

# Adaptive Information Belief Space Planning

Moran Barenboim<sup>1</sup> and Vadim Indelman<sup>2</sup>

<sup>1</sup>Technion Autonomous Systems Program

<sup>2</sup>Department of Aerospace Engineering

Technion - Israel Institute of Technology, Haifa 32000, Israel

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

## A Proofs

### A.1 Lemma 1

The proof is provided for continuous state space; The discrete case obtained similarly by changing integrals to summations.

*Proof.*

$$\begin{aligned} & \sum_{n=1}^{N_o} \bar{\mathbb{P}}(o^n | H^-) \mathcal{H}(\bar{b}) = \\ & - \sum_{n=1}^{N_o} \bar{\mathbb{P}}(o^n | H^-) \int_s \bar{\mathbb{P}}(s | H) \cdot \log(\bar{\mathbb{P}}(s | H)) \end{aligned} \quad (1)$$

applying Bayes' rule for  $\bar{\mathbb{P}}(s | H)$ ,

$$\begin{aligned} & - \sum_{n=1}^{N_o} \int_s \bar{Z}(o^n | s) \mathbb{P}(s | H^-) \\ & \cdot \log \left( \frac{\bar{Z}(o^n | s) \mathbb{P}(s | H^-)}{\int_{s'} \bar{Z}(o^n | s') \mathbb{P}(s' | H^-)} \right) \end{aligned} \quad (2)$$

Splitting summation to follow the partitioning of the abstract observation model,

$$\begin{aligned} & - \sum_{c=1}^C \sum_{k=K(c-1)+1}^{Kc} \int_s \bar{Z}(o^k | s) \mathbb{P}(s | H^-) \\ & \cdot \log \left( \frac{\bar{Z}(o^k | s) \mathbb{P}(s | H^-)}{\int_{s'} \bar{Z}(o^k | s') \mathbb{P}(s' | H^-)} \right) \end{aligned} \quad (3)$$

By construction,  $\bar{Z}(o | s)$  has uniform distribution for  $o^k$ , where  $k \in [K(c-1)+1, Kc]$ . Thus,

$$- \sum_{c=1}^C \left[ \sum_{k=K(c-1)+1}^{Kc} 1 \right] \int_s \bar{Z}(o^{Kc} | s) \mathbb{P}(s | H^-) \quad (4)$$

$$\begin{aligned} & \cdot \log \left( \frac{\bar{Z}(o^{Kc} | s) \mathbb{P}(s | H^-)}{\int_{s'} \bar{Z}(o^{Kc} | s') \mathbb{P}(s' | H^-)} \right) = \\ & - \sum_{c=1}^C K \cdot \int_s \bar{Z}(o^{Kc} | s) \mathbb{P}(s | H^-) \end{aligned} \quad (5)$$

$$\begin{aligned} & \cdot \log \left( \frac{\bar{Z}(o^{Kc} | s) \mathbb{P}(s | H^-)}{\int_{s'} \bar{Z}(o^{Kc} | s') \mathbb{P}(s' | H^-)} \right) = \\ & K \cdot \sum_{c=1}^C \bar{\mathbb{P}}(o^{Kc} | H^-) \mathcal{H}(\bar{b}) \end{aligned}$$

which concludes the proof.  $\square$

## A.2 Lemma 2

*Proof.* We begin with,

$$\bar{\mathbb{E}}_o [\mathbb{E}_{s \sim \bar{b}} [r(s, a)]] \quad (6)$$

by definition of  $\bar{\mathbb{E}}_o [\cdot]$ ,  $\bar{b}(s)$ ,

$$\sum_{n=1}^{N_o} \bar{\mathbb{P}}(o^n | H^-) \left[ \sum_{s \in S} \bar{\mathbb{P}}(s | o^n, H^-) r(s, a) \right] \quad (7)$$

applying chain rule,

$$\begin{aligned} \sum_{n=1}^{N_o} \sum_{s \in S} \bar{\mathbb{P}}(s, o^n | H^-) r(s, a) = & \quad (8) \\ \sum_{n=1}^{N_o} \sum_{s \in S} \bar{Z}(o^n | s) b^-(s) r(s, a) \end{aligned}$$

we split the sum over the observations to comply with the abstraction partitioning and use the the abstract observation model definition, (4),

$$\sum_{s \in S} \sum_{c=1}^C \sum_{k=K(c-1)+1}^{Kc} \frac{\sum_{m=K(c-1)+1}^{Kc} Z(o^m | s)}{K} b^-(s) r(s, a) \quad (9)$$

we then arrive at the desired result,

$$\sum_{s \in S} \sum_{n=1}^{N_o} Z(o^n | s) b^-(s) r(s, a) = \quad (10)$$

$$\bar{\mathbb{E}}_o [\mathbb{E}_{s \sim \bar{b}} [r(s, a)]] \quad (11)$$

□

## A.3 Theorem 1

*Proof.* For clarity, we omit the time index in the derivation, the result holds for any time step. We use  $H^-$  to denote past history while excluding last observation. We also use  $b$  and  $\mathbb{P}(s | o, H^-)$  interchangeably. Rearranging the abstraction from (1),

$$\sum_{k=1}^K \bar{Z}(o^k | s) = K \cdot \bar{Z}(o^b | s) \doteq \sum_{k=1}^K Z(o^k | s) \quad \forall b \in [1, K]$$

Plugging it to the expected entropy term, (11),

$$\bar{\mathbb{E}}_o [\mathcal{H}(\bar{b})] - \mathbb{E}_o [\mathcal{H}(b)] = \quad (12)$$

$$\sum_{i=1}^{N_o} \bar{\mathbb{P}}(o_i | H^-) \mathcal{H}(\bar{b}) - \sum_{i=1}^{N_o} \mathbb{P}(o_i | H^-) \mathcal{H}(b) \quad (13)$$

expanding the entropy term,

$$\begin{aligned} - \sum_{i=1}^{N_o} \bar{\mathbb{P}}(o_i | H^-) \int_s \bar{\mathbb{P}}(s | o_i, H^-) \log(\bar{b}) & \quad (14) \\ + \sum_{i=1}^{N_o} \mathbb{P}(o_i | H^-) \int_s \mathbb{P}(s | o_i, H^-) \log(b) \end{aligned}$$

by Bayes' rule,

$$\begin{aligned} - \sum_{i=1}^{N_o} \int_s \bar{Z}(o_i | s) \mathbb{P}(s | H^-) \log(\bar{b}) & \quad (15) \\ + \sum_{i=1}^{N_o} \int_s Z(o_i | s) \mathbb{P}(s | H^-) \log(b) \end{aligned}$$

a change in the order of summation and integral and a split of  $N_o = C \cdot K$  result in,

$$\begin{aligned} - \int_s \sum_{c=1}^C \sum_{k=K(c-1)+1}^{Kc} \bar{Z}(o_k | s) \mathbb{P}(s | H^-) \log(\bar{b}) & \quad (16) \\ + \int_s \sum_{c=1}^C \sum_{k=K(c-1)+1}^{Kc} Z(o_k | s) \mathbb{P}(s | H^-) \log(b) \end{aligned}$$

By plugging-in the definition of the abstract model,

$$- \int_s \sum_{c=1}^C \left[ \sum_{k=K(c-1)+1}^{Kc} \frac{\sum_{\bar{k}=K(c-1)+1}^{Kc} Z(o_{\bar{k}} | s)}{K} \right] \quad (17)$$

$$\begin{aligned} & \mathbb{P}(s | H^-) \log(\bar{b}) \\ & + \int_s \sum_{c=1}^C \sum_{k=K(c-1)+1}^{Kc} Z(o_k | s) \mathbb{P}(s | H^-) \log(b) = \\ & - \int_s \sum_{c=1}^C \sum_{k=K(c-1)+1}^{Kc} Z(o_k | s) \mathbb{P}(s | H^-) \log(\bar{b}) \quad (18) \end{aligned}$$

$$\begin{aligned} & + \int_s \sum_{c=1}^C \sum_{k=K(c-1)+1}^{Kc} Z(o_k | s) \mathbb{P}(s | H^-) \log(b) = \\ & \sum_{i=1}^{N_o} \mathbb{P}(o_i | H^-) \int_s b \cdot \log\left(\frac{b}{\bar{b}}\right) = \quad (19) \\ & E_o [\mathcal{D}_{KL}(b || \bar{b})] \geq 0 \end{aligned}$$

(19) obtained by applying similar steps in reverse order. The last equality holds since KL-divergence is non-negative and so is its expectation. It is left to prove the upper bound; Applying Bayes rule to the nominator and denominator of (19),

$$\begin{aligned} & \sum_{i=1}^{N_o} \mathbb{P}(o_i | H^-) \int_s b \log\left(\frac{Z(o_i | s)}{\bar{Z}(o_i | s)}\right) \quad (20) \\ & + \sum_{i=1}^{N_o} \mathbb{P}(o_i | H^-) \log\left(\frac{\bar{\mathbb{P}}(o_i | H^-)}{\mathbb{P}(o_i | H^-)}\right) \int_s b ds \end{aligned}$$

By construction of the abstract observation model,

$$\begin{aligned} & \sum_{c=1}^C \sum_{k=K(c-1)+1}^{Kc} \mathbb{P}(o_k | H^-) \int_s b \cdot \log\left(\frac{Z(o_k | s) \cdot K}{\sum_{\bar{k}=K(c-1)+1}^{Kc} Z(o_{\bar{k}} | s)}\right) ds \\ & + \sum_{c=1}^C \sum_{k=K(c-1)+1}^{Kc} \mathbb{P}(o_k | H^-) \log\left(\frac{\bar{\mathbb{P}}(o_k | H^-)}{\mathbb{P}(o_k | H^-)}\right) \\ & \leq \log(K) \sum_{c=1}^C \sum_{k=K(c-1)+1}^{Kc} \mathbb{P}(o_k | H^-) \int_s b ds + 0 = \log(K). \end{aligned}$$

The inequality is due to positiveness of the denominator in the first term and Jensen's inequality in the second term. we end up with,

$$0 \leq \bar{\mathbb{E}}_o [\mathcal{H}(\bar{b})] - \mathbb{E}_o [\mathcal{H}(b)] \leq \log(K). \quad (21)$$

□

#### A.4 Corollary 1.1

*Proof.* From Lemma 2 it is clear that the expected state-dependent reward is unaffected by the abstraction, and thus will not affect the value function. For the sake of conciseness and clarity, we prove the case that the value function depends only on the entropy. The general case derived similarly by applying Lemma 2 instead of the expected state-dependent reward.

$$V^\pi(b_t) = \sum_{n=1}^{N_o} \mathbb{P}(o_{t+1}^n | H_{t+1}^-) [-\mathcal{H}(b_{t+1}) + V^\pi(b_{t+1})]$$

expanding the value function,

$$\begin{aligned} & - \sum_{n=1}^{N_o} \mathbb{P}(o_{t+1}^n | H_{t+1}^-) [\mathcal{H}(b_{t+1}) \\ & + \sum_{n'=1}^{N_o} \mathbb{P}(o_{t+2}^{n'} | H_{t+2}^-) \mathcal{H}(b_{t+2}) + \dots] \quad (22) \end{aligned}$$

by linearity of expectation,

$$- \mathbb{E}_{o_{t+1}} [\mathcal{H}(b_{t+1})] + \mathbb{E}_{o_{t+1}} [\mathbb{E}_{o_{t+2}} [\mathcal{H}(b_{t+2})]] + \dots \quad (23)$$

using Theorem 1 for each of the expected entropy terms separately until time-step  $\mathcal{T} - 1$ ,

$$V^\pi(b_t) \geq - [\bar{\mathbb{E}}[\mathcal{H}(\bar{b}_{t+1})] + \log(K)] \quad (24)$$

$$\begin{aligned} & - \mathbb{E}_{o_{t+1}} [\bar{\mathbb{E}}_{o_{t+2}} [\mathcal{H}(\bar{b}_{t+2})] + \log(K)] \dots \\ & = - \bar{\mathbb{E}}_{o_{t+1}} [\mathcal{H}(\bar{b}_{t+1})] \quad (25) \\ & - \mathbb{E}_{o_{t+1}} \bar{\mathbb{E}}_{o_{t+2}} [\mathcal{H}(\bar{b}_{t+2})] \dots + \mathcal{T} \cdot \log(K) \end{aligned}$$

applying similar steps in reverse order yields the abstract value function,

$$\begin{aligned} & \bar{V}^\pi(b_t) + \mathcal{T} \cdot \log(K) \\ & \implies \bar{V}^\pi(b_t) - V^\pi(b_t) \leq \mathcal{T} \cdot \log(K). \end{aligned}$$

Following the same derivation and applying the other side of the inequality of Theorem 1, completes the derivation for the entropy as reward. Using the more general reward, (1), and applying Lemma 2, yields the proof for corollary 1.1,

$$0 \leq \bar{V}^\pi(b_t) - V^\pi(b_t) \leq \mathcal{T} \cdot \omega_2 \log(K).$$

□

## A.5 Expected Entropy Estimation

We derive an estimator to the expected differential entropy with continuous observation space. The discrete state or observation spaces follows similar derivation by replacing integrals with summations.

$$\mathbb{E}[\mathcal{H}(b_t)] = - \int_{o_t} p(o_t | H_t^-) \int_{s_t} p(s_t | o_t, H_t^-) \cdot \log(p(s_t | o_t, H_t^-)) \quad (26)$$

applying Bayes' rule,

$$\begin{aligned} \mathbb{E}[\mathcal{H}(b_t)] &= - \int_{o_t} \int_{s_t} p(s_t, o_t | H_t^-) \quad (27) \\ &\cdot \log \left( Z(o_t | s_t) \int_{s_{t-1}} T(s_t | s_{t-1}, a_{t-1}) b(s_{t-1}) \right) \\ &+ \int_{o_t} \int_{s_t} p(s_t, o_t | H_t^-) \cdot \log(p(o_t | H_t^-)) \end{aligned}$$

by chain rule and marginalization,

$$\begin{aligned} \mathbb{E}[\mathcal{H}(b_t)] &= - \int_{o_t} \int_{s_t} Z(o_t | s_t) b^-(s_t) \quad (28) \\ &\cdot \log \left( Z(o_t | s_t) \int_{s_{t-1}} T(s_t | s_{t-1}, a_{t-1}) b(s_{t-1}) \right) \\ &+ \int_{o_t} \int_{s_t} Z(o_t | s_t) b^-(s_t) \\ &\cdot \log \left( \int_{s_t} Z(o_t | s_t) b^-(s_t) \right) \end{aligned}$$

using particle filter, the belief represented as a set of weighted particles,  $\{(s^1, q^1), \dots, (s^i, q^i), \dots, (s^n, q^n)\}$ . Where  $q^i$  denotes the weight of particle  $i$ .

$$\begin{aligned} \mathbb{E}[\mathcal{H}(b_t)] &\approx - \int_{o_t} \eta_t \sum_{i=1}^n Z(o_t | s_t^i) q_{t-1}^i \quad (29) \\ &\cdot \log \left( Z(o_t | s_t^i) \sum_{j=1}^n p(s_t^i | s_{t-1}^j, a_{t-1}) q_{t-1}^j \right) \\ &+ \int_{o_t} \eta_t \sum_{i=1}^n Z(o_t | s_t^i) q_{t-1}^i \\ &\cdot \log \left( \sum_i Z(o_t | s_t^i) q_{t-1}^i \right) \end{aligned}$$

where  $\eta_t = \int_{o_t} \sum_{i=1}^n Z(o_t | s_t^i) q_{t-1}^i$  normalizes the estimator for the probability function so that it sums to 1. Then, we approximate expectation over the observation space using observation samples, and query the likelihood model conditioned on the state samples,  $Z(o^m | s^i) \quad \forall o^m \in$

$\{o^1, \dots, o^M\}$ ,

$$\hat{\mathbb{E}}[\mathcal{H}(\hat{b}_t)] = -\bar{\eta}_t \sum_{m=1}^M \sum_{i=1}^n Z(o_t^m | s_t^i) q_{t-1}^i \quad (30)$$

$$\begin{aligned} &\cdot \log \left( Z(o_t^m | s_t^i) \sum_{j=1}^n T(s_t^i | s_{t-1}^j, a_{t-1}) q_{t-1}^j \right) \\ &+ \bar{\eta}_t \sum_{m=1}^M \left[ \sum_{i=1}^n Z(o_t^m | s_t^i) q_{t-1}^i \right] \\ &\cdot \log \left( \sum_{i'=1}^n Z(o_t^m | s_t^{i'}) q_{t-1}^{i'} \right) \\ \bar{\eta}_t &= \frac{1}{\sum_{m=1}^M \sum_{i=1}^n Z(o_t^m | s_t^i) q_{t-1}^i} \quad (31) \end{aligned}$$

which concludes the derivation.

## A.6 Theorem 2

Note that the reward function, (1), is built of two terms, state dependent reward and entropy,

$$R(b, a, b') = \omega_1 \mathbb{E}_{s \sim b'} [r(s, a)] + \omega_2 \mathcal{H}(b'). \quad (32)$$

For clarity, we divide the proof into two parts,

$$\hat{\mathbb{E}}_o [R(\hat{b}, a, \hat{b}')] - \hat{\mathbb{E}}_o [R(\hat{b}, a, \hat{b}')] = \quad (33)$$

$$\omega_1 \left( \hat{\mathbb{E}}_o \left[ \mathbb{E}_{s \sim \hat{b}'} [r(s, a)] \right] - \hat{\mathbb{E}}_o \left[ \mathbb{E}_{s \sim \hat{b}'} [r(s, a)] \right] \right) \quad (34)$$

$$+ \omega_2 \left( \hat{\mathbb{E}}_o \left[ \mathcal{H}(\hat{b}') \right] - \hat{\mathbb{E}}_o \left[ \mathcal{H}(\hat{b}') \right] \right).$$

The first is about the difference in expected entropy, which is similar in spirit to the proof of Theorem 1. The second follows the claim and proof of Lemma (2). We begin with the difference of the expected entropy. For clarity, we derive the upper and lower bounds separately. For the upper bound,

*Proof.* In the following we directly plug-in the expected entropy estimator, with both the abstract observation model and the original observation model. For clarity, we split the expression into two parts and deal with each separately.

$$\begin{aligned} & \hat{\mathbb{E}}_o \left[ \mathcal{H}(\hat{b}') \right] - \hat{\mathbb{E}}_o \left[ \mathcal{H}(\hat{b}') \right] \\ = & - \underbrace{\bar{\eta}_t \sum_{m=1}^M \sum_{i=1}^n \bar{Z}(o_t^m | s_t^i) q_{t-1}^i}_{(a)} \quad (35) \\ & \cdot \underbrace{\log \left( \bar{Z}(o_t^m | s_t^i) \sum_{j=1}^n T(s_t^i | s_{t-1}^j, a_{t-1}) q_{t-1}^j \right)}_{(a)} \\ & + \underbrace{\bar{\eta}_t \sum_{m=1}^M \sum_{i=1}^n \bar{Z}(o_t^m | s_t^i) q_{t-1}^i \cdot \log \left( \sum_{i=1}^n \bar{Z}(o_t^m | s_t^i) q_{t-1}^i \right)}_{(b)} \\ & + \underbrace{\bar{\eta}_t \sum_{m=1}^M \sum_{i=1}^n Z(o_t^m | s_t^i) q_{t-1}^i}_{(a)} \\ & \cdot \underbrace{\log \left( Z(o_t^m | s_t^i) \sum_{j=1}^n T(s_t^i | s_{t-1}^j, a_{t-1}) q_{t-1}^j \right)}_{(a)} \\ & - \underbrace{\bar{\eta}_t \sum_{m=1}^M \sum_{i=1}^n Z(o_t^m | s_t^i) q_{t-1}^i \cdot \log \left( \sum_{i=1}^n Z(o_t^m | s_t^i) q_{t-1}^i \right)}_{(b)} \end{aligned}$$

In the first expression we start by splitting the summation to sum over its clusters and sum over the components of each

cluster,

$$(a) = \bar{\eta}_t \sum_{c=1}^C \sum_{k=K(c-1)+1}^{K \cdot c} \sum_{i=1}^n Z(o_t^k | s_t^i) q_{t-1}^i \quad (36)$$

$$\cdot \log \left( \frac{Z(o_t^k | s_t^i) \sum_{j=1}^n p(s_t^i | s_{t-1}^j, a_{t-1}) q_{t-1}^j}{\bar{Z}(o_t^k | s_t^i) \sum_{j=1}^n p(s_t^i | s_{t-1}^j, a_{t-1}) q_{t-1}^j} \right)$$

$$(a) = \bar{\eta}_t \sum_{c=1}^C \sum_{k=K(c-1)+1}^{K \cdot c} \sum_{i=1}^n Z(o_t^k | s_t^i) q_{t-1}^i \quad (37)$$

$$\cdot \log \left( \frac{Z(o_t^k | s_t^i)}{\bar{Z}(o_t^k | s_t^i)} \right)$$

using the abstract model, (4),

$$(a) = \bar{\eta}_t \sum_{c=1}^C \sum_{k=K(c-1)+1}^{K \cdot c} \sum_{i=1}^n Z(o_t^k | s_t^i) q_{t-1}^i \quad (38)$$

$$\cdot \log \left( \frac{K \cdot Z(o_t^k | s_t^i)}{\sum_{k=K(c-1)+1}^{K \cdot c} Z(o_t^k | s_t^i)} \right). \quad (39)$$

since the denominator within the log is a sum of positive values, the following clearly holds,

$$(a) \leq \bar{\eta}_t \sum_{c=1}^C \sum_{k=K(c-1)+1}^{K \cdot c} \sum_{i=1}^n Z(o_t^k | s_t^i) q_{t-1}^i \cdot \log(K) \quad (40)$$

by taking the constant  $\log(K)$  out of the summation, the rest sums to one, so  $(a) \leq \log(K)$ . Next we bound the second expression from above,

$$(b) = \bar{\eta}_t \sum_{c=1}^C \sum_{k=K(c-1)+1}^{K \cdot c} \sum_{i=1}^n Z(o_t^k | s_t^i) q_{t-1}^i \quad (41)$$

$$\cdot \log \left( \frac{\sum_{i=1}^n \bar{Z}(o_t^k | s_t^i) q_{t-1}^i / \bar{\eta}_t}{\sum_{i=1}^n Z(o_t^k | s_t^i) q_{t-1}^i / \bar{\eta}_t} \right)$$

applying Jensen's inequality,

$$(b) \leq \log \left( \bar{\eta}_t \sum_{c=1}^C \sum_{k=K(c-1)+1}^{K \cdot c} \sum_{i=1}^n \bar{Z}(o_t^k | s_t^i) q_{t-1}^i \right) \quad (42)$$

by recalling the definition of the normalizer, we end up with  $\log(1) = 0$   $\square$

Last, we provide a proof for the lower bound,

*Proof.*

$$\begin{aligned} & \hat{\mathbb{E}}_o \left[ \mathcal{H} \left( \hat{b} \right) \right] - \hat{\mathbb{E}}_o \left[ \mathcal{H} \left( \hat{b} \right) \right] = \\ & -\bar{\eta}_t \sum_{c=1}^C \sum_{k=K(c-1)+1}^{K \cdot c} \sum_{i=1}^n Z \left( o_t^k \mid s_t^i \right) q_{t-1}^i \\ & \cdot \log \left[ \frac{\bar{Z} \left( o_t^k \mid s_t^i \right) \sum_{j=1}^n T \left( s_t^i \mid s_{t-1}^j, a_{t-1} \right) q_{t-1}^j}{\sum_{i=1}^n \bar{Z} \left( o_t^k \mid s_t^i \right) q_{t-1}^i} \right] \\ & \cdot \frac{\sum_{i=1}^n Z \left( o_t^k \mid s_t^i \right) q_{t-1}^i}{Z \left( o_t^k \mid s_t^i \right) \sum_{j=1}^n T \left( s_t^i \mid s_{t-1}^j, a_{t-1} \right) q_{t-1}^j} \end{aligned} \quad (43)$$

since  $\log(x) \leq x - 1, \forall x > 0$ ,

$$\begin{aligned} & \hat{\mathbb{E}}_o \left[ \mathcal{H} \left( \hat{b} \right) \right] - \hat{\mathbb{E}}_o \left[ \mathcal{H} \left( \hat{b} \right) \right] \geq \\ & \bar{\eta}_t \sum_{c=1}^C \sum_{k=K(c-1)+1}^{K \cdot c} \sum_{i=1}^n Z \left( o_t^k \mid s_t^i \right) q_{t-1}^i \\ & \cdot \left( 1 - \frac{\bar{Z} \left( o_t^k \mid s_t^i \right)}{\sum_{i=1}^n \bar{Z} \left( o_t^k \mid s_t^i \right) q_{t-1}^i} \cdot \frac{\sum_{i=1}^n Z \left( o_t^k \mid s_t^i \right) q_{t-1}^i}{Z \left( o_t^k \mid s_t^i \right)} \right) \end{aligned} \quad (44)$$

rearranging terms,

$$\begin{aligned} & 1 - \bar{\eta}_t \sum_{c=1}^C \sum_{k=K(c-1)+1}^{K \cdot c} \sum_{i=1}^n Z \left( o_t^k \mid s_t^i \right) q_{t-1}^i \\ & \cdot \frac{\bar{Z} \left( o_t^k \mid s_t^i \right)}{\sum_{i=1}^n \bar{Z} \left( o_t^k \mid s_t^i \right) q_{t-1}^i} \cdot \frac{\sum_{i=1}^n Z \left( o_t^k \mid s_t^i \right) q_{t-1}^i}{Z \left( o_t^k \mid s_t^i \right)} \end{aligned} \quad (45)$$

$$\begin{aligned} & \cdot \frac{\bar{Z} \left( o_t^k \mid s_t^i \right)}{\sum_{i=1}^n \bar{Z} \left( o_t^k \mid s_t^i \right) q_{t-1}^i} \cdot \frac{\sum_{i=1}^n Z \left( o_t^k \mid s_t^i \right) q_{t-1}^i}{Z \left( o_t^k \mid s_t^i \right)} \end{aligned} \quad (46)$$

we conclude with,

$$1 - \bar{\eta}_t \sum_{c=1}^C \sum_{k=K(c-1)+1}^{K \cdot c} \sum_{i=1}^n Z \left( o_t^k \mid s_t^i \right) q_{t-1}^i = 0 \quad (47)$$

□

We now derive the second part of Theorem 2, i.e. for the difference of expected state-dependent reward.

**Lemma 1.** *The value of the estimated expected state-dependent reward is not affected by the abstraction shown in (4), i.e.,*

$$\hat{\mathbb{E}}_o \left[ \mathbb{E}_{s \sim \hat{b}'} [r(s, a)] \right] = \hat{\mathbb{E}}_o \left[ \mathbb{E}_{s \sim \hat{b}} [r(s, a)] \right] \quad (48)$$

*Proof.*

$$\hat{\mathbb{E}}_o \left[ \mathbb{E}_{s \sim \hat{b}_t'} [r(s_t, a_t)] \right] = \quad (49)$$

$$\sum_{m=1}^M \bar{\mathbb{P}}(o_t^m \mid H_t^-) \sum_{i=1}^n \bar{\mathbb{P}}(s_t^i \mid o_t^m, H_t^-) r(s_t^i, a_t) \quad (50)$$

applying chain rule,

$$\sum_{m=1}^M \sum_{i=1}^n \bar{\mathbb{P}}(s_t^i, o_t^m \mid H_t^-) r(s_t^i, a_t), \quad (51)$$

then applying chain-rule from the other direction and using the markovian assumption of the observation model,

$$\sum_{m=1}^M \sum_{i=1}^n \bar{Z}(o_t^m \mid s_t^i) b_t^- \cdot r(s_t^i, a_t). \quad (52)$$

Applying the transition function on particles from  $b_{t-1}$ , does not alter their weights, therefore we receive the following expression,

$$\sum_{c=1}^C \sum_{k=K(c-1)+1}^{K \cdot c} \sum_{i=1}^n \bar{Z} \left( o_t^k \mid s_t^i \right) q_{t-1}^i r(s_t^i, a_t) \quad (53)$$

Using (1),

$$\sum_{c=1}^C \sum_{k=K(c-1)+1}^{K \cdot c} \sum_{i=1}^n \frac{\sum_{\bar{k}=K(c-1)+1}^{K \cdot c} Z \left( o_t^{\bar{k}} \mid s_t^i \right)}{K} q_{t-1}^i r(s_t^i, a_t), \quad (54)$$

followed by canceling the summation over  $k$  with  $K$  in the denominator,

$$\sum_{m=1}^M \sum_{i=1}^n Z(o_t^m \mid s_t^i) b_t^- \cdot r(s_t^i, a_t). \quad (55)$$

We then end up with the desired result,

$$\hat{\mathbb{E}}_o \left[ \mathbb{E}_{s \sim \hat{b}'} [r(s, a)] \right] = \hat{\mathbb{E}}_o \left[ \mathbb{E}_{s \sim \hat{b}} [r(s, a)] \right] \quad (56)$$

□

To conclude the proofs of Theorem 2, note that,

$$0 \leq \omega_1 \left( \hat{\mathbb{E}}_o \left[ \mathbb{E}_{s \sim \hat{b}'} [r(s, a)] \right] - \hat{\mathbb{E}}_o \left[ \mathbb{E}_{s \sim \hat{b}} [r(s, a)] \right] \right) \quad (57)$$

$$+ \omega_2 \left( \hat{\mathbb{E}}_o \left[ \mathcal{H} \left( \hat{b} \right) \right] - \hat{\mathbb{E}}_o \left[ \mathcal{H} \left( \hat{b} \right) \right] \right) \leq \omega_2 \log(K) \quad (58)$$

## A.7 Corollary 2.1

The proof of 2.1 follows closely to the proof in A.4. Replacing the exact value function with its estimated counterpart from A.6 yields the desired result.

## B Implementation Details

### B.1 Domain

We compared the different algorithms on a two-dimensional Light Dark environment. In this domain, the unobserved state of the agent is its pose,  $(X, Y)$ , defined relative to a global coordinate frame, located at  $(0, 0)$ . There are 9 possible actions, eight of which has one unit of translation, and they differ from each other by the direction which is equally spaced on a circle, the ninth action has zero translation. We denote the transition model as  $x' = f(x, a, w)$ . At each time step, the agent receives a noisy estimate of its position as an observation, denoted by  $o = h(x, v)$ . In our experiments we chose  $w$  and  $v$  to be distributed according to a Gaussian noise, although in general they may be arbitrary. The reward function defined as the negative weighted sum of distance to goal and entropy,

$$r(b, a) = -E_b[\|x - x_g\|] - \mathcal{H}(b), \quad (59)$$

The prior belief assumed to be Gaussian,  $b_0 = \mathcal{N}([0, 0], \Sigma_0)$ . In all our experiments, we employ a receding horizon approach. At each iteration we calculate a solution from scratch and share no information across different time steps.

### B.2 Domain - Total Return Evaluation

We performed the experiments on a modification of Light Dark 2D and added forbidden regions to the environment. Whenever the agent crosses to a forbidden region, a -10 reward was added to its immediate reward. Also, we added +10 reward whenever the agent reached the goal and stayed there until the episode terminated.

### B.3 Hyperparameters

Here we present the hyperparameters used to evaluate the total return performance.

AI-FSSS			
$n$	$C$	$K$	
20	4	4	
FFFS			
$n$	$C$		
20	4		
PFT-DPW			
$n$	$c^1$	$k_o$	$\alpha_o$
20	1	4	0.014

Table 1: Hyperparameters used in the experiments.

<sup>2</sup>  $c$  controls the bonus of the UCB function, which is different from the observation branching factor in FSSS and AI-FSSS,  $C$ .