# Appendices

## A. Definitions

We repeat the relevant definitions in our paper.

### A1. Safe Space: For more details, see Turchetta et al. (2016).

Set of the states identified as safe up to some confidence level of $\epsilon_g$:

$$R^{\text{safe}}_{\epsilon_g}(X) = X \cup \{s \in \mathcal{S} \mid \exists s' \in X : g(s') - \epsilon_g - Ld(s, s') \geq h\}.$$

Set of states with reachability from $X$:

$$R_{\text{reach}}(X) = X \cup \{s \in \mathcal{S} \mid \exists s' \in X, a \in \mathcal{A}(s') : s = f(s', a)\}.$$

Set of states with returnability to $X$:

$$
\begin{aligned}
R_{\text{ret}}(X, \bar{X}) &= \bar{X} \cup \{s \in X \mid \exists a \in \mathcal{A} : f(s, a) \in \bar{X}\}, \\
R^n_{\text{ret}}(X, \bar{X}) &= R_{\text{ret}}(X, R^{n-1}_{\text{ret}}(X, \bar{X})), \text{with } R^1_{\text{ret}}(X, \bar{X}) = R_{\text{ret}}(X, \bar{X}), \\
\bar{R}_{\text{ret}}(X, \bar{X}) &= \lim_{n \to \infty} R^n_{\text{ret}}(X, \bar{X}).
\end{aligned}
$$

Set of safe states with reachability and returnability:

$$
\begin{aligned}
R_{\epsilon_g}(X) &= R^{\text{safe}}_{\epsilon_g}(X) \cap R_{\text{reach}}(X) \cap R_{\text{ret}}(R^{\text{safe}}_{\epsilon_g}(X), X), \\
R_{\epsilon_g}(X) &= R_{\epsilon_g}(R^{n-1}_{\epsilon_g}(X)), \text{with } R^1_{\epsilon_g}(X) = R_{\epsilon_g}(X), \\
\bar{R}_{\epsilon_g}(X) &= \lim_{n \to \infty} R^n_{\epsilon_g}(X).
\end{aligned}
$$

Pessimistic safe space:

$$
\begin{aligned}
S^-_t &= \{s \in \mathcal{S} \mid \exists s' \in \mathcal{X}^-_{t-1} : l_t(s') - L \cdot d(s, s') \geq h\}, \\
\mathcal{X}^-_t &= \{s \in S^-_t \mid s \in R_{\text{reach}}(\mathcal{X}^-_{t-1}) \cap \bar{R}_{\text{ret}}(S^-_t, \mathcal{X}^-_{t-1})\}.
\end{aligned}
$$

Optimistic safe space:

$$
\begin{aligned}
S^+_t &= \{s \in \mathcal{S} \mid \exists s' \in \mathcal{X}^+_{t-1} : u_t(s') - L \cdot d(s, s') \geq h\}, \\
\mathcal{X}^+_t &= \{s \in S^+_t \mid s \in R_{\text{reach}}(\mathcal{X}^+_{t-1}) \cap \bar{R}_{\text{ret}}(S^+_t, \mathcal{X}^+_{t-1})\}.
\end{aligned}
$$

### A2. Optimization of Cumulative Reward

For optimal policy:

$$V^*_{\mathcal{M}}(s_t) = \max_{s_{t+1} \in R_{\epsilon_g}(S_0)} \left[ r(s_{t+1}) + \gamma V^*_{\mathcal{M}}(s_{t+1}) \right].$$

For balancing exploration and exploitation (neither $\mathsf{ES}^2$ nor $\mathsf{P\text{-}ES}^2$ is used):

$$
\begin{aligned}
U_t(s) &= \mu^r_t(s) + \alpha^{1/2}_{t+1} \cdot \sigma^r_t(s), \\
J^*_{\mathcal{X}}(s_t, b^r_t, b^g_t) &= \max_{s_{t+1} \in \mathcal{X}^-_{t*}} \left[ U_t(s_{t+1}) + \gamma J^*_{\mathcal{X}}(s_{t+1}, b^r_t, b^g_t) \right].
\end{aligned}
$$

## A3. $\mathsf{ES}^2$ Algorithm

For checking whether the termination condition is satisfied:

$$V_{\mathcal{M}_y}(s_t) = \max_{s_{t+1} \in \mathcal{X}_t^+} [\, r'(s_{t+1}) + \gamma V_{\mathcal{M}_y}(s_{t+1}) \,],$$

$$\mathcal{Y}_t = \{ s' \in \mathcal{S}^+ \mid \forall s \in \mathcal{X}_t^- : s' = f(s, \pi_y^*(a \mid s)) \},$$

$$\mathcal{Y}_t \subseteq \mathcal{X}_t^-.$$

For balancing exploration and exploitation in terms of reward:

$$J_{\mathcal{Y}}^*(s_t, b_t^r, b_t^g) = \max_{s_{t+1} \in \mathcal{Y}_t} [\, U_t(s_{t+1}) + \gamma J_{\mathcal{Y}}^*(s_{t+1}, b_t^r, b_t^g) \,].$$

## A4. P-$\mathsf{ES}^2$ Algorithm

For checking whether the termination condition is satisfied:

$$V_{\mathcal{M}_z}(s_t) = \max_{s_{t+1} \in \mathcal{X}_t^+} [\, P^z \cdot \{ r'(s_{t+1}) + \gamma V_{\mathcal{M}_z}(s_{t+1}) \} \,],$$

$$\mathcal{Z}_t = \{ s' \in \mathcal{S}^+ \mid \forall s \in \mathcal{X}_t^- : s' = f(s, \pi_z^*(a \mid s)) \},$$

$$\mathcal{Z}_t \subseteq \mathcal{X}_t^-.$$

For balancing exploration and exploitation in terms of the reward:

$$J_{\mathcal{Z}}^*(s_t, b_t^r, b_t^g) = \max_{s_{t+1} \in \mathcal{Z}_t} [\, U_t(s_{t+1}) + \gamma J_{\mathcal{Z}}^*(s_{t+1}, b_t^r, b_t^g) \,].$$

# B. Preliminary Lemma

**Lemma 3.** *For two arbitrary functions $f_1(x)$ and $f_2(x)$, the following inequality holds:*

$$\max_x f_1(x) - \max_x f_2(x) \geq \min_x (f_1(x) - f_2(x)).$$

*Proof.* For two arbitrary functions $f_4(x)$ and $f_5(x)$, the following inequality holds:

$$\max_x f_4(x) + \max_x f_5(x) \geq \max_x \{ f_4(x) + f_5(x) \}.$$

Let $f_2(x) = f_4(x) + f_5(x)$ and $f_3(x) = -f_4(x)$. Then,

$$\max_x \{ -f_3(x) \} + \max_x \{ f_2(x) + f_3(x) \} \geq \max_x f_2(x),$$

$$\max_x \{ f_2(x) + f_3(x) \} - \max_x f_2(x) \geq -\max_x \{ -f_3(x) \},$$

$$\max_x \{ f_2(x) + f_3(x) \} - \max_x f_2(x) = \min_x f_3(x).$$

Finally, let $f_1(x) = f_2(x) + f_3(x)$. Then, the desired lemma is obtained. $\qquad\square$

# C. Near-optimality

**Lemma 4.** *Let $J_{\mathcal{X}}^*(s_t, b_t^r, b_t^g)$ be the value function calculated by* SNO-MDP *without the* $\mathsf{ES}^2$ *algorithm. Then, $J_{\mathcal{X}}^*(s_t, b_t^r, b_t^g)$ satisfies the following inequality:*

$$J_{\mathcal{X}}^*(s_t, b_t^r, b_t^g) \geq V^*(s_t).$$

*Proof.* Consider a state $s_t$ and beliefs $b_t^r$ and $b_t^g$. Also, let $I$ denote the following safety indicator function:

$$I(s) := \begin{cases} 1 & \text{if } s \in \bar{R}_{\epsilon_g}(S_0), \\ 0 & \text{otherwise.} \end{cases} \tag{5}$$

Then, the following chain of equations and inequalities holds:

$$J_{\mathcal{X}}^*(\boldsymbol{s}_t, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) - V^*(\boldsymbol{s}_t)$$

$$= \max_{\boldsymbol{s}_{t+1} \in \mathcal{X}_{t*}^-} \left[\, U_t(\boldsymbol{s}_{t+1}) + \gamma J_{\mathcal{X}}^*(\boldsymbol{s}_{t+1}, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) \,\right] - \max_{\boldsymbol{s}_{t+1} \in \bar{R}_{\epsilon_g}(S_0)} \left[\, r(\boldsymbol{s}_{t+1}) + \gamma V_{\mathcal{M}}^*(\boldsymbol{s}_{t+1}) \,\right]$$

$$\geq \max_{\boldsymbol{s}_{t+1} \in \bar{R}_{\epsilon_g}(S_0)} \left[\, U_t(\boldsymbol{s}_{t+1}) + \gamma J_{\mathcal{X}}^*(\boldsymbol{s}_{t+1}, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) \,\right] - \max_{\boldsymbol{s}_{t+1} \in \bar{R}_{\epsilon_g}(S_0)} \left[\, r(\boldsymbol{s}_{t+1}) + \gamma V_{\mathcal{M}}^*(\boldsymbol{s}_{t+1}) \,\right]$$

$$= \max_{a_t} \left[\, I(\boldsymbol{s}_{t+1}) \cdot \{U_t(\boldsymbol{s}_{t+1}) + \gamma J_{\mathcal{X}}^*(\boldsymbol{s}_{t+1}, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g)\} \,\right] - \max_{a_t} \left[\, I(\boldsymbol{s}_{t+1}) \cdot \{r(\boldsymbol{s}_{t+1}) + \gamma V_{\mathcal{M}}^*(\boldsymbol{s}_{t+1})\} \,\right]$$

$$\geq \min_{a_t} \left[\, I(\boldsymbol{s}_{t+1}) \cdot \{U_t(\boldsymbol{s}_{t+1}) - r(\boldsymbol{s}_{t+1})\} + \gamma I(\boldsymbol{s}_{t+1}) J_{\mathcal{X}}^*(\boldsymbol{s}_{t+1}, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) - \gamma I(\boldsymbol{s}_{t+1}) V^*(\boldsymbol{s}_{t+1}) \,\right]$$

$$= \min_{a_t} \left[\, I(\boldsymbol{s}_{t+1}) \cdot \{U_t(\boldsymbol{s}_{t+1}) - r(\boldsymbol{s}_{t+1})\} + \gamma I(\boldsymbol{s}_{t+1}) \{J_{\mathcal{X}}^*(\boldsymbol{s}_{t+1}, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) - V^*(\boldsymbol{s}_{t+1})\} \,\right].$$

The third line follows from $\mathcal{X}_{t*}^- \supseteq \bar{R}_{\epsilon_g}(S_0)$ in Theorem 1. Also, the fourth line follows from the definition of $I$, and the fifth line follows from Lemma 3. Because $\boldsymbol{s}$ is arbitrary in the above derivation, we have

$$\min_{\boldsymbol{s}_t} [\, J_{\mathcal{X}}^*(\boldsymbol{s}_t, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) - V^*(\boldsymbol{s}_t) \,] \geq \min_{\boldsymbol{s}_{t+1}} [\, I(\boldsymbol{s}_{t+1}) \{U_t(\boldsymbol{s}_{t+1}) - r(\boldsymbol{s}_{t+1})\} + \gamma I(\boldsymbol{s}_{t+1}) \{J^*(\boldsymbol{s}_{t+1}, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) - V^*(\boldsymbol{s}_{t+1})\} \,].$$

By Lemma 2, the following equation holds with probability at least $1 - \Delta^r$:

$$\min_{\boldsymbol{s}_t} [\, J_{\mathcal{X}}^*(\boldsymbol{s}_t, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) - V^*(\boldsymbol{s}_t, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) \,] \geq \gamma \cdot \min_{\boldsymbol{s}_{t+1}} [I(\boldsymbol{s}_{t+1}) \{J_{\mathcal{X}}^*(\boldsymbol{s}_{t+1}, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) - V^*(\boldsymbol{s}_{t+1})\} \,]$$

Repeatedly applying this equation proves the desired lemma. Therefore, we have

$$J_{\mathcal{X}}^*(\boldsymbol{s}_t, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) \geq V^*(\boldsymbol{s}_t)$$

with high probability. $\qquad\square$

**Lemma 5. (Generalized induced inequality)** *Let $\boldsymbol{b}^r, \boldsymbol{b}^g, r$ and $\hat{\boldsymbol{b}}^r, \hat{\boldsymbol{b}}^g, \hat{r}$ be the beliefs (over reward and safety, respectively) and reward functions (including the exploration bonus) that are identical on some set of states $\Omega$ — i.e., $\boldsymbol{b}^r = \hat{\boldsymbol{b}}^r$, $\boldsymbol{b}^g = \hat{\boldsymbol{b}}^g$, and $r = \hat{r}$ for all $\boldsymbol{s} \in \Omega$. Let $P(A_\Omega)$ be the probability that a state not in $\Omega$ is generated when starting from state $\boldsymbol{s}$ and following a policy $\pi$. If the value is bound in $[0, V_{\max}]$, then*

$$V^\pi(\boldsymbol{s}, \boldsymbol{b}^r, \boldsymbol{b}^g, r) \geq V^\pi(\boldsymbol{s}, \hat{\boldsymbol{b}}_r, \hat{\boldsymbol{b}}_g, \hat{r}) - V_{\max} P(A_\Omega),$$

*where we now make explicit the dependence of the value function on the reward.*

*Proof.* The lemma follows from Lemma 8 in Strehl & Littman (2005). $\qquad\square$

**Lemma 6.** *Assume that the reward function $r$ satisfies $\|r\|_k^2 \leq B^r$, and that the noise $n_t^r$ is $\sigma_r$-sub-Gaussian. If $\alpha_t = B^r + \sigma_r \sqrt{2(\Gamma_{t-1}^r + 1 + \log(1/\Delta^r))}$ and $C_r = 8/\log(1 + \sigma_r^{-2})$, then the following holds:*

$$\frac{1}{2} \sqrt{\frac{C_r \alpha_{t*} \Gamma_{t*}^r}{t^*}} \geq \alpha_{t*}^{1/2} \sigma_{t*}^r(\boldsymbol{s}),$$

*with probability at least $1 - \Delta^r$.*

*Proof.* The lemma follows from Lemma 4 in Chowdhury & Gopalan (2017). $\qquad\square$

# D. ES$^2$ algorithm

**Lemma 7.** *Assume that $\mathcal{Y}_t \subseteq \mathcal{X}_t^-$ holds. Suppose that we obtain the optimal policy, $\pi_y^*$ on the basis of $J_{\mathcal{Y}}^*(\boldsymbol{s}_t, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) = \max_{\boldsymbol{s}_{t+1} \in \mathcal{Y}_t} [U_t(\boldsymbol{s}_{t+1}) + \gamma J_{\mathcal{Y}}^*(\boldsymbol{s}_{t+1}, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g)]$. Then, for all $t$, the following holds:*

$$\boldsymbol{s}_t \in \mathcal{Y}_t \implies \boldsymbol{s}_{t+1} \in \mathcal{Y}_t.$$

*Proof.* When $\mathcal{Y}_t \subseteq \mathcal{X}_t^-$ holds, we have

$$\{s' \in \mathcal{S}^+ \mid \forall s \in \mathcal{Y}_t : s' = f(s, \pi_y^*(a \mid s))\} \subseteq \{s' \in \mathcal{S}^+ \mid \forall s \in \mathcal{X}_t^- : s' = f(s, \pi_y^*(a \mid s))\}$$
$$= \mathcal{Y}_t.$$

This means that the next state $s_{t+1}$ will be within $\mathcal{Y}_t$ if the agent is in $\mathcal{Y}_t$ and decides the action based on $\pi_y^*$. Therefore, we have the desired lemma. $\qquad\square$

**Lemma 8.** *Assume that $\mathcal{Y}_t \subseteq \mathcal{X}_t^-$ holds, and let $J_{\mathcal{Y}}^*(s_t, b_t^r, b_t^g)$ be the value function calculated by* SNO-MDP *with the* ES$^2$ *algorithm. Then, for all $s_t \in \mathcal{X}_t^-$, $J_{\mathcal{Y}}^*(s_t, b_t^r, b_t^g)$ satisfies the following equation:*

$$J_{\mathcal{Y}}^*(s_t, b_t^r, b_t^g) \geq V^*(s_t).$$

*Proof.* Consider a state $s_t \in \mathcal{X}_t^-$ and beliefs $b^r$ and $b^g$. Also, we define the function $I$ as in (5). Then, the following chain of the equations and inequalities holds:

$$
\begin{aligned}
&J_{\mathcal{Y}}^*(s_t, b_t^r, b_t^g) - V^*(s_t) \\
&= \max_{s_{t+1} \in \mathcal{Y}_t} \big[\, U_t(s_{t+1}) + \gamma J_{\mathcal{Y}}^*(s_{t+1}, b_t^r, b_t^g) \,\big] - \max_{a_t} \big[\, I(s_{t+1}) \cdot \{r(s_{t+1}) + \gamma V_{\mathcal{M}}^*(s_{t+1})\} \,\big] \\
&= \max_{s_{t+1} \in \mathcal{Y}_t} \big[\, U_t(s_{t+1}) + \gamma J_{\mathcal{Y}}^*(s_{t+1}, b_t^r, b_t^g) \,\big] - \max_{s_{t+1} \in \mathcal{X}_t^+} \big[\, I(s_{t+1}) \cdot \{r(s_{t+1}) + \gamma V_{\mathcal{M}}^*(s_{t+1})\} \,\big] \\
&= \max_{s_{t+1} \in \mathcal{Y}_t} \big[\, U_t(s_{t+1}) + \gamma J_{\mathcal{Y}}^*(s_{t+1}, b_t^r, b_t^g) \,\big] - \max_{s_{t+1} \in \mathcal{Y}_t} \big[\, I(s_{t+1}) \cdot \{r(s_{t+1}) + \gamma V_{\mathcal{M}}^*(s_{t+1})\} \,\big] \\
&\geq \min_{s_{t+1} \in \mathcal{Y}_t} \big[\, U_t(s_{t+1}) + \gamma J_{\mathcal{Y}}^*(s_{t+1}, b_t^r, b_t^g) - I(s_{t+1}) \cdot \{r(s_{t+1}) + \gamma V_{\mathcal{M}}^*(s_{t+1})\} \,\big] \\
&\geq \min_{s_{t+1} \in \mathcal{Y}_t} \big[\, U_t(s_{t+1}) + \gamma J_{\mathcal{Y}}^*(s_{t+1}, b_t^r, b_t^g) - \{r(s_{t+1}) + \gamma V_{\mathcal{M}}^*(s_{t+1})\} \,\big] \\
&= \min_{s_{t+1} \in \mathcal{Y}_t} \big[\, U_t(s_{t+1}) - r(s_{t+1}) + \gamma J_{\mathcal{Y}}^*(s_{t+1}, b_t^r, b_t^g) - \gamma V_{\mathcal{M}}^*(s_{t+1}) \,\big].
\end{aligned}
$$

The second and third lines follow from the definitions of $I$ and $V_{\mathcal{M}}^*$. The forth line follows from the definition of $\mathcal{Y}$ and the assumption of $\mathcal{Y}_t \subseteq \mathcal{X}_t^-$. The fifth line follows from Lemma 3.

Then, by Lemma 2, the following equation holds with probability at least $1 - \Delta^r$:

$$
\begin{aligned}
\min_{s_t \in \mathcal{X}_t^-} \big[\, J_{\mathcal{Y}}^*(s_t, b_t^r, b_t^g) - V^*(s_t)\} \,\big] &\geq \gamma \cdot \min_{s_{t+1} \in \mathcal{Y}_t} \big[\, J_{\mathcal{Y}}^*(s_{t+1}, b_t^r, b_t^g) - V_{\mathcal{M}}^*(s_{t+1}) \,\big] \\
&\geq \gamma^2 \cdot \min_{s_{t+2} \in \mathcal{Y}_t} \big[\, J_{\mathcal{Y}}^*(s_{t+2}, b_t^r, b_t^g) - V_{\mathcal{M}}^*(s_{t+2}) \,\big].
\end{aligned}
$$

The second line follows from Lemma 7. Repeatedly applying this equation proves the desired lemma. Therefore, for all $s_t \in \mathcal{X}_t^-$, we have

$$J_{\mathcal{Y}}^*(s_t, b_t^r, b_t^g) \geq V^*(s_t).$$

$\qquad\square$

## E. Main Theoretical Results

**Theorem 1.** *Assume that the safety function $g$ satisfies $\|g\|_k^2 \leq B^g$ and is $L$-Lipschitz continuous. Also, assume that $S_0 \neq \emptyset$ and $g(s) \geq h$ for all $s \in S_0$. Fix any $\epsilon_g > 0$ and $\Delta^g \in (0, 1)$. Suppose that we conduct the stage of "exploration of safety" with the noise $n_t^g$ being $\sigma_g$-sub-Gaussian, and that $\beta_t = B^g + \sigma_g \sqrt{2(\Gamma_{t-1}^g + 1 + \log(1/\Delta^g))}$ until $\max_{s \in G_t} w_t(s) < \epsilon_g$ is achieved. Finally, let $t^*$ be the smallest integer satisfying*

$$\frac{t^*}{\beta_{t^*} \Gamma_{t^*}^g} \geq \frac{C_g |\bar{R}_0(S_0)|}{\epsilon_g^2} \cdot D(\mathcal{M}),$$

*with $C_g = 8 / \log(1 + \sigma_g^{-2})$. Then, the following statements jointly hold with probability at least $1 - \Delta^g$:*

- $\forall t \geq 1$, $g(\boldsymbol{s}_t) \geq h$,

- $\exists t_0 \leq t^*$, $\bar{R}_{\epsilon_g}(S_0) \subseteq \mathcal{X}_{t_0}^- \subseteq \bar{R}_0(S_0)$.

*Proof.* This is an extension of Theorem 1 in Turchetta et al. (2016) to our settings, where $t$ represents not the number of samples but the number of actions. $\qquad\square$

**Theorem 2.** *Assume that the reward function $r$ satisfies $\|r\|_k^2 \leq B^r$, and that the noise is $\sigma_r$-sub-Gaussian. Let $\pi_t$ denote the policy followed by* SNO-MDP *at time $t$, and let $\boldsymbol{s}_t$ and $\boldsymbol{b}_t^r, \boldsymbol{b}_t^g$ be the corresponding state and beliefs, respectively. Let $t^*$ be the smallest integer satisfying $\frac{t^*}{\beta_{t^*} \Gamma_{t^*}^g} \geq \frac{C_g |\bar{R}_0(S_0)|}{\epsilon_g^2} D(\mathcal{M})$, and fix any $\Delta^r \in (0, 1)$. Finally, set $\alpha_t = B^r + \sigma_r \sqrt{2(\Gamma_{t-1}^r + 1 + \log(1/\Delta^r))}$ and*

$$\epsilon_V^* = V_{\max} \cdot (\Delta^g + \Sigma_{t^*}^r / R_{\max}),$$

*with $\Sigma_{t^*}^r = \frac{1}{2}\sqrt{\frac{C_r \alpha_{t^*} \Gamma_{t^*}^r}{t^*}}$. Then, with high probability,*

$$V^{\pi_t}(\boldsymbol{s}_t, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) \geq V^*(\boldsymbol{s}_t) - \epsilon_V^*$$

*— i.e., the algorithm is $\epsilon_V^*$-close to the optimal policy — for all but $t^*$ time steps, while guaranteeing safety with probability at least $1 - \Delta^g$.*

*Proof.* Define $\tilde{r}$ as the reward function (including the exploration bonus) that is used by SNO-MDP. Let $\hat{r}$ be a reward function equal to $r$ on $\Omega$ and equal to $\tilde{r}$ elsewhere. Furthermore, let $\tilde{\pi}$ be the policy followed by SNO-MDP at time $t$, that is, the policy calculated on the basis of the current beliefs, (i.e., $\boldsymbol{b}_t^r$ and $\boldsymbol{b}_t^g$) and the reward $\tilde{r}$. Finally, let $A_\Omega$ be the event in which $\tilde{\pi}$ escapes from $\Omega$. Then,

$$V^{\pi_t}(r, \boldsymbol{s}_t, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) \geq V^{\tilde{\pi}}(\hat{r}, \boldsymbol{s}_t, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) - V_{\max} P(A_\Omega)$$

by Lemma 5. In addition, note that, for all $t \geq t^*$, because $\hat{r}$ and $\tilde{r}$ differ by at most $\alpha_{t^*}^{1/2} \sigma_{t^*}^r$ at each state,

$$|V^{\tilde{\pi}}(\hat{r}, \boldsymbol{s}_t, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g) - V^{\tilde{\pi}}(\tilde{r}, \boldsymbol{s}_t, \boldsymbol{b}_t^r, \boldsymbol{b}_t^g)| \leq \frac{1}{1-\gamma} \cdot \alpha_{t^*}^{1/2} \sigma_{t^*}^r(\boldsymbol{s})$$
$$\leq V_{\max}/R_{\max} \cdot \Sigma_{t^*}^r. \tag{6}$$

For the above inequality, we used Lemma 6. Here, consider the case of $\Omega = \mathcal{X}_{t^*}^-$. Once the safe region is fully explored, $P(A_\Omega) \leq \Delta^g$ holds after $t^*$ time steps. Then, the following chain of equations and inequalities holds:

$$V^{\pi_t}(R, \boldsymbol{s}, \boldsymbol{b}) \geq V^{\tilde{\pi}}(\hat{R}, \boldsymbol{s}, \boldsymbol{b}) - V_{\max} \cdot P(A_\Omega)$$
$$= V^{\tilde{\pi}}(\hat{R}, \boldsymbol{s}, \boldsymbol{b}) - V_{\max} \cdot P(A_{\mathcal{X}^-})$$
$$\geq V^{\tilde{\pi}}(\hat{R}, \boldsymbol{s}, \boldsymbol{b}) - V_{\max} \cdot \Delta^g$$
$$\geq V^{\tilde{\pi}}(\tilde{R}, \boldsymbol{s}, \boldsymbol{b}) - V_{\max} \cdot (\Delta^g + \Sigma_{t^*}^r / R_{\max})$$
$$= J_{\mathcal{X}}^*(\tilde{R}, \boldsymbol{s}, \boldsymbol{b}) - V_{\max} \cdot (\Delta^g + \Sigma_{t^*}^r / R_{\max})$$
$$\geq V^*(R, \boldsymbol{s}) - V_{\max} \cdot (\Delta^g + \Sigma_{t^*}^r / R_{\max}).$$

In this derivation, the second line follows from the assumption of $\Omega = \mathcal{X}^-$, the third line follows from $P(A_{\mathcal{X}^-}) \leq \Delta^g$, the fourth line follows from (6), the fifth line follows from the fact that $\tilde{\pi}$ is precisely the optimal policy for $\tilde{R}$ and $\boldsymbol{b}$, and the final line follows from Lemma 4. $\qquad\square$

**Theorem 3.** *Assume that the reward function $r$ satisfies $\|r\|_k^2 \leq B^r$, and that the noise is $\sigma_r$-sub-Gaussian. Let $\pi_t$ denote the policy followed by* SNO-MDP *with the the* ES$^2$ *algorithm at time $t$, and let $\boldsymbol{s}_t$ and $\boldsymbol{b}_t^r, \boldsymbol{b}_t^g$ be the corresponding state and beliefs, respectively. Let $\tilde{t}$ be the smallest integer for which (4) holds, and fix any $\Delta^r \in (0, 1)$. Finally, set $\alpha_t = B^r + \sigma_r \sqrt{2(\Gamma_{t-1}^r + 1 + \log(1/\Delta^r))}$ and*

$$\tilde{\epsilon}_V = V_{\max} \cdot (\Delta^g + \Sigma_{\tilde{t}}^r / R_{\max}),$$

*with $\Sigma_{\tilde{t}}^r = \frac{1}{2}\sqrt{\frac{C_r \alpha_{\tilde{t}} \Gamma_{\tilde{t}}^r}{\tilde{t}}}$. Then, with high probability,*

$$V^{\pi_t}(s_t, b_t^r, b_t^g) \geq V^*(s_t) - \tilde{\epsilon}_V$$

*— i.e., the algorithm is $\tilde{\epsilon}_V$-close to the optimal policy — for all but $\tilde{t}$ time steps while guaranteeing safety with probability at least $1 - \Delta^g$.*

*Proof.* The proof of Theorem 3 is analogous to that of Theorem 2. Define $\tilde{r}$ as the reward function (including the exploration bonus) that is used by SNO-MDP. Let $\hat{r}$ be a reward function equal to $r$ on $\mathcal{Y}$ and equal to $\tilde{r}$ elsewhere. Furthermore, let $\tilde{\pi}$ be the policy followed by SNO-MDP with the ES$^2$ algorithm at time $t$, that is, the policy calculated on the basis of the current beliefs, (i.e., $b_t^r$ and $b_t^g$) and the reward $\tilde{r}$. Finally, let $A_{\mathcal{Y}}$ be the event in which $\tilde{\pi}$ escapes from $\mathcal{Y}$. Then,

$$V^{\pi_t}(r, s_t, b_t^r, b_t^g) \geq V^{\tilde{\pi}}(\hat{r}, s_t, b_t^r, b_t^g) - V_{\max} P(A_{\mathcal{Y}})$$

by Lemma 5. In addition, note that, for all $t \geq \tilde{t}$, because $\hat{r}$ and $\tilde{r}$ differ by at most $\alpha_{\tilde{t}}^{1/2} \sigma_{\tilde{t}}^r$ at each state,

$$|V^{\tilde{\pi}}(\hat{r}, s_t, b_t^r, b_t^g) - V^{\tilde{\pi}}(\tilde{r}, s_t, b_t^r, b_t^g)| \leq \frac{1}{1-\gamma} \cdot \alpha_{\tilde{t}}^{1/2} \sigma_t^r(s)$$

$$\leq V_{\max}/R_{\max} \cdot \Sigma_{\tilde{t}}^r. \tag{7}$$

For the above inequalities, we used Lemma 6. Then, the following chain of equations and inequalities holds:

$$
\begin{aligned}
V^{\pi_t}(R, s, b) &= V^{\tilde{\pi}}(\hat{R}, s, b) - V_{\max} \cdot P(A_{\mathcal{Y}}) \\
&\geq V^{\tilde{\pi}}(\hat{R}, s, b) - V_{\max} \cdot \Delta^g \\
&\geq V^{\tilde{\pi}}(\tilde{R}, s, b) - V_{\max} \cdot (\Delta^g + \Sigma_{\tilde{t}}^r/R_{\max}) \\
&= J_{\mathcal{Y}}^*(\tilde{R}, s, b) - V_{\max} \cdot (\Delta^g + \Sigma_{\tilde{t}}^r/R_{\max}) \\
&\geq V^*(R, s) - V_{\max} \cdot (\Delta^g + \Sigma_{\tilde{t}}^r/R_{\max}).
\end{aligned}
$$

In this derivation, the second line follows from $P(A_{\mathcal{Y}}) \leq \Delta^g$, the third line follows from (7), the fourth line follows from the fact that $\tilde{\pi}$ is precisely the optimal policy for $\tilde{R}$ and $b$, and the final line follows from Lemma 8. $\qquad \square$