How to do MCTS on a POMDP?

$Q[(s,a)]$

$N[(s,a)]$

# Online POMDP Methods

# Approximate POMDP Solutions

# Approximate POMDP Solutions

## Numerical Approximations

(approximately solve original problem)

# Approximate POMDP Solutions

## Numerical Approximations

(approximately solve original problem)

## Offline

# Approximate POMDP Solutions

## Numerical Approximations

(approximately solve original problem)

Offline

Online

# Approximate POMDP Solutions

**Numerical Approximations**

(approximately solve original problem)

**Formulation Approximations**

(solve a slightly different problem)

**Offline**

Previously

**Online**

# Approximate POMDP Solutions

**Numerical Approximations**

(approximately solve original problem)

**Formulation Approximations**

(solve a slightly different problem)

Last Time

**Offline**

Previously

**Online**

# Approximate POMDP Solutions

QMDP
Certainty Equivalence

## Numerical Approximations
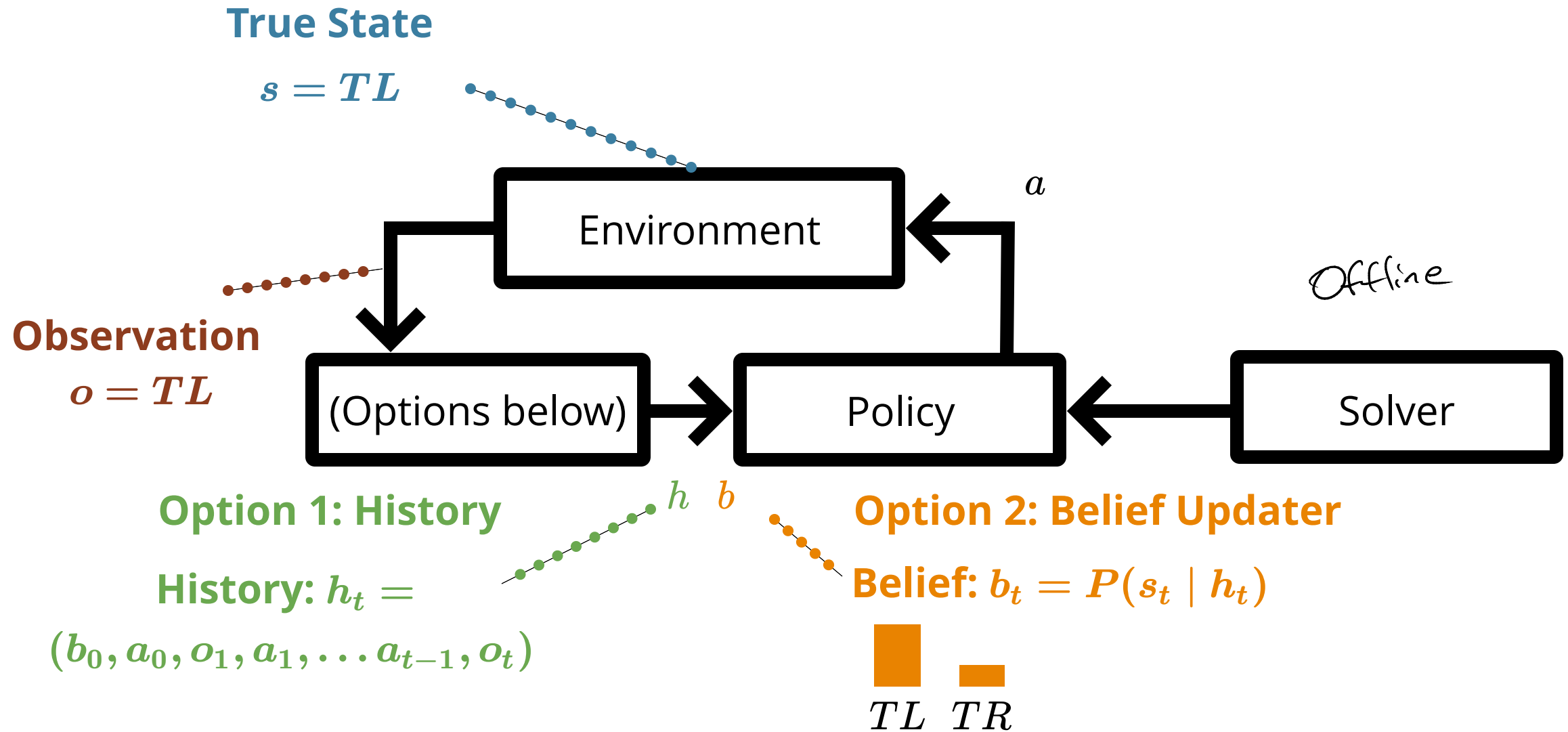
(approximately solve original problem)

## Formulation Approximations

(solve a slightly different problem)

Last Time

## Offline

Previously

SARSOP

## Online

Today!

# POMDP Sense-Plan-Act Loop

**True State**

$$s = TL$$

Environment

$a$

*Offline*

**Observation**

$$o = TL$$

(Options below)

Policy

Solver

**Option 1: History**

$h$   $b$

**Option 2: Belief Updater**

**History:** $h_t =$

$$(b_0, a_0, o_1, a_1, \ldots a_{t-1}, o_t)$$

**Belief:** $b_t = P(s_t \mid h_t)$

$TL$   $TR$

# POMDP Sense-Plan-Act Loop

**True State**

$s = TL$

Environment

$a$

**Observation**

$o = TL$

(Options below)

Planner

$b$

**Option 2: Belief Updater**

**Belief:** $b_t = P(s_t \mid h_t)$

$TL$ $TR$

# Belief-Space Tree Search: AEMS



G

$b$

$b' = \tau(b, a^i, o^i)$

$a^i$

$o^1$ $o^2$

$b^i$

while time remains

$\rightarrow b^* = \underset{b \in \text{fringe}(G)}{\text{argmax}} E(b)$

$\rightarrow$ expand $(b^*)$

backup $(b^*)$

$$E(b) = \gamma^d P(b^d) \hat{\varepsilon}(b^d)$$

$$\hat{\varepsilon}(b) = U(b) - L(b) \quad \leftarrow \text{problem -specific}$$

$$P(b^d) = \prod_{i=0}^{d-1} P(o^i | b^i, a^i) P(a^i | b^i)$$

T, Z

$$P(a|b) = \frac{U(a,b) - L(b)}{U(b) - L(b)} \quad \text{AEMS 1}$$

$$P(a|b) = \begin{cases} 1 & \text{if } a = \underset{a'}{\text{argmax}} \, U(a',b) \\ 0 & \text{o.w.} \end{cases} \quad \text{AEMS 2}$$

4

# Monte Carlo Tree Search (MCTS/UCT)
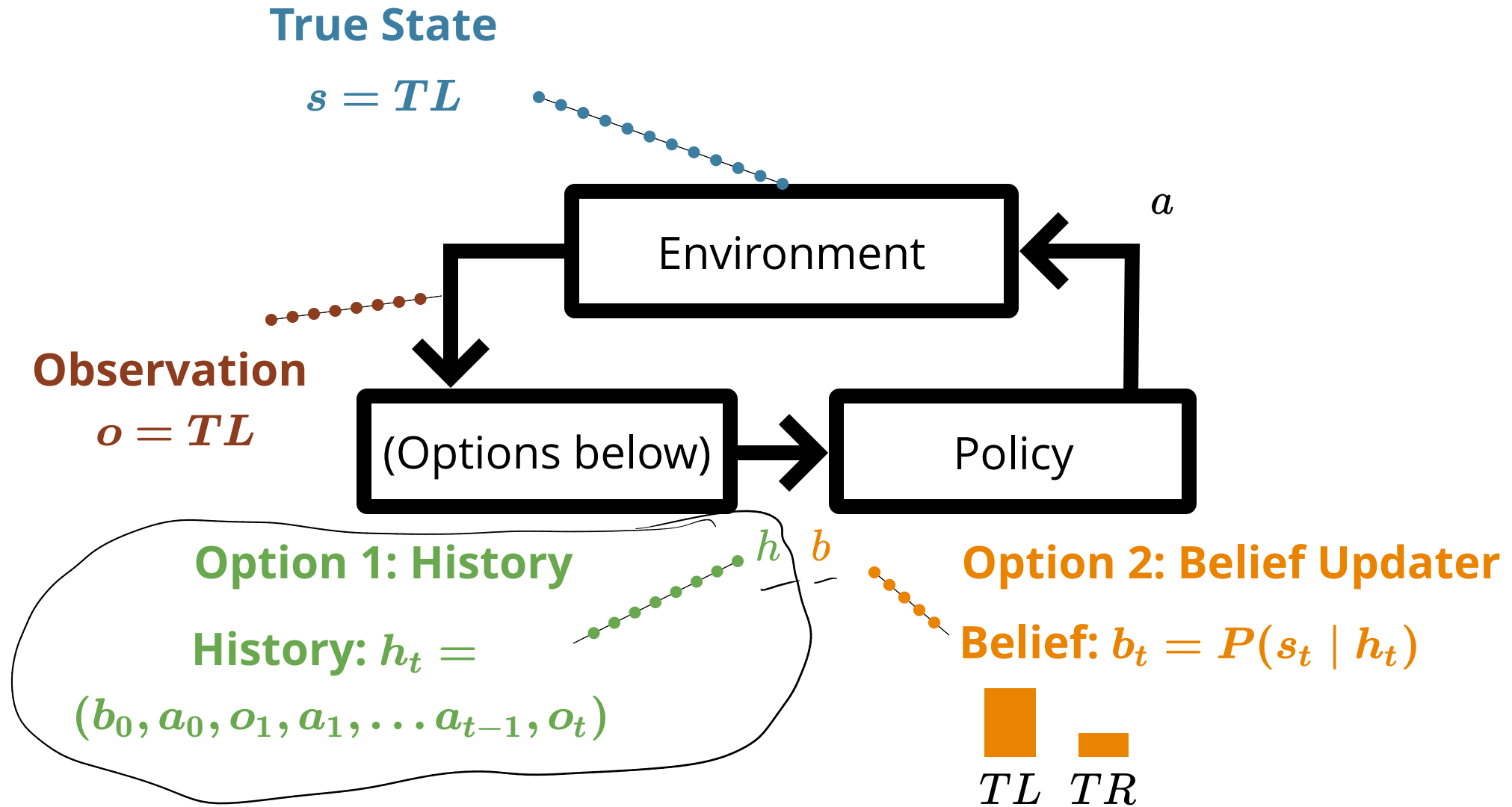
Search

Expansion

Rollout
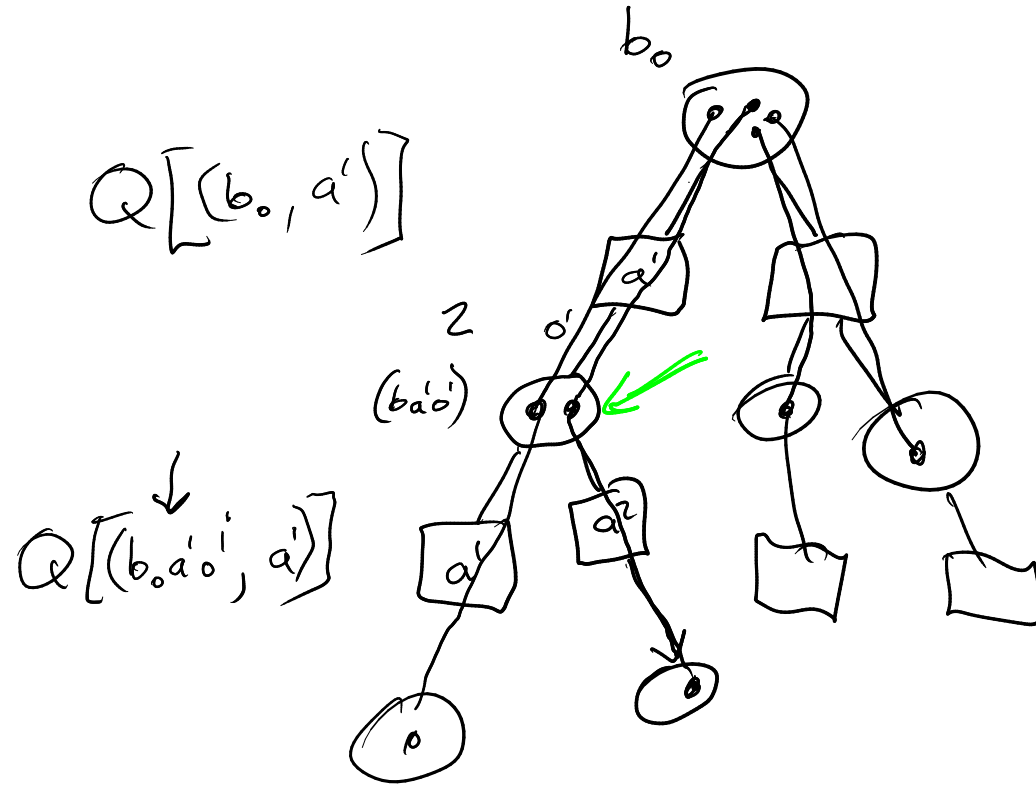
Backup



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

low $N(s,a)/N(s)$ = high bonus
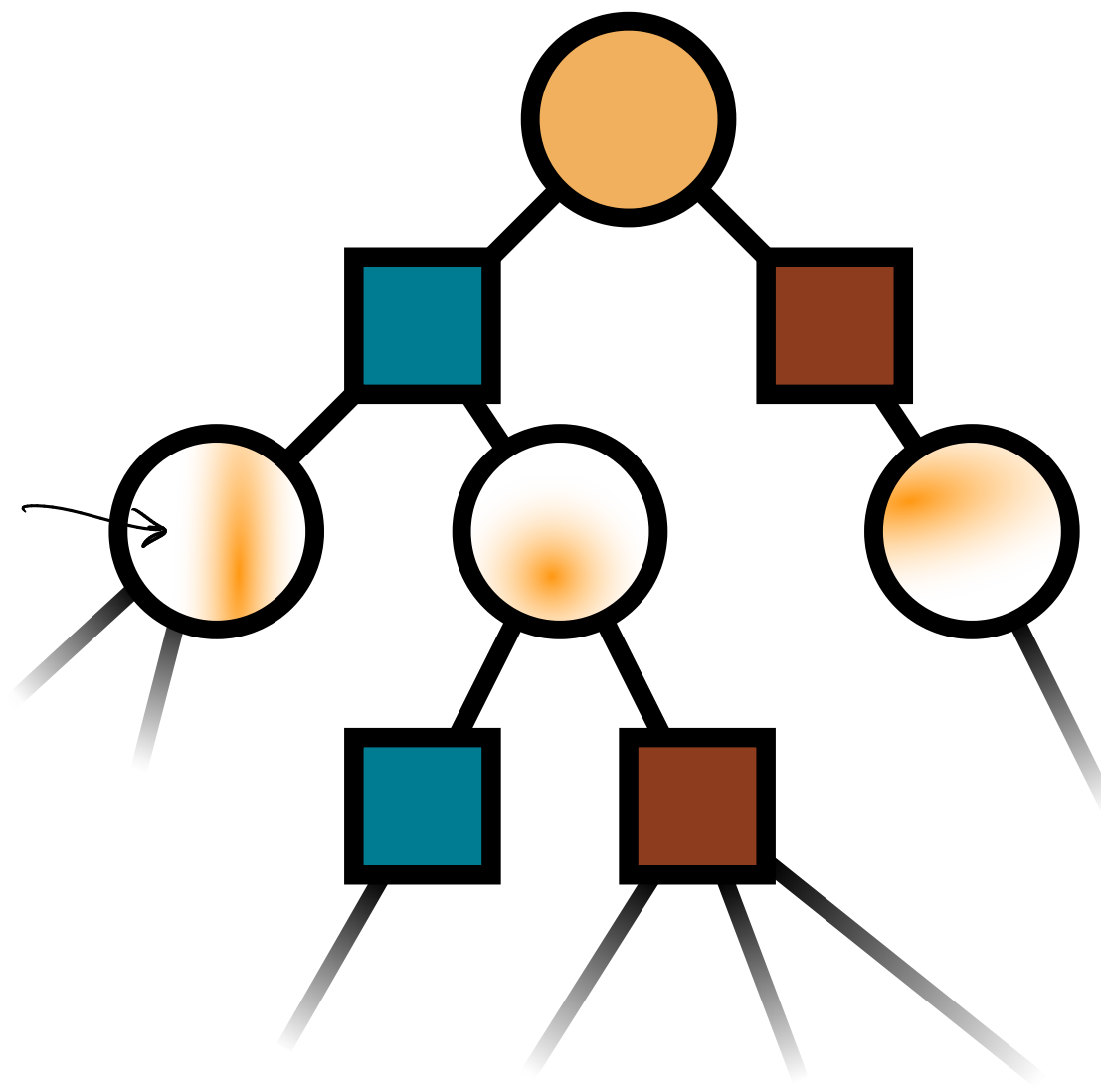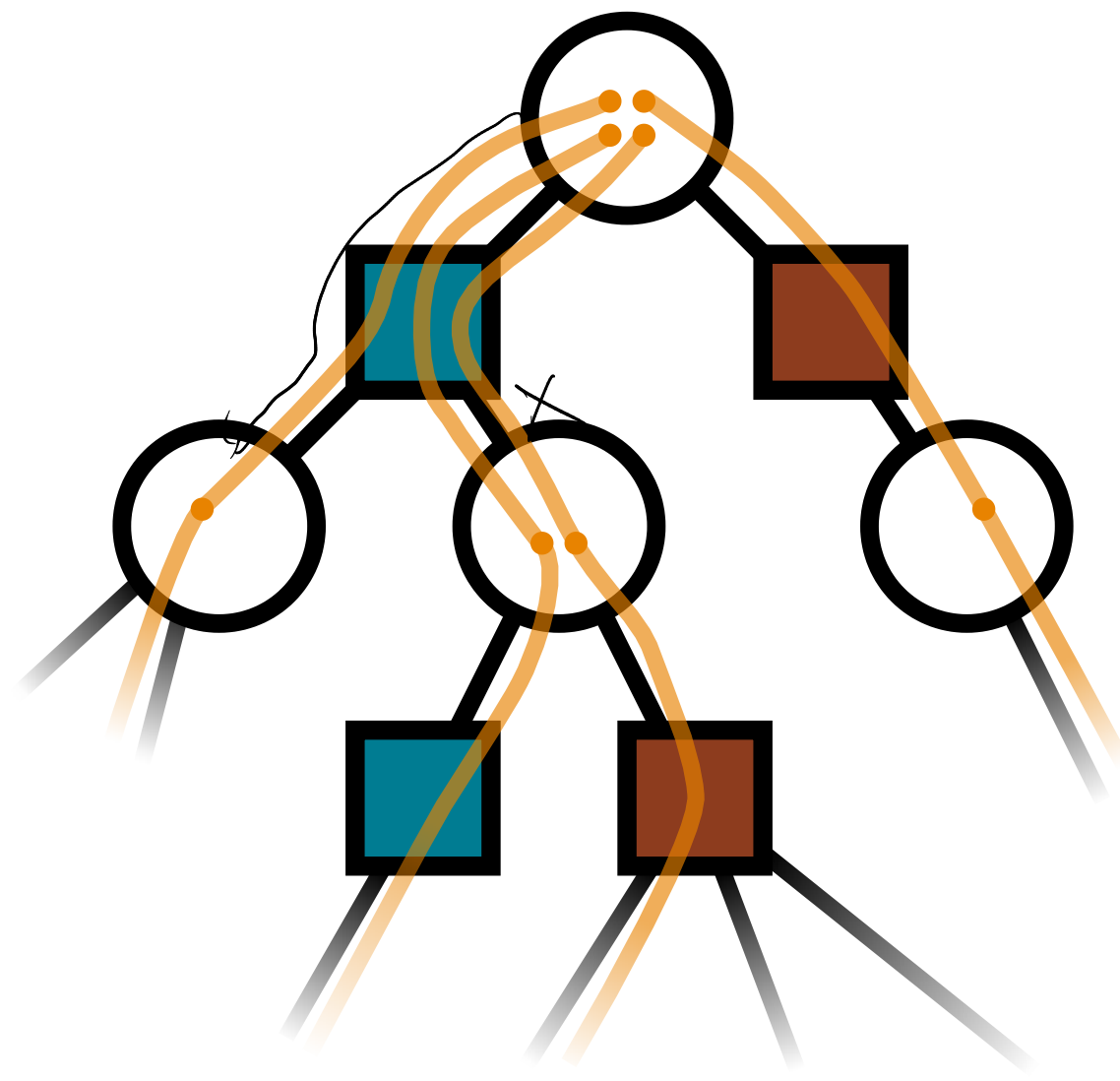
start with $c = 2(\bar{V} - \underline{V})$

update $Q(s,a)$

$U(s) = \hat{u}$

$\hat{\pi}$  $\hat{u}$↑

# How should we adapt MCTS for POMDPs?



**True State**

$$s = TL$$

Environment

$a$

**Observation**

$$o = TL$$

(Options below)

Policy

**Option 1: History**

**History:** $h_t = (b_0, a_0, o_1, a_1, \ldots a_{t-1}, o_t)$

$h$  $b$

**Option 2: Belief Updater**

**Belief:** $b_t = P(s_t \mid h_t)$

$TL$  $TR$

# MCTS on Histories

# DESPOT



Somani, A., Ye, N., Hsu, D., & Lee, W. "DESPOT : Online POMDP Planning with Regularization." *Journal of Artificial Intelligence Research*, 2017

# DESPOT



- Determinized Scenarios

Somani, A., Ye, N., Hsu, D., & Lee, W. "DESPOT : Online POMDP Planning with Regularization." *Journal of Artificial Intelligence Research*, 2017

# DESPOT



- Determinized Scenarios
- Guided by Lower and Upper Bounds

Somani, A., Ye, N., Hsu, D., & Lee, W. "DESPOT : Online POMDP Planning with Regularization." *Journal of Artificial Intelligence Research*, 2017

# POMCP

# POMCPOW

# DESPOT-$\alpha$
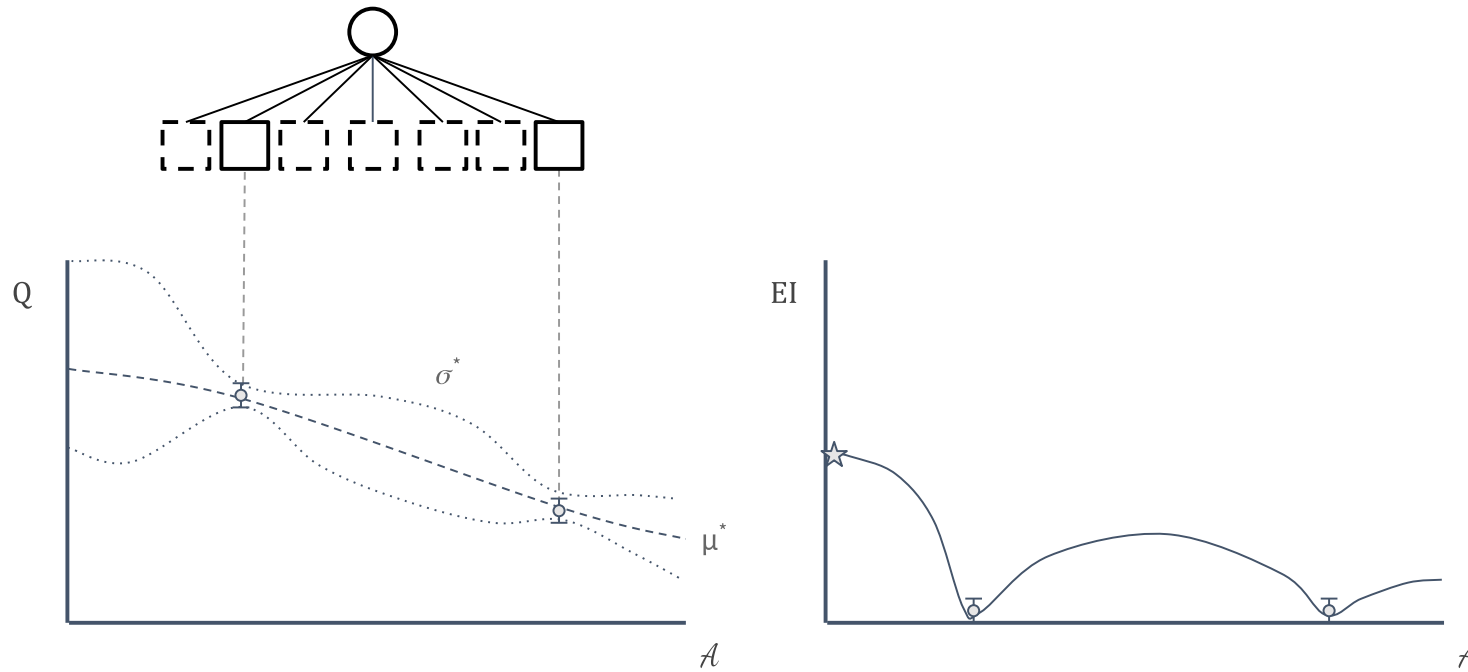
# BOMCP

**Bayesian Optimized Action Branching**



**[Mern, Sunberg, et al. AAAI 2021]**

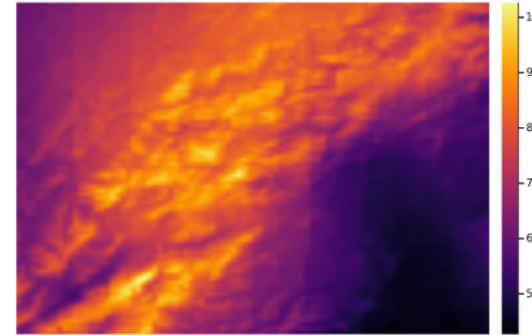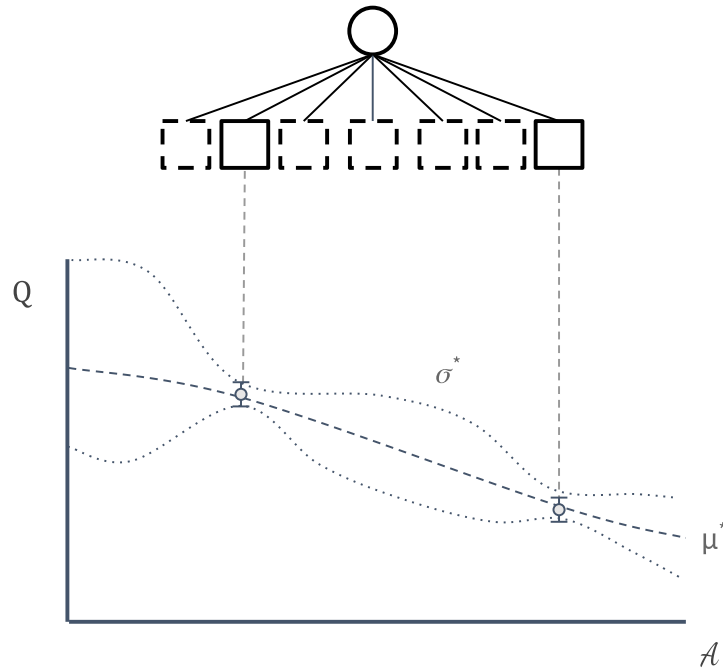# BOMCP

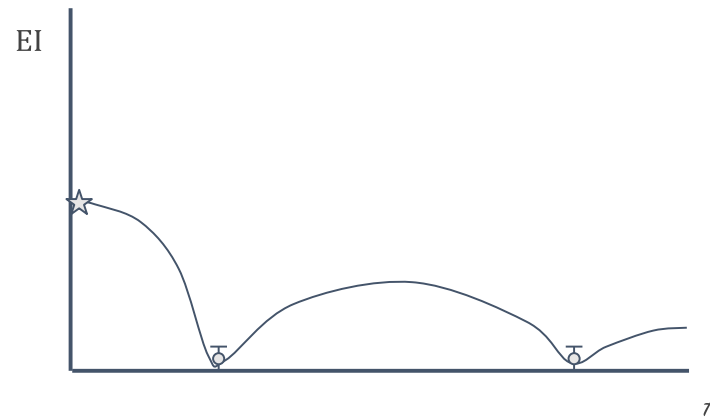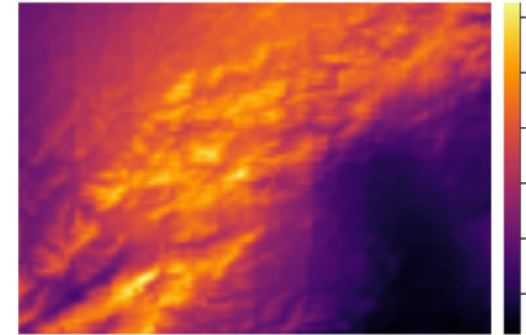**Bayesian Optimized Action Branching**



Figure 2: Wind Map. Figure shows wind map for Altamont Pass, CA at 100m altitude. The colors represent the average annual wind speed in m/s.

**[Mern, Sunberg, et al. AAAI 2021]**
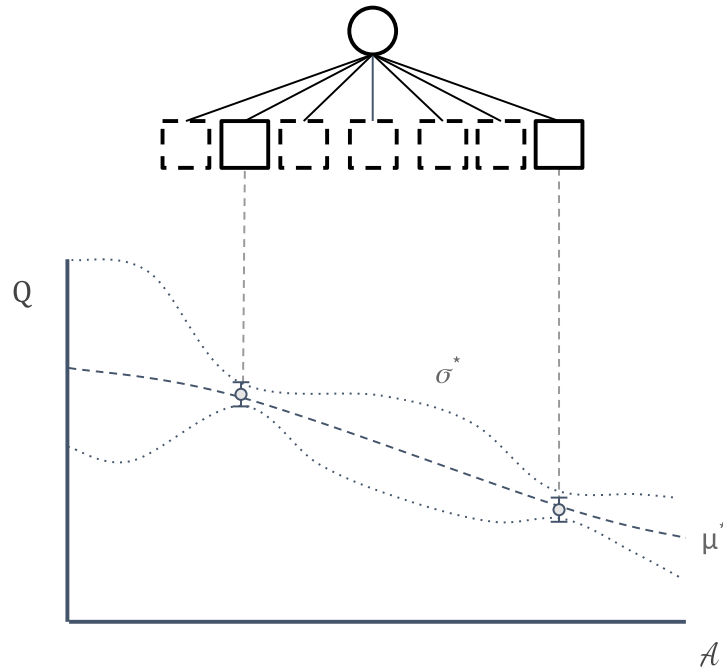
# BOMCP

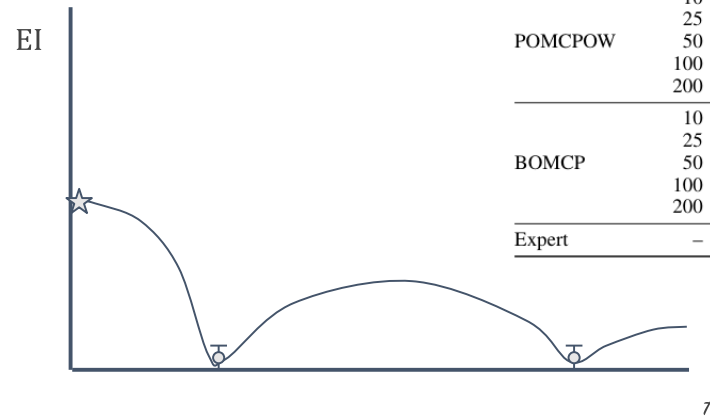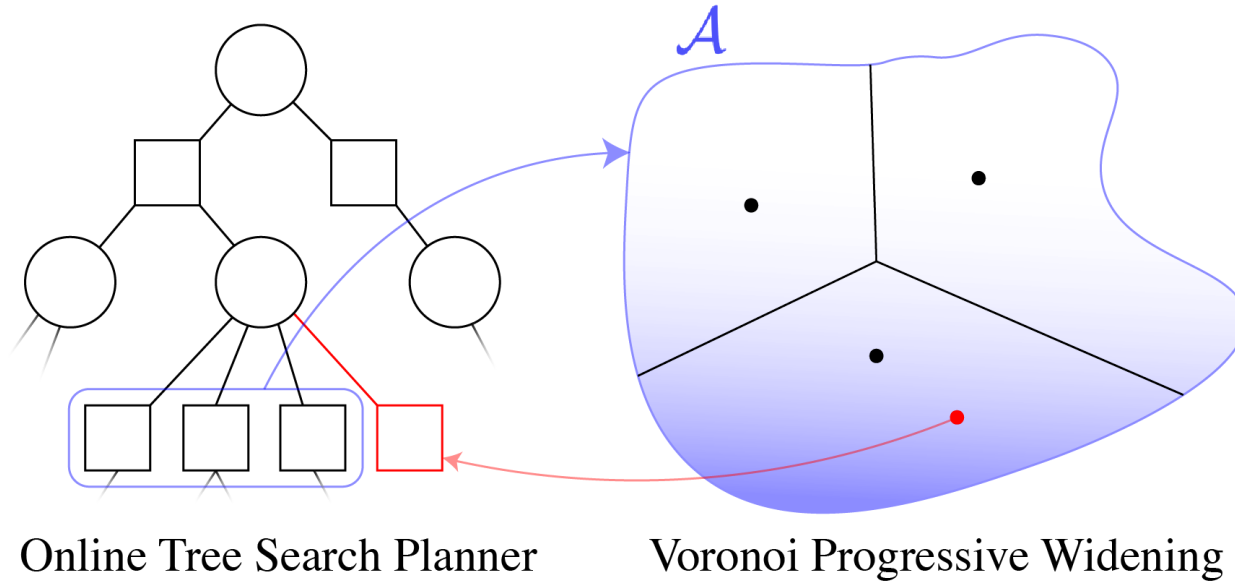**Bayesian Optimized Action Branching**



Figure 2: Wind Map. Figure shows wind map for Altamont Pass, CA at 100m altitude. The colors represent the average annual wind speed in m/s.

| Algorithm | Queries | Score | Time (seconds) |
|---|---|---|---|
| POMCPOW | 10 | 15708 ± 229 | 2.25 ± 0.07 |
| | 25 | 16234 ± 217 | 4.80 ± 0.07 |
| | 50 | 16374 ± 212 | 6.27 ± 0.08 |
| | 100 | 16018 ± 262 | 11.98 ± 0.07 |
| | 200 | 15787 ± 233 | 20.67 ± 0.09 |
| BOMCP | 10 | 18095 ± 183 | 2.55 ± 0.08 |
| | 25 | 18154 ± 158 | 5.21 ± 0.07 |
| | 50 | 18015 ± 163 | 6.71 ± 0.06 |
| | 100 | 18225 ± 119 | 13.39 ± 0.07 |
| | 200 | 18113 ± 157 | 25.14 ± 0.08 |
| Expert | – | 8130 ± 51 | – |

**[Mern, Sunberg, et al. AAAI 2021]**

# Voronoi Progressive Widening
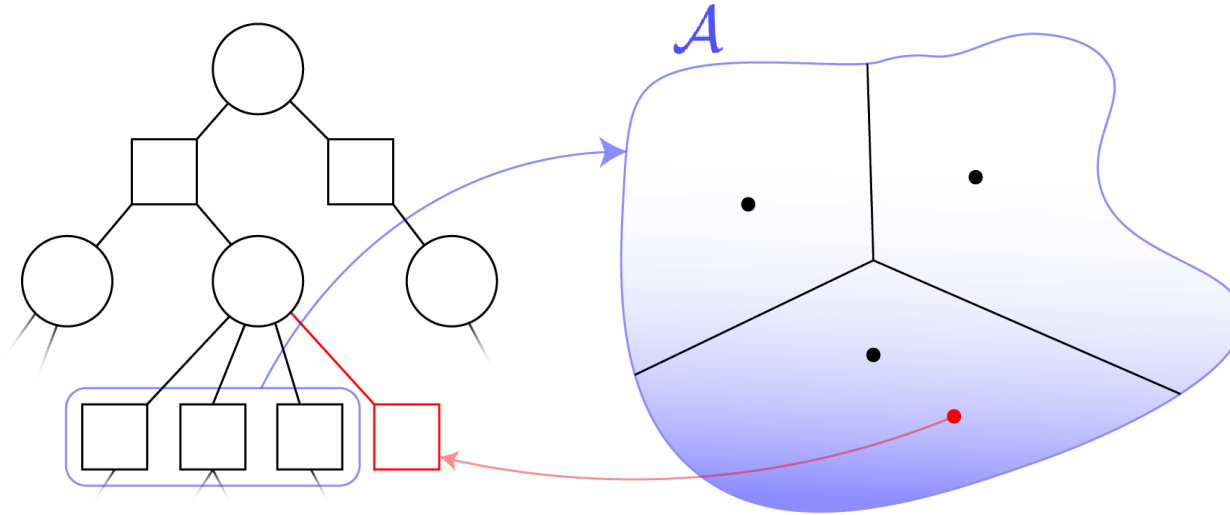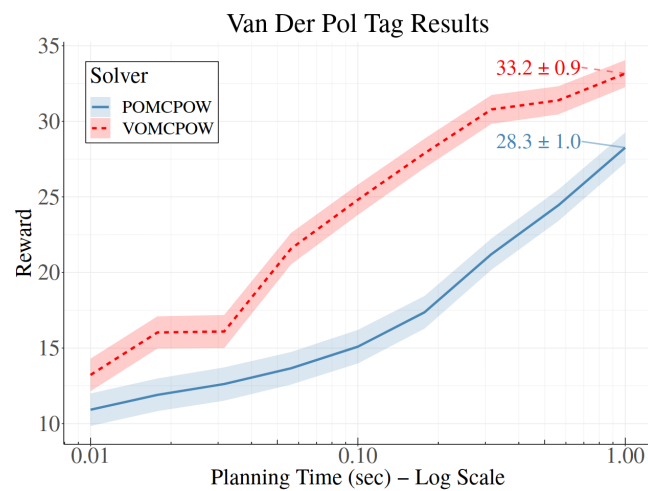


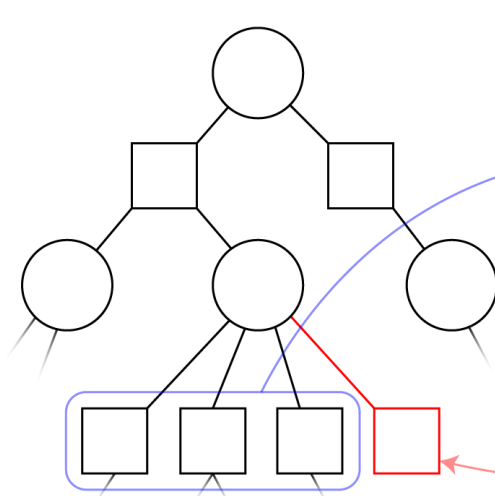Online Tree Search Planner     Voronoi Progressive Widening

**[Lim, Tomlin, & Sunberg CDC 2021]**

# Voronoi Progressive Widening
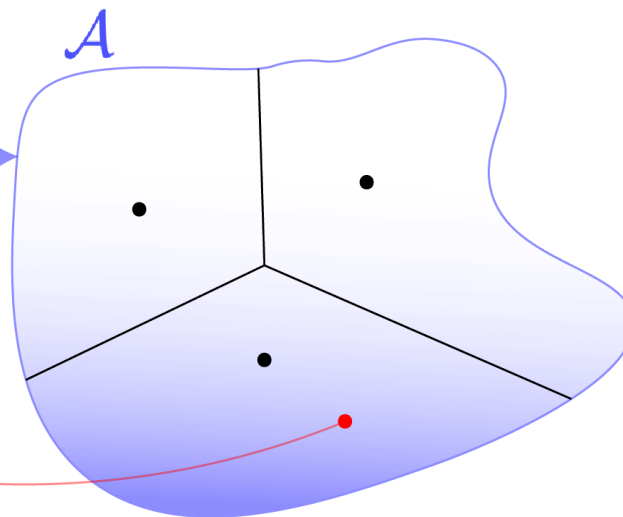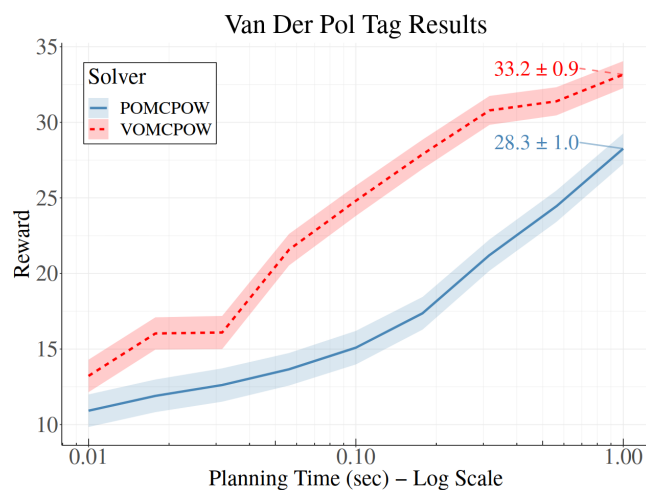


Online Tree Search Planner

Voronoi Progressive Widening

Van Der Pol Tag Results

**[Lim, Tomlin, & Sunberg CDC 2021]**

13

# Voronoi Progressive Widening



Online Tree Search Planner

Voronoi Progressive Widening



Van Der Pol Tag Results

**Theorem 2** (VOWSS Inequality). *Given the action sampling width of $C_a$ and state sampling width of $C_s$ at every height of the tree that follow the intermediate concentration bounds in the form of POWSS (Lim, Tomlin, and Sunberg 2020) and regret bounds in the form of VOO (Kim et al. 2020), the following bounds for the VOWSS estimator $\hat{V}^{C_a}_{\text{VOWSS},d}(b)$ hold for all $d \in [0, D-1]$ in expectation:*

$$\left| V^\star_d(b) - \hat{V}^{C_a}_{\text{VOWSS},d}(b) \right| \leq \eta + \alpha$$

**[Lim, Tomlin, & Sunberg CDC 2021]**

13