)$ for some $L > 0$.
(iii) (Shrinkage ratio of the Voronoi cells) Consider any point $a'$ inside the Voronoi cell $C$ generated by the point $a_0$, and denote $d_0 = D(a', a_0)$. If we randomly sample a point $a_1$ from $C$, we have $\mathbb{E}[\min(d_0, D(a', a_1))] \leq \lambda d_0$ for $\lambda \in (0,1)$.
(iv) (Well-shaped Voronoi cells) There exists $\eta > 0$ such that for any Voronoi cell generated by $a$ with expected diameter $d_0$ contains a ball of radius $\eta d_0$ centered at $a$.
(v) (Local symmetry near optimum) The set of global optima $A_*$ consists of finite number of disjoint and connected components $\{A_*^{(l)}\}_{l=1}^k, k < \infty$. For each component, there exists an open ball $B_{v_l}(a_*^{(l)})$ for some $a_*^{(l)} \in A_*^{(l)}$ such that $D(a, a_*^{(l)}) \leq D(a', a_*^{(l)})$ implies $Q(a) \geq Q(a')$ for any $a, a' \in B_{v_l}(a_*^{(l)})$.

Here, we define $\bar{\mu}_B(r) = \mu(B_r(\cdot))/\mu(A)$, with $\mu$ the Borel measure on $A$, $\delta_{\max}$ the largest distance between two points in $A$, and $v_{\min} = \min_{l \in [k]} v_l$. Then, the action sampling width $C_a$ is chosen such that the following holds:

$$\omega \geq \frac{1 - \lambda^{1/k}}{\bar{\mu}_B(v_{\min}) + 1 - \bar{\mu}_B(\eta \cdot \lambda v_{\min})}, \quad C_a \geq \max_{d=0,\cdots,D-1} \left[ \log\left( \frac{\eta_{C_a}(d) - \gamma \eta_{C_a}(d+1)}{2L\delta_{\max}C_{\max}} \right) \cdot \min(G_{\lambda,\omega}, K_{v,\omega,\lambda}) \right]. \tag{6}$$

Here, $C_{\max}, G_{\lambda,\omega}, K_{v,\omega,\lambda}$ are problem specific quantities/functions that are explicitly defined in [12]. Satisfying these constraints for a decreasing sequence $\{\eta_{C_a}(d)\}$ will allow us to use the VOO type bound $V_d^\star(b) - \hat{V}_d^\star(b) \leq \eta_{C_a}(d)$ that holds in expectation.

*3) Proof of VOWSS Inequality:*

*Theorem 1 (VOWSS Inequality):* Suppose we choose the action sampling width $C_a$ and state sampling width $C_s$ such that under the union of regularity conditions specified by [13] and [12], the intermediate POWSS bounds and VOO bounds in Lemma 1 are satisfied at every depth of the tree. Then, the following bounds for the VOWSS estimator $\hat{V}_{\text{VOWSS},d}(\bar{b})$ hold for all $d \in [0, D-1]$ in expectation:

$$\left| V_d^\star(b) - \hat{V}_{\text{VOWSS},d}(\bar{b}) \right| \leq \eta_{C_a} + \alpha_{C_s}. \tag{7}$$

In order to prove Theorem 1, we first prove an intermediate lemma which will allow us to obtain the bound through triangle inequality. We introduce and prove the following lemma first. All of the following calculations are done in expectation. We also denote $\bar{b}$ for a particle representation of belief $b$ that POWSS and VOWSS take as an argument.

*Lemma 1 (VOWSS Intermediate Inequality):* Suppose with our notation, the POWSS estimators at all depths $d$ are within $\varepsilon$ of their mean values with probability $1 - p$, and the VOO agents have regret bounds of $\mathscr{R}_{C_a}$. The following inequalities hold for all $d \in [0, D-1]$ in expectation:

$$\left| V_d^\star(b) - \hat{V}_d^{C_a}(b) \right| \leq \eta_{C_a}(d), \tag{8}$$

$$\left| \hat{V}_d^{C_a}(b) - \hat{V}_{\text{VOWSS},d}(\bar{b}) \right| \leq \alpha_{C_s}(d). \tag{9}$$

$\eta_{C_a}(d)$ and $\alpha_{C_s}(d)$ are sequences that satisfy the following properties:

$$\eta_{C_a}(d) \geq \gamma \cdot \eta_{C_a}(d+1) + \mathscr{R}_{C_a}, \ \eta_{C_a}(D) = 0, \tag{10}$$

$$\eta_{C_a} \equiv \max_{d=0,\cdots,D-1} \eta_{C_a}(d) < +\infty, \tag{11}$$

$$\alpha_{C_s}(d) \equiv (1+\gamma)(\varepsilon + 2p \cdot V_{\max}) + \gamma(\alpha_{C_s}(d+1) + 2p \cdot V_{\max}), \tag{12}$$

$$\alpha_{C_s}(D-1) = \varepsilon + 2p \cdot V_{\max}, \tag{13}$$

$$\alpha_{C_s} \equiv \max_{d=0,\cdots,D-1} \alpha_{C_s}(d) < +\infty. \tag{14}$$

*Proof:* This proof proceeds through induction by assuming that the bounds hold for all depths from $d+1$ to $D-1$, and then proving they also hold for depth $d$ (induction hypothesis is omitted as the bounds trivially hold as per the definitions of VOOT and POWSS). We divide the main inequality into VOO bound and POWSS bound by introducing an intermediate term $\hat{V}_d^{C_a}(b)$:

$$\left| V_d^\star(b) - \hat{V}_{\text{VOWSS},d}(b) \right| \leq \underbrace{\left| V_d^\star(b) - \hat{V}_d^{C_a}(b) \right|}_{\text{VOO bound}} + \underbrace{\left| \hat{V}_d^{C_a}(b) - \hat{V}_{\text{VOWSS},d}(\bar{b}) \right|}_{\text{POWSS bound}} \tag{15}$$

Essentially, we have two main layers of inequality, caused by the stochastic nature of VOO action selection, and the uncertainties in state transition and observation. We will first analyze the VOO bound, then the POWSS bound.

*a) Step 1: VOO Bound:* The VOO bound can once again be further decomposed into the following terms as [12] do in their work:

$$\left| V_d^\star(b) - \hat{V}_d^{C_a}(b) \right| \leq \underbrace{\left| V_d^\star(b) - \hat{V}_d(b) \right|}_{\text{VOO Recursive bound}} + \underbrace{\left| \hat{V}_d(b) - \hat{V}_d^{C_a}(b) \right|}_{\text{VOO Regret bound}} \tag{16}$$

Here, the intermediate random variables are defined in the following manner:

$$V_d^\star(b) \equiv \max_{a \in A} Q_d^\star(b,a) = \max_{a \in A} \left\{ R(b,a) + \gamma \mathbb{E}[V_{d+1}^\star(bao)|b] \right\} \tag{17}$$

$$\hat{V}_d(b) \equiv \max_{a \in A} \hat{Q}_d(b,a) = \max_{a \in A} \left\{ R(b,a) + \gamma \mathbb{E}[\hat{V}_{d+1}^{C_a}(bao)|b] \right\} \tag{18}$$

$$\hat{V}_d^{C_a}(b) \equiv \max_{a \in VOO(A,C_a)} \hat{Q}_d(b,a) = \max_{a \in VOO(A,C_a)} \left\{ R(b,a) + \gamma \mathbb{E}[\hat{V}_{d+1}^{C_a}(bao)|b] \right\} \tag{19}$$

As a notation, $a \in VOO(A, C_a)$ indicates that the actions $a$ are chosen sequentially through the VOO algorithm over the action space $A$ for $C_a$ iterations of VOO. Here, we compare the two quantities with a reference action $a^\star = \arg\max Q_d^\star(b,a)$, which results in a looser bound but allows us to directly compare the quantities inside the max operations. We closely follow

the calculations from the proof of Lemma 5 in [12]:

$$\underbrace{\left|V_d^\star(b) - \hat{V}_d(b)\right|}_{\text{VOO Recursive bound}} = \left|\max_{a\in A}\left\{R(b,a) + \gamma\mathbb{E}[V_{d+1}^\star(bao)|b]\right\} - \max_{a\in A}\left\{R(b,a) + \gamma\mathbb{E}[\hat{V}_{d+1}^{C_a}(bao)|b]\right\}\right| \tag{20}$$

$$\leq \left|R(b,a^\star) + \gamma\mathbb{E}[V_{d+1}^\star(ba^\star o)|b] - R(b,a^\star) + \gamma\mathbb{E}[\hat{V}_{d+1}^{C_a}(ba^\star o)|b]\right| \qquad (a^\star = \arg\max Q_d^\star(b,a))$$

$$\leq \gamma\mathbb{E}\left[\left|V_{d+1}^\star(bao) - \hat{V}_{d+1}^{C_a}(bao)\right|\,\Big|\,b\right] \tag{21}$$

$$\leq \gamma\cdot\eta_{C_a}(d+1) \tag{22}$$

$$\underbrace{\left|\hat{V}_d(b) - \hat{V}_d^{C_a}(b)\right|}_{\text{VOO Regret bound}} \leq \mathscr{R}_{C_a}. \tag{23}$$

In the VOO Regret bound, the difference cannot be less than zero since the global maximum is at least as big as the VOO maximum, so the absolute value disappears and we get the regret of VOO. Thus, with our choice of $C_a$ that is designed to satisfy the recurrence relation:

$$\left|V_d^\star(b) - \hat{V}_d^{C_a}(b)\right| \leq \underbrace{\left|V_d^\star(b) - \hat{V}_d(b)\right|}_{\text{VOO Recursive bound}} + \underbrace{\left|\hat{V}_d(b) - \hat{V}_d^{C_a}(b)\right|}_{\text{VOO Regret bound}} \tag{24}$$

$$\leq \gamma\cdot\eta_{C_a}(d+1) + \mathscr{R}_{C_a} \leq \eta_{C_a}(d). \tag{25}$$

*b) Step 2: POWSS Bound:* The POWSS bound can also be further decomposed into the following terms:

$$\left|\hat{V}_d^{C_a}(b) - \hat{V}_{\text{VOWSS},d}(b)\right| \leq \underbrace{\left|\hat{V}_d^{C_a}(b) - \tilde{V}_{\text{VOWSS},d}(b)\right|}_{\text{POWSS Concentration bound}} + \underbrace{\left|\tilde{V}_{\text{VOWSS},d}(b) - \hat{V}_{\text{VOWSS},d}(b)\right|}_{\text{POWSS Recursive bound}}. \tag{26}$$

Here, the extra intermediate random variable and the VOWSS estimator are defined in the following manner:

$$\tilde{V}_{\text{VOWSS},d}(b) \equiv \max_{a\in VOO(A,C_a)}\tilde{Q}_{\text{VOWSS},d}(b,a) = \max_{a\in VOO(A,C_a)}\left\{\frac{\sum_{i=1}^{C_s}w_{d,i}(r_{d,i} + \gamma\hat{V}_{d+1}^{C_a}(bao_i))}{\sum_{i=1}^{C_s}w_{d,i}}\right\}, \tag{27}$$

$$\hat{V}_{\text{VOWSS},d}(\bar{b}) \equiv \max_{a\in VOO(A,C_a)}\hat{Q}_{\text{VOWSS},d}(\bar{b},a) = \max_{a\in VOO(A,C_a)}\left\{\frac{\sum_{i=1}^{C_s}w_{d,i}(r_{d,i} + \gamma\hat{V}_{\text{VOWSS},d+1}(\overline{bao_i}))}{\sum_{i=1}^{C_s}w_{d,i}}\right\}. \tag{28}$$

We compare the two quantities by picking a reference action to directly compare the quantities inside the max operations.
For the concentration bound term:

$$\underbrace{\left|\hat{V}_d^{C_a}(b) - \tilde{V}_{\text{VOWSS},d}(b)\right|}_{\text{POWSS Concentration bound}} \leq \left|\max_{a\in VOO(A,C_a)}\left\{R(b,a) + \gamma\mathbb{E}[\hat{V}_{d+1}^{C_a}(bao)|b]\right\}\right. \tag{29}$$

$$\left. - \max_{a\in VOO(A,C_a)}\left\{\frac{\sum_{i=1}^{C_s}w_{d,i}(r_{d,i} + \gamma\hat{V}_{d+1}^{C_a}(bao_i))}{\sum_{i=1}^{C_s}w_{d,i}}\right\}\right| \tag{30}$$

$$\leq \left|R(b,a^*) + \gamma\mathbb{E}[\hat{V}_{d+1}^{C_a}(ba^*o)|b] - \frac{\sum_{i=1}^{C_s}w_{d,i}(r_{d,i} + \gamma\hat{V}_{d+1}^{C_a}(ba^*o_i))}{\sum_{i=1}^{C_s}w_{d,i}}\right|. \qquad (a^* = \arg\max_{VOO}\hat{Q}_d(b,a))$$

As a small note, picking the maximizing action with respect to $\hat{Q}_d$ does not guarantee that we have bounded the term in this case. We also need to consider the case of picking the maximizing action with respect to $\tilde{Q}_{\text{VOWSS},d}$, since $\tilde{Q}_{\text{VOWSS},d}$ could achieve larger values than $\hat{Q}_d$ due to sparse sampling of states. However, our overall result does not change, since choosing the other maximizing action will only result in a change in sign within the absolute value, which we can still bound with the recursive bound. Thus, in the following calculations whenever we pick the reference action for difference in quantities that do not have strict magnitude hierarchy, we do not consider this possibility as the terms can still be bounded with their respective recursive bounds by choosing an appropriate reference action.
Decomposing the quantity into the reward difference and the next-step value difference:

$$\underbrace{\left|\hat{V}_d^{C_a}(b) - \tilde{V}_{\text{VOWSS},d}(b)\right|}_{\text{POWSS Concentration bound}} \leq \underbrace{\left|R(b,a^*) - \frac{\sum_{i=1}^{C_s}w_{d,i}r_{d,i}}{\sum_{i=1}^{C_s}w_{d,i}}\right|}_{\text{Reward difference}} + \gamma\underbrace{\left|\mathbb{E}[\hat{V}_{d+1}^{C_a}(ba^*o)|b] - \frac{\sum_{i=1}^{C_s}w_{d,i}\hat{V}_{d+1}^{C_a}(ba^*o_i)}{\sum_{i=1}^{C_s}w_{d,i}}\right|}_{\text{Value difference}}. \tag{31}$$

For the reward difference, we can crudely upper bound it by the POWSS $Q$-value estimate concentration bound in Theorem 2 of [13], since this is effectively the same structure as the leaf node estimate. In the POWSS concentration bound in Theorem 2, the difference has an upper bound of $\lambda/(1-\gamma)$, which we will define as $\varepsilon$, and we denote the corresponding probability $\delta$ as $p$. This bounds the quantity by $(\varepsilon + 2p \cdot V_{\max})$ using the expectation version of the POWSS concentration bound.

$$\underbrace{\left| R(b,a^*) - \frac{\sum_{i=1}^{C_s} w_{d,i} r_{d,i}}{\sum_{i=1}^{C_s} w_{d,i}} \right|}_{\text{Reward difference}} \leq \varepsilon + 2p \cdot V_{\max}. \tag{32}$$

We could instead use Lemma 1 of [13], the leaf node estimate concentration bound, but we use the more general Theorem 2. This allows us to consistently use the same theorem throughout this analysis and effectively combine terms. Similarly, while the worst case bound is $2R_{\max}$ since we are taking the difference of two reward functions, we crudely upper bound that with $2V_{\max}$ for algebraic convenience of combining it with the value difference term.

$$\underbrace{\left| \mathbb{E}[\hat{V}_{d+1}^{C_a}(ba^*o)|b] - \frac{\sum_{i=1}^{C_s} w_{d,i} \hat{V}_{d+1}^{C_a}(ba^*o_i)}{\sum_{i=1}^{C_s} w_{d,i}} \right|}_{\text{Value difference}} \leq \varepsilon + 2p \cdot V_{\max}. \tag{33}$$

On the other hand, for the value difference, the POWSS concentration bound also turns out to be an upper bound, but with more sophisticated calculations. During this part of the proof, we will refer heavily back to the continued proof of Lemma 2 in the Appendix C of [13] and give a general overview of how the steps apply here.

*c) Step 2-i: Value Difference:* While Lemma 2 in [13] is calculated with respect to the theoretically optimal value function $V_{d+1}^*$, the calculation steps themselves and the theorems and lemmas used there can apply exactly the same way for $\hat{V}_{d+1}^{C_a}$. We will briefly illustrate how the steps are parallel to the proof of Lemma 2 in the Appendix C of [13].

The value difference corresponds to the difference between the expected value of $\hat{V}_{d+1}^{C_a}$ and the self-normalized importance sampling estimator of the expected value. This value difference specifically corresponds to the first two error terms in the continued proof of Lemma 2 in the Appendix C of [13]. The decomposition looks like the following:

$$\left| \mathbb{E}[\hat{V}_{d+1}^{C_a}(ba^*o)|b] - \frac{\sum_{i=1}^{C_s} w_{d,i} \hat{V}_{d+1}^{C_a}(ba^*o_i)}{\sum_{i=1}^{C_s} w_{d,i}} \right| \tag{34}$$

$$\leq \underbrace{\left| \mathbb{E}[\hat{V}_{d+1}^{C_a}(ba^*o)|b] - \frac{\sum_{i=1}^{C_s} w_{d,i} \hat{\mathbf{V}}_{d+1}^{C_a}(s_{d,i},b,a^*)}{\sum_{i=1}^{C_s} w_{d,i}} \right|}_{\text{Importance sampling error}} + \underbrace{\left| \frac{\sum_{i=1}^{C_s} w_{d,i}(\hat{\mathbf{V}}_{d+1}^{C_a}(s_{d,i},b,a^*) - \hat{V}_{d+1}^{C_a}(ba^*o_i))}{\sum_{i=1}^{C_s} w_{d,i}} \right|}_{\text{MC next-step integral approximation error}}. \tag{35}$$

The next-step marginal integral $\hat{\mathbf{V}}_{d+1}^{C_a}(s_{d,i},b,a^*)$ is defined as

$$\hat{\mathbf{V}}_{d+1}^{C_a}(s_{d,i},b,a^*) \equiv \int_S \int_O \hat{V}_{d+1}^{C_a}(ba^*o) \mathscr{Z}(o|a^*,s_{d+1}) \mathscr{T}(s_{d+1}|s_{d,i},a^*) ds_{d+1} do. \tag{36}$$

With this analogous definition, the self-normalized estimator identities in [13] can be exactly applied to this setting once again, now instead for the function $\hat{V}_d^{C_a}$ because the algebraic steps taken in the proof should hold for our $V$ function estimator as well. Intuitively, the next-step marginal integral is defined to be the marginal random variable for $\hat{V}_d^{C_a}$, instead of the optimal $V$ function as it is done in the works of [13].

First, following [13], we define the following notation for products of transition and observation densities:

$$\mathscr{T}_{1:d}^i \equiv \prod_{n=1}^{d} \mathscr{T}(s_{n,i}|s_{n-1,i},a_n), \tag{37}$$

$$\mathscr{Z}_{1:d}^{i,j} \equiv \prod_{n=1}^{d} \mathscr{Z}(o_{n,j}|a_n,s_{n,i}). \tag{38}$$

Specifically, $i$ denotes the index of the state sample, and $j$ denotes the index of the observation sample. Absence of any of the indices $i$ or $j$ means that the state trajectory $\{s_n\}$ or the observation history $\{o_n\}$ appear as regular variables, mostly for the purposes of integration.

*d) Step 2-ii: Importance Sampling Error:* For the importance sampling error, note that for a belief $b$ at depth $d$,

$$\mathbb{E}[\hat{V}_{d+1}^{C_a}(ba^*o)|b] = \int_S \int_S \int_O \hat{V}_{d+1}^{C_a}(ba^*o)(\mathscr{Z}_{d+1})(\mathscr{T}_{d,d+1})b \cdot ds_{d:d+1}do \tag{39}$$

$$= \int_S \hat{\mathbf{V}}_{d+1}^{C_a}(s_{d,i},b,a^*)b \cdot ds_d \tag{40}$$

$$= \frac{\int_{S^d} \hat{\mathbf{V}}_{d+1}^{C_a}(s_{d,i},b,a^*)(\mathscr{Z}_{1:d})(\mathscr{T}_{1:d})b_0 ds_{0:d}}{\int_{S^d}(\mathscr{Z}_{1:d})(\mathscr{T}_{1:d})b_0 ds_{0:d}}. \tag{41}$$

Consequently, the weighted average of the next-step marginal integral $\hat{\mathbf{V}}_{d+1}^{C_a}(s_{d,i},b,a^*)$ is a self-normalized importance sampling estimator of $\hat{V}_{d+1}^{C_a}(ba^*o)$ given $b$, and we can apply the augmented self-normalized estimator concentration bound in the same way as [13] to get the concentration bound of $\lambda/3$.

*e) Step 2-iii: Monte Carlo Next-Step Integral Approximation Error:* For the MC next-step integral approximation error, generating estimates of $\hat{V}_{d+1}^{C_a}(ba^*o_i)$ for a given $(s_{d,i},b,a^*)$ also results in an unbiased Monte Carlo estimator of the next-step marginal integral, and the following difference is mean zero conditioned on $(s_{d,i},b,a^*)$:

$$\Delta_{d+1}(s_{d,i},b,a^*) \equiv \hat{\mathbf{V}}_{d+1}^{C_a}(s_{d,i},b,a^*) - \hat{V}_{d+1}^{C_a}(ba^*o_i). \tag{42}$$

Thus, the same calculation steps hold. Note that in the works of [13], the MC next-step integral approximation error is further crudely bound by $\frac{2}{3\gamma}\lambda$, when the actual bounds also hold for $\frac{2}{3}\lambda$, for the convenience of being able to combine the $\gamma$ multiplied terms in the main proof of Lemma 2. Thus, we can bound the MC next-step integral approximation error with the stricter bound $\frac{2}{3}\lambda$.

In our work, we used the variable $\varepsilon$ to denote the POWSS concentration bound, which corresponds to $\varepsilon = \lambda/(1-\gamma) \geq \lambda$. Since the sum of importance sampling error and MC next-step integral approximation error are bounded by $(1/3 + 2/3)\lambda = \lambda$, this can also further be crudely bounded by the POWSS concentration inequality with the upper bound $\varepsilon$.

We have now obtained bounds for the value difference term using the expectation version of the POWSS concentration inequality where the extreme probability event is once again bounded with the term $2p \cdot V_{\max}$:

$$\left| \mathbb{E}[\hat{V}_{d+1}^{C_a}(ba^*o)|b] - \frac{\sum_{i=1}^{C_s} w_{d,i}\hat{V}_{d+1}^{C_a}(ba^*o_i)}{\sum_{i=1}^{C_s} w_{d,i}} \right| \leq \varepsilon + 2p \cdot V_{\max}. \tag{43}$$

Finally, we can bound the POWSS Concentration bound term:

$$\underbrace{\left| \hat{V}_d^{C_a}(b) - \tilde{V}_{\text{VOWSS},d}(b) \right|}_{\text{POWSS Concentration bound}} \leq (\varepsilon + 2p \cdot V_{\max}) + \gamma(\varepsilon + 2p \cdot V_{\max}) = (1+\gamma)(\varepsilon + 2p \cdot V_{\max}). \tag{44}$$

For the POWSS Recursive bound term, we simply apply the inductive hypothesis for step $d+1$:

$$\underbrace{\left| \tilde{V}_{\text{VOWSS},d}(b) - \hat{V}_{\text{VOWSS},d}(\bar{b}) \right|}_{\text{POWSS Recursive bound}} \leq \left| \max_{a \in VOO(A,C_a)} \left\{ \frac{\sum_{i=1}^{C_s} w_{d,i}(r_{d,i} + \gamma \hat{V}_{d+1}^{C_a}(bao_i))}{\sum_{i=1}^{C_s} w_{d,i}} \right\} \right. \tag{45}$$

$$\left. - \max_{a \in VOO(A,C_a)} \left\{ \frac{\sum_{i=1}^{C_s} w_{d,i}(r_{d,i} + \gamma \hat{V}_{\text{VOWSS},d+1}(\overline{bao_i}))}{\sum_{i=1}^{C_s} w_{d,i}} \right\} \right| \tag{46}$$

$$\leq \left| \gamma \frac{\sum_{i=1}^{C_s} w_{d,i}(\hat{V}_{d+1}^{C_a}(b\tilde{a}o_i) - \hat{V}_{\text{VOWSS},d+1}(\overline{b\tilde{a}o_i}))}{\sum_{i=1}^{C_s} w_{d,i}} \right| \qquad (\tilde{a} = \arg\max_{VOO} \tilde{Q}_{\text{VOWSS},d}(b,a))$$

$$\leq \gamma \frac{\sum_{i=1}^{C_s} w_{d,i} \left| \hat{V}_{d+1}^{C_a}(b\tilde{a}o_i) - \hat{V}_{\text{VOWSS},d+1}(\overline{b\tilde{a}o_i}) \right|}{\sum_{i=1}^{C_s} w_{d,i}} \tag{47}$$

$$\leq \gamma \cdot \alpha_{C_s}(d+1). \tag{48}$$

Putting the POWSS components together, we prove the POWSS bound by induction:

$$\left| \hat{V}_d^{C_a}(b) - \hat{V}_{\text{VOWSS},d}(\bar{b}) \right| \leq \underbrace{\left| \hat{V}_d^{C_a}(b) - \tilde{V}_{\text{VOWSS},d}(b) \right|}_{\text{POWSS Concentration bound}} + \underbrace{\left| \tilde{V}_{\text{VOWSS},d}(b) - \hat{V}_{\text{VOWSS},d}(\bar{b}) \right|}_{\text{POWSS Recursive bound}} \tag{49}$$

$$\leq (1+\gamma)(\varepsilon + 2p \cdot V_{\max}) + \gamma \cdot \alpha_{C_s}(d+1) = \alpha_{C_s}(d). \blacksquare \tag{50}$$

*Proof for Theorem 1.* Finally, we prove the main theorem by combining the two terms:

$$\left|V_d^\star(b) - \hat{V}_{\text{VOWSS},d}(\bar{b})\right| \leq \underbrace{\left|V_d^\star(b) - \hat{V}_d^{C_a}(b)\right|}_{\text{VOO bound}} + \underbrace{\left|\hat{V}_d^{C_a}(b) - \hat{V}_{\text{VOWSS},d}(b)\right|}_{\text{POWSS bound}} \leq \eta_{C_a}(d) + \alpha_{C_s}(d) \leq \eta_{C_a} + \alpha_{C_s}. \quad \blacksquare \tag{51}$$

Here, the $\eta_{C_a}(d)$ corresponds to the VOO bound in [12]. Thus, we can always find a corresponding $C_a$ for an arbitrary decreasing sequence of $\eta_{C_a}(d)$ which satisfies the properties in our lemma, and this sequence can be made closer to 0 for each depth by choosing bigger $C_a$ values.

On the other hand, $\alpha_{C_s}(d)$ is comprised of POWSS bound components in [13], which decreases in both $p$ and $\varepsilon$ with more samples $C_s$. Thus, it can also be made to decrease to 0 as we increase $C_s$.

## B. Experiment Hyperparameters

Hyperparameters were taken from the references if they were given, and otherwise the set of hyperparmaters for a system was obtained by first training POMCPOW, and then training VOMCPOW and BOMCP centered around POMCPOW hyperparameters as initial estimates. This usually augmented the performances of VOMCPOW and BOMCP compared to directly borrowing the POMCPOW hyperparameters. For BOMCP, buffer and $k$ were fixed at 100 and 5, respectively, as it was done in the original paper by [14]. Table III shows the final hyperparameters used in the experiments. Specifically for the lunar lander problem, we have chosen the dynamics time step $dt = 0.4$, which gave us the most consistent performance of BOMCP with 100 queries over 1000 iterations that is in line with the results of [14].

We give a brief intuitive explanation for each of the parameters utilized. $c$ is the critical factor or the Upper Confidence Bound exploration parameter, which governs how much we should explore among the action samples we have collected. Since VPW needs to balance local and global search, on top of estimating relatively faithful $Q$-values, it usually worked better to set the $c$ on the equal magnitude/lower than the $c$ value for POMCPOW, usually around 10-100. $k_a$ is the constant factor of the action widening parameter, which sets the overall width of the action widening. $\alpha_a$ is the exponential factor of the action widening parameter, which determines how adaptive we want to progressively widen the action width. $k_o, \alpha_o$ are the state/observation widening parameters, which function similarly to the action widening parameters. $\omega$ is the VOO exploration probability, which governs how much we should explore in the action space. Since in global search, we need to see enough samples to explore the action space sufficiently, it usually worked better to set $\omega$ to be relatively high around 0.7-0.9. Lastly, $\Sigma$ is the Gaussian covariance matrix for VOO rejection sampling.

The only hyperparameter that was manually picked is the Gaussian covariance matrix for VOO rejection sampling. When picking this covariance matrix, we usually found that it was most effective to use a diagonal matrix with entries that are around 10 to 20 times less than the maximum action space bounds for each dimension. In theory, the diagonal covariance matrix can also be fitted via CEM, but we chose to handpick these values after some inspection in order to reduce the number of hyperparameters that needed to be fit.

To limit the number of rejection sampling iterations, we set the maximum number of rejection sampling iterations to 20 and automatically choose the closest action we have sampled when we reach 20 iterations. We also manually set automatic sample acceptance regions that was usually on the order of a tenth of the sampling radius, which we found was not strictly necessary in conjunction with the sampling iteration limit.

The open source code is available at `github.com/michael-lim/VOOTreeSearch.jl`.

TABLE III

SUMMARY OF HYPERPARAMETERS USED IN EXPERIMENTS FOR POMCPOW, VOMCPOW, BOMCP, AND VOWSS.

| | $c$ | $k_a$ | $\alpha_a$ | $k_o$ | $\alpha_o$ | $C_s$ | $C_a$ | $\omega$ | $\Sigma = \text{diag}(\sigma^2)$ | $l$ | $\lambda$ | $\gamma_a$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **LQG (Depth 3)** | | | | | | | | | | | | |
| POMCPOW | 65.0 | 30.0 | $\frac{1}{2.5}$ | 30.0 | $\frac{1}{4}$ | | | | | | | |
| VOMCPOW | 60.0 | 25.0 | $\frac{1}{5.5}$ | 25.0 | $\frac{1}{2.5}$ | | | 0.8 | [0.5, 0.5] | | | |
| BOMCP | 135.0 | 30.0 | $\frac{1}{4}$ | 20.0 | $\frac{1}{4}$ | | | | | log(15) | 0.4 | |
| VOWSS | | | | | | 10 | 200 | 0.8 | [0.5, 0.5] | | | 0.4 |
| **VDP Tag (Depth 10)** | | | | | | | | | | | | |
| POMCPOW | 110.0 | 30.0 | $\frac{1}{30}$ | 5.0 | $\frac{1}{100}$ | | | | | | | |
| VOMCPOW | 85.0 | 30.0 | $\frac{1}{30}$ | 2.5 | $\frac{1}{100}$ | | | 0.7 | [0.1] | | | |
| **Lunar Lander (Depth 250)** | | | | | | | | | | | | |
| POMCPOW | 10.0 | 3.0 | $\frac{1}{4}$ | 2.0 | $\frac{1}{10}$ | | | | | | | |
| VOMCPOW | 30.0 | 4.0 | $\frac{1}{4}$ | 1.5 | $\frac{1}{5}$ | | | 0.9 | [0.2, 0.5, 0.05] | | | |
| BOMCP | 10.0 | 3.0 | $\frac{1}{4}$ | 2.0 | $\frac{1}{10}$ | | | | | log(15) | 0.5 | |

*C. Voronoi Optimistic Sparse Sampling (VOSS) for Stochastic MDPs*

Voronoi Optimistic Sparse Sampling (VOSS) is an application of VPW to the sparse sampling algorithm [15] to tackle the stochastic MDP case. Like VOWSS, it can be defined by an ESTIMATEQ function that estimates the $Q$-value with next-step state samples. Since VOSS does not have observation uncertainty, it only needs to take the arithmetic average instead of the observation likelihood weighted average. The ESTIMATEQ function that takes in a state $s$ instead of a belief particle set $\bar{b}$ is outlined in Algorithm 5. Note that this is analogous definition to the $Q$-function estimation algorithm in [15]. The ESTIMATEV function should function similarly, except working with $s$ instead of $\bar{b}$.

---

**Algorithm 5** Value estimation algorithm for VOSS

---

**Algorithm:** ESTIMATEQ$(s,a,d)$
**Input:** State $s$, action $a$, depth $d$.
**Output:** A scalar $\hat{Q}_d^\star(s,a)$ that is an estimate of $Q_d^\star(s,a)$.
1: For $i = 1, \cdots, C_s$, generate $s_i', r = G(s,a)$.
2: Return the $Q$-value estimate:

$$\hat{Q}_d^\star(s,a) = r + \gamma \frac{1}{C_s} \sum_{i=1}^{C_s} \text{ESTIMATEV}(s_i', d+1).$$

---

The conditions required for the proof is the union of regularity conditions required for sparse sampling and VOO. The sparse sampling conditions are the subset of POWSS conditions. Namely, we only require the conditions (i), (iii), (iv), only for the state and action spaces. Since in this proof we simply use the Chernoff bound for the sparse sampling type bound, we just need to make sure that $p \geq \exp(-\varepsilon^2 C_s/(2V_{\max})^2)$ holds when picking $C_s$.

*Theorem 2 (VOSS Inequality):* Suppose we choose the action sampling width $C_a$ and state sampling width $C_s$ such that under the union of regularity conditions specified by [15] and [12], the intermediate sparse sampling bounds and VOO bounds in Lemma 2 are satisfied at every depth of the tree. Then, the following bounds for the VOSS estimator $\hat{V}_{\text{VOSS},d}(s)$ hold for all $d \in [0, D-1]$ in expectation:

$$\left| V_d^\star(s) - \hat{V}_{\text{VOSS},d}(s) \right| \leq \eta_{C_a} + \alpha_{C_s}. \tag{52}$$

Similar to proving Theorem 1, to prove Theorem 2, we first prove an intermediate lemma which will allow us to obtain the bound through triangle inequality. The proof is easier than that of Theorem 1, since we do not explicitly need to deal with the observation uncertainty. We introduce and prove the following lemma first. All of the following calculations are done in expectation.

*Lemma 2 (VOSS Intermediate Inequality):* Suppose with our notation, the sparse sampling estimators at all depths $d$ are within $\varepsilon$ of their mean values with probability $1 - p$, and the VOO agents have regret bounds of $\mathscr{R}_{C_a}$. The following inequalities hold for all $d \in [0, D-1]$ in expectation:

$$\left| V_d^\star(s) - \hat{V}_d^{C_a}(s) \right| \leq \eta_{C_a}(d), \tag{53}$$

$$\left| \hat{V}_d^{C_a}(s) - \hat{V}_{\text{VOSS},d}(s) \right| \leq \alpha_{C_s}(d). \tag{54}$$

$\eta_{C_a}(d)$ and $\alpha_{C_s}(d)$ are sequences that satisfy the following properties:

$$\eta_{C_a}(d) \geq \gamma \cdot \eta_{C_a}(d+1) + \mathscr{R}_{C_a}, \ \eta_{C_a}(D) = 0, \tag{55}$$

$$\eta_{C_a} \equiv \max_{d=0,\cdots,D-1} \eta_{C_a}(d) < +\infty, \tag{56}$$

$$\alpha_{C_s}(d) \equiv \gamma(\alpha_{C_s}(d+1) + \varepsilon + 2p \cdot V_{\max}), \ \alpha_{C_s}(D-1) = \varepsilon + 2p \cdot V_{\max}, \tag{57}$$

$$\alpha_{C_s} \equiv \max_{d=0,\cdots,D-1} \alpha_{C_s}(d) < +\infty. \tag{58}$$

*Proof:* This proof proceeds through induction by assuming that the bounds hold for all depths from $d+1$ to $D-1$, and then proving they also hold for depth $d$ (induction hypothesis is omitted as the bounds trivially hold as per the definitions of VOOT and sparse sampling). We divide the main inequality into VOO bound and sparse sampling bound (SS bound) by introducing an intermediate term $\hat{V}_d^{C_a}(s)$:

$$\left| V_d^\star(s) - \hat{V}_{\text{VOSS},d}(s) \right| \leq \underbrace{\left| V_d^\star(s) - \hat{V}_d^{C_a}(s) \right|}_{\text{VOO bound}} + \underbrace{\left| \hat{V}_d^{C_a}(s) - \hat{V}_{\text{VOSS},d}(s) \right|}_{\text{SS bound}}. \tag{59}$$

We now have two main layers of inequality, caused by the stochastic nature of VOO action selection and the uncertainty in state transition. We will first analyze the VOO bound, then the SS bound.

*a) Step 1: VOO Bound:* The VOO bound can be further decomposed into the following terms as [12] do in their work:

$$\left| V_d^\star(s) - \hat{V}_d^{C_a}(s) \right| \leq \underbrace{\left| V_d^\star(s) - \hat{V}_d(s) \right|}_{\text{VOO Recursive bound}} + \underbrace{\left| \hat{V}_d(s) - \hat{V}_d^{C_a}(s) \right|}_{\text{VOO Regret bound}}. \tag{60}$$

This step is very close to our previous procedure in VOWSS. Here, the intermediate random variables are defined in the following manner:

$$V_d^\star(s) \equiv \max_{a \in A} Q_d^\star(s,a) = \max_{a \in A} \left\{ R(s,a) + \gamma \mathbb{E}_{s' \sim T(s,a)}[V_{d+1}^\star(s')] \right\}, \tag{61}$$

$$\hat{V}_d(s) \equiv \max_{a \in A} \hat{Q}_d(s,a) = \max_{a \in A} \left\{ R(s,a) + \gamma \mathbb{E}_{s' \sim T(s,a)}[\hat{V}_{d+1}^{C_a}(s')] \right\}, \tag{62}$$

$$\hat{V}_d^{C_a}(s) \equiv \max_{a \in VOO(A,C_a)} \hat{Q}_d(s,a) = \max_{a \in VOO(A,C_a)} \left\{ R(s,a) + \gamma \mathbb{E}_{s' \sim T(s,a)}[\hat{V}_{d+1}^{C_a}(s')] \right\}. \tag{63}$$

Repeating the procedures in Lemma 1, we closely follow the calculations from Lemma 5 in [12]:

$$\underbrace{\left| V_d^\star(s) - \hat{V}_d(s) \right|}_{\text{VOO Recursive bound}} = \left| \max_{a \in A} \left\{ R(s,a) + \gamma \mathbb{E}_{s' \sim T(s,a)}[V_{d+1}^\star(s')] \right\} - \max_{a \in A} \left\{ R(s,a) + \gamma \mathbb{E}_{s' \sim T(s,a)}[\hat{V}_{d+1}^{C_a}(s')] \right\} \right| \tag{64}$$

$$\leq \left| R(s,a^\star) + \gamma \mathbb{E}_{s' \sim T(s,a^\star)}[V_{d+1}^\star(s')] - R(s,a^\star) - \gamma \mathbb{E}_{s' \sim T(s,a^\star)}[\hat{V}_{d+1}^{C_a}(s')] \right| \qquad (a^\star = \arg\max Q_d^\star(s,a))$$

$$\leq \gamma \mathbb{E}_{s' \sim T(s,a^\star)} \left| V_{d+1}^\star(s') - \hat{V}_{d+1}^{C_a}(s') \right| \tag{65}$$

$$\leq \gamma \cdot \eta_{C_a}(d+1) \tag{66}$$

$$\underbrace{\left| \hat{V}_d(s) - \hat{V}_d^{C_a}(s) \right|}_{\text{VOO Regret bound}} \leq \mathscr{R}_{C_a}. \tag{67}$$

Thus, with our choice of $C_a$ that is designed to satisfy the recurrence relation, we obtain the bound:

$$\left| V_d^\star(s) - \hat{V}_d^{C_a}(s) \right| \leq \underbrace{\left| V_d^\star(s) - \hat{V}_d(s) \right|}_{\text{VOO Recursive bound}} + \underbrace{\left| \hat{V}_d(s) - \hat{V}_d^{C_a}(s) \right|}_{\text{VOO Regret bound}} \tag{68}$$

$$\leq \gamma \cdot \eta_{C_a}(d+1) + \mathscr{R}_{C_a} \leq \eta_{C_a}(d). \tag{69}$$

*b) Step 2: Sparse Sampling Bound:* Similar to Lemma 1, the SS bound can also be further decomposed into the following terms:

$$\left| \hat{V}_d^{C_a}(s) - \hat{V}_{\text{VOSS},d}(s) \right| \leq \underbrace{\left| \hat{V}_d^{C_a}(s) - \tilde{V}_{\text{VOSS},d}(s) \right|}_{\text{SS Concentration bound}} + \underbrace{\left| \tilde{V}_{\text{VOSS},d}(s) - \hat{V}_{\text{VOSS},d}(s) \right|}_{\text{SS Recursive bound}}. \tag{70}$$

Here, the extra intermediate random variable and the VOSS estimator are defined in the following manner:

$$\tilde{V}_{\text{VOSS},d}(s) \equiv \max_{a \in VOO(A,C_a)} \tilde{Q}_{\text{VOSS},d}(s,a) = \max_{a \in VOO(A,C_a)} \left\{ R(s,a) + \gamma \frac{1}{C_s} \sum_{i=1}^{C_s} \hat{V}_{d+1}^{C_a}(s_i') \right\}, \tag{71}$$

$$\hat{V}_{\text{VOSS},d}(s) \equiv \max_{a \in VOO(A,C_a)} \hat{Q}_{\text{VOSS},d}(s,a) = \max_{a \in VOO(A,C_a)} \left\{ R(s,a) + \gamma \frac{1}{C_s} \sum_{i=1}^{C_s} \hat{V}_{\text{VOSS},d+1}(s_i') \right\}. \tag{72}$$

We now apply the sparse sampling bound as well as the recursive bound in order to bound the SS bound components. In our case, since our intermediate sparse sampling term $\tilde{V}_{\text{VOSS},d}(s)$ is merely swapping out the expectation of $\hat{V}_{d+1}^{C_a}$ with a sample average under the appropriate sampling density, this turns out to be simply the Chernoff bound. We transform the sparse sampling concentration bound to an expected value bound by using the fact that the difference of $V$ functions/estimators is bounded above by $2V_{\max}$, setting the concentration bound to be $\varepsilon$ as per the Lemma statement and assigning the worst case result with probability $p$. Once again, we compare the two quantities by picking a reference action to directly compare the

quantities inside the max operations:

$$\underbrace{\left| \hat{V}_d^{C_a}(s) - \tilde{V}_{\text{VOSS},d}(s) \right|}_{\text{SS Concentration bound}} \leq \left| \max_{a \in VOO(A,C_a)} \{ R(s,a) + \gamma \mathbb{E}_{s' \sim T(s,a)}[\hat{V}_{d+1}^{C_a}(s')] \} - \max_{a \in VOO(A,C_a)} \{ R(s,a) + \gamma \frac{1}{C_s} \sum_{i=1}^{C_s} \hat{V}_{d+1}^{C_a}(s_i') \} \right| \tag{73}$$

$$\leq \gamma \left| \mathbb{E}_{s' \sim T(s,a^*)}[\hat{V}_{d+1}^{C_a}(s')] - \frac{1}{C_s} \sum_{i=1}^{C_s} \hat{V}_{d+1}^{C_a}(s_i') \right| \qquad (a^* = \arg\max_{VOO} \hat{Q}_d(s,a)) $$

$$\leq \gamma \cdot ((1-p)\varepsilon + 2p \cdot V_{\max}) \leq \gamma \cdot (\varepsilon + 2p \cdot V_{\max}). \tag{74}$$

For the SS Recursive bound, we proceed with the similar recursive calculation as SS Recursive bound term done in Lemma 1 by using the inductive hypothesis for step $d+1$:

$$\underbrace{\left| \tilde{V}_{\text{VOSS},d}(s) - \hat{V}_{\text{VOSS},d}(s) \right|}_{\text{SS Recursive bound}} \leq \left| \max_{a \in VOO(A,C_a)} \left\{ R(s,a) + \gamma \frac{1}{C_s} \sum_{i=1}^{C_s} \hat{V}_{d+1}^{C_a}(s_i') \right\} \right. \tag{75}$$

$$\left. - \max_{a \in VOO(A,C_a)} \left\{ R(s,a) + \gamma \frac{1}{C_s} \sum_{i=1}^{C_s} \hat{V}_{\text{VOSS},d+1}(s_i') \right\} \right| \tag{76}$$

$$\leq \gamma \frac{1}{C_s} \sum_{i=1}^{C_s} \left| \hat{V}_{d+1}^{C_a}(s_i') - \hat{V}_{\text{VOSS},d+1}(s_i') \right| \qquad (\tilde{a} = \arg\max_{VOO} \tilde{Q}_{\text{VOSS},d}(s,a))$$

$$\leq \gamma \cdot \alpha_{C_s}(d+1). \tag{77}$$

Putting the sparse sampling components together by applying the recursive definition of $\alpha_{C_s}(d)$, we obtain the recurring concentration inequality by induction:

$$\left| \hat{V}_d^{C_a}(s) - \hat{V}_{\text{VOSS},d}(s) \right| \leq \underbrace{\left| \hat{V}_d^{C_a}(s) - \tilde{V}_{\text{VOSS},d}(s) \right|}_{\text{SS Concentration bound}} + \underbrace{\left| \tilde{V}_{\text{VOSS},d}(s) - \hat{V}_{\text{VOSS},d}(s) \right|}_{\text{SS Recursive bound}} \tag{78}$$

$$\leq \gamma \cdot (\varepsilon + 2p \cdot V_{\max}) + \gamma \cdot \alpha_{C_s}(d+1) = \alpha_{C_s}(d). \ \blacksquare \tag{79}$$

*Proof for Theorem 2.* Finally, we prove the main theorem by combining the two terms:

$$\left| V_d^\star(s) - \hat{V}_{\text{VOSS},d}(s) \right| \leq \underbrace{\left| V_d^\star(s) - \hat{V}_d^{C_a}(s) \right|}_{\text{VOO bound}} + \underbrace{\left| \hat{V}_d^{C_a}(s) - \hat{V}_{\text{VOSS},d}(s) \right|}_{\text{SS bound}} \leq \eta_{C_a}(d) + \alpha_{C_s}(d) \leq \eta_{C_a} + \alpha_{C_s}. \ \blacksquare \tag{80}$$

Here, the $\eta_{C_a}(d)$ corresponds to the VOO bound in [12]. Thus, we can always find a corresponding $C_a$ for an arbitrary decreasing sequence of $\eta_{C_a}(d)$ which satisfies the properties in our lemma, and this sequence can be made closer to 0 for each depth by choosing bigger $C_a$ values.

On the other hand, $\alpha_{C_s}(d)$ corresponds to the sparse sampling bound/Chernoff bound which decreases in both $p$ and $\varepsilon$ with more samples $C_s$. Thus, it can also be made to decrease to 0 as we increase $C_s$.