

SACBP: Belief Space Planning for Continuous-Time Dynamical Systems via Stochastic Sequential Action Control

Journal Title
XX(X):1–24
©The Author(s) 2019
Reprints and permission:
sagepub.co.uk/journalsPermissions.nav
DOI: 10.1177/ToBeAssigned
www.sagepub.com/

SAGE

Haruki Nishimura¹ and Mac Schwager¹

Abstract

We propose a novel **belief space planning technique for continuous dynamics** by viewing the belief system as a **hybrid dynamical system with time-driven switching**. Our approach is based on the perturbation theory of differential equations and extends Sequential Action Control [Ansari and Murphey \(2016\)](#) to stochastic belief dynamics. The resulting algorithm, which we name SACBP, does not require discretization of spaces or time and synthesizes control signals in near real-time. SACBP is an anytime algorithm that can handle general parametric Bayesian filters under certain assumptions. We demonstrate the effectiveness of our approach in an active sensing scenario and a model-based Bayesian reinforcement learning problem. In these challenging problems, we show that the algorithm significantly outperforms other existing solution techniques including approximate dynamic programming and local trajectory optimization.

Keywords

Belief Space Planning, Active Sensing, Mobile Robots, Optimization and Optimal Control, Probabilistic Reasoning, Vision and Sensor-based Control

1 Introduction

Planning under uncertainty still remains as a challenge for robotic systems. Various types of uncertainty, including unmodeled dynamics, stochastic disturbances, and imperfect sensing, significantly complicate problems that are otherwise easy. For example, suppose that a robot needs to manipulate an object from some initial state to a desired goal. If the mass properties of the object are not known beforehand, the robot needs to simultaneously estimate these parameters and perform control, while taking into account the effects of their uncertainty; the exploration and exploitation trade-off needs to be resolved [Slade et al. \(2017\)](#). On the other hand, uncertainty is quite fundamental in motivating some problems. For instance, a noisy sensor may encourage the robot to carefully plan a trajectory so the observations taken along it are sufficiently informative. This type of problem concerns pure information gathering and is often referred to as active sensing [Mihaylova et al. \(2002\)](#), active perception [Bajcsy \(1988\)](#), or informative motion planning [Hollinger and Sukhatme \(2014\)](#).

A principled approach to address all those problems is to form plans in the belief space, where the planner chooses sequential control inputs based on the evolution of the belief state. This approach enables the robot to appropriately execute controls under stochasticity and partial observability since they are both incorporated into the belief state. Belief space planning is also well suited for generating information gathering actions [Platt et al. \(2010\)](#).

This paper proposes a novel online belief space planning algorithm. It does not require discretization of the state space or the action space, and can directly handle continuous-time system dynamics. The algorithm optimizes the expected value of the first-order cost reduction with

respect to a nominal control policy at every re-planning time, proceeding in a receding horizon fashion. We are inspired by the Sequential Action Control (SAC) algorithm recently proposed in [Ansari and Murphey \(2016\)](#) for model-based deterministic optimal control problems. SAC is an online method to synthesize control signals in real time for challenging (but deterministic) physical systems such as a cart pendulum and a spring-loaded inverted pendulum. Based on the concept of SAC, this paper develops an algorithmic framework to control stochastic belief systems whose dynamics are governed by parametric Bayesian filters.

1.1 Related Work in Belief Space Planning

Greedy Strategies Belief space planning is known to be challenging for a couple of reasons. First, the belief state is continuous and can be high-dimensional even if the underlying state space is small or discrete. Second, the dynamics that govern the belief state transitions are stochastic due to unknown future observations. Greedy approaches alleviate the complexity by ignoring long-term effects and solve single-shot decision making problems. Despite their suboptimality for long-term planning, these methods are often employed to find computationally tractable solutions and achieve reasonable performance in different problems [Bourgault et al. \(2002\)](#); [Seekircher et al.](#)

¹ Stanford University, USA

Corresponding author:

Haruki Nishimura, Department of Aeronautics and Astronautics, Stanford University, Stanford, CA 94305, USA.

Email: hnishimura@stanford.edu

(2011); Schwager et al. (2017), especially in the active sensing domain.

Trajectory Optimization Methods In contrast to the greedy approaches, trajectory optimization methods take into account multiple timesteps at once and find non-myopic solutions. In doing so, it is often assumed that the maximum likelihood observation (MLO) will always occur at the planning phase Platt et al. (2010); Erez and Smart (2010); Patil et al. (2014). This heuristic assumption results in a deterministic optimal control problem, whereby various nonlinear trajectory optimization algorithms are applicable. However, ignoring the effects of stochastic future observations can degrade the performance van den Berg et al. (2012). Other methods van den Berg et al. (2012); Rafieisakhaei et al. (2017) that do not rely on the MLO assumption are advantageous in that regard. In particular, belief iLQG van den Berg et al. (2012) performs iterative local optimization in the Gaussian belief space by quadratically approximating the value function and linearizing the dynamics to obtain a time-variant linear feedback policy. However, this method as well as many other solution techniques in this category result in multiple iterations of intensive computation and can require a significant amount of time until convergence.

Belief MDP and POMDP Approaches Belief space planning can be modeled as a Markov decision process (MDP) in the belief space, given that the belief state transition is Markovian. If in addition the reward (or cost) is defined as an explicit function of the state and the control, the problem is equivalent to a partially observable Markov decision process (POMDP) Kaelbling et al. (1998). A key challenge in POMDPs and belief MDPs has been to address problems with large state spaces. This is particularly important in belief MDPs since the state space for a belief MDP is a continuous belief space. To handle continuous spaces, Couëtoux et al. (2011) introduce double progressive widening (DPW) for Monte Carlo Tree Search (MCTS) Browne et al. (2012). In Slade et al. (2017), this MCTS-DPW algorithm is run in the belief space to solve the object manipulation problem mentioned above. We have also presented a motion-based communication algorithm in our prior work, which uses MCTS-DPW for active intent inference with monocular vision Nishimura and Schwager (2018a).

While MCTS-DPW as well as other general purpose POMDP methods Somani et al. (2013); Sunberg and Kochenderfer (2017) are capable of handling continuous state spaces, their algorithmic concepts are rooted in dynamic programming and tree search, requiring a sufficient amount of exploration in the tree. The tree search technique also implicitly assumes discrete-time transition models. In fact, most prior works discussed above are intended for discrete-time systems. There still remains a need for an efficient and high-performance belief space planning algorithm that is capable of directly handling systems with inherently continuous-space, continuous-time dynamics, such as maneuvering micro-aerial vehicles, or autonomous cars at freeway speeds.

1.2 Contributions

Our approach presented in this paper is significantly different than the previous approaches discussed above. We view the stochastic belief dynamics as a hybrid system with time-driven switching Heemels et al. (2009), where the controls are applied in continuous time and the observations are made in discrete time. A discrete-time observation creates a jump discontinuity in the belief state trajectory due to a sudden Bayesian update of the belief state. This view of belief space planning yields a continuous-time optimal control problem of a high-dimensional hybrid system. We then propose a model-based control algorithm to efficiently compute the control signals in a receding-horizon fashion. The algorithm is based on Sequential Action Control (SAC) Ansari and Murphey (2016). SAC in its original form is a deterministic, model-based hybrid control algorithm, which “perturbs” a nominal control trajectory in a structured way so that the cost functional is optimally reduced up to the first order. The key to this approach is the use of the perturbation theory of differential equations that is often discussed in the mode scheduling literature Egerstedt et al. (2006); Wardi and Egerstedt (2012). As a result, SAC derives the optimal perturbation in closed form and synthesizes control signals at a high frequency to achieve a significant improvement over other optimal control methods based on local trajectory optimization.

We apply the perturbation theory to parametric Bayesian filters and derive the optimal control perturbation using the framework of SAC. To account for stochasticity, we also extend the original algorithm by incorporating Monte Carlo sampling of nominal belief trajectories. Our key contribution is the resulting continuous belief space planning algorithm, which we name SACBP. The algorithm has the following desirable properties:

1. SACBP optimizes the expected value of the first-order reduction of the cost functional with respect to some nominal control in near real-time.
2. SACBP does not require discretization of the state space, the observation space, or the control space. It also does not require discretization of time other than for numerical integration purposes.
3. General nonlinear parametric Bayesian filters can be used for state estimation as long as the system is control-affine and the control cost is quadratic.
4. Stochasticity in the future observations are fully considered.
5. SACBP is an anytime algorithm. Furthermore, the Monte Carlo sampling part of the algorithm is naturally parallelizable.
6. Even though SACBP is inherently suboptimal for the original stochastic optimal control problem, empirical results suggest that it is highly sample-efficient and outperforms other approaches when near real-time performance is required.

Although there exists prior work Mavrommati et al. (2018) that uses SAC for active sensing, its problem

formulation relies on the ergodic control framework, which is significantly different from the belief space planning framework we propose here. We show that our SACBP outperforms projection-based ergodic trajectory optimization, MCTS-DPW, and a greedy method on a multi-target tracking example. We also show that SACBP outperforms belief iLQG and MCTS-DPW on a manipulation scenario.

This paper is an extension of the theory and results previously presented by the authors in Nishimura and Schwager (2018b). Compared to the conference version, we provide a more detailed derivation of the algorithm (Section 2) as well as a thorough mathematical analysis of the control perturbation for stochastic hybrid systems (Section 3 and Appendix A). This analysis leads to a guarantee for SACBP that, with an appropriate choice of the perturbation duration, the algorithm is expected to perform no worse than the nominal policy. Since the nominal policy can be arbitrary, one could even provide an approximately optimal discrete POMDP policy derived offline as a nominal policy to “warm-start” the planning.

In the next section we derive relevant equations and present the SACBP algorithm along with a discussion on computational complexity. Section 3 provides the key results of the mathematical analysis. Section 4 summarizes the simulation results. Conclusions and future work are presented in Section 5.

2 SACBP Algorithm

We first consider the case where some components of the state are fully observable. We begin with this mixed observability case as it is simpler to explain, yet still practically relevant. For example, this is a common assumption in various active sensing problems Schwager et al. (2017); Le Ny and Pappas (2009); Popovi et al. (2017) where the state of the robot is perfectly known, but some external variable of interest (e.g. a target’s location) is stochastic. In addition, deterministic state transitions are often assumed for the robot. Therefore, in Section 2.1 we derive the SACBP control update formulae for this case. The general belief space planning where none of the state is fully observable or deterministically controlled is discussed in Section 2.2. An extension to use a closed-loop policy as the nominal control is presented in Section 2.3. The computation time complexity is discussed in Section 2.4.

2.1 Problems with Mixed Observability

Suppose that a robot can fully observe and deterministically control some state $p(t) \in \mathbb{R}^{n_p}$. Other states are not known to the robot and are estimated with the belief vector $b(t) \in \mathbb{R}^{n_b}$. This belief vector characterizes a probability distribution that the robot uses for state estimation. If the belief is Gaussian, for example, the covariance matrix can be vectorized column-wise and stacked all together with the mean to form the belief vector. We define the augmented state as $s \triangleq (p^T, b^T)^T \in \mathbb{R}^{n_s}$.

2.1.1 Dynamics Model The physical state p is described by the following ODE:

$$\dot{p}(t) = f(p(t), u(t)), \quad (1)$$

where $u(t) \in \mathbb{R}^m$ is the control signal. On the other hand, suppose that the belief state only changes in discrete time upon arrival of a new observation from the sensors. This is the usual case for discrete-time Bayesian filtering. We will discuss the more general continuous-discrete time filtering in Section 2.2. Let t_k be the time when the k -th observation becomes available to the robot. The belief state transition is given by

$$\begin{cases} b(t_k) = g(p(t_k^-), b(t_k^-), y_k) \\ b(t) = b(t_k) \end{cases} \quad \forall t \in [t_k, t_{k+1}), \quad (2)$$

where t_k^- is infinitesimally smaller than t_k . Nonlinear function g corresponds to a discrete-time, parametric Bayesian filter (e.g., Kalman filter, extended Kalman filter, discrete Bayesian filter, etc.) that forward-propagates the belief for prediction, takes the new observation $y_k \in \mathbb{R}^q$, and returns the updated belief state. The concrete choice of the filter depends on the instance of the problem.

Equations (1) and (2) constitute a hybrid system with time-driven switching Heemels et al. (2009). This hybrid system representation is practical since it captures the fact that the observation updates occur less frequently than the control actuation in general, due to expensive information processing of sensor readings. Furthermore, with this representation one can naturally handle agile systems as they are without coarse discretization in time.

Given the initial state $s_0 \triangleq (p(t_0)^T, b(t_0)^T)^T$ and a control trajectory from t_0 to t_f denoted as u , the system evolves stochastically according to the hybrid dynamics equations. The stochasticity is due to a sequence of stochastic future observations that will be taken by t_f . In this paper we assume that the observation interval $t_{k+1} - t_k \triangleq \Delta t_o$ is fixed, and the control signals are recomputed when a new observation is incorporated in the belief, although this assumption is not critical.

2.1.2 Perturbed Dynamics The control synthesis of SACBP begins with a given nominal control u . Suppose that the nominal control is applied to the system and a sequence of T observations (y_1, \dots, y_T) is obtained. Conditioned on the observation sequence, the augmented state evolves deterministically. Let $s = (p^T, b^T)^T$ be the nominal trajectory of the augmented state induced by (y_1, \dots, y_T) .

Now let us consider perturbing the nominal trajectory at a fixed time $\tau < t_1$ for a short duration ϵ . The perturbed control u^ϵ is defined as

$$u^\epsilon(t) \triangleq \begin{cases} v & \text{if } t \in (\tau - \epsilon, \tau] \\ u(t) & \text{otherwise.} \end{cases} \quad (3)$$

Therefore, the control perturbation is determined by the nominal control u , the tuple (τ, v) , and ϵ . Given (τ, v) , the resulting perturbed system trajectory can be written as

$$\begin{cases} p^\epsilon(t) \triangleq p(t) + \epsilon \Psi_p(t) + o(\epsilon) \\ b^\epsilon(t) \triangleq b(t) + \epsilon \Psi_b(t) + o(\epsilon), \end{cases} \quad (4)$$

where $\Psi_p(t)$ and $\Psi_b(t)$ are the state variations that are linear in the perturbation duration ϵ :

$$\Psi_p(t) = \left. \frac{\partial_+}{\partial \epsilon} p^\epsilon(t) \right|_{\epsilon=0} \triangleq \lim_{\epsilon \rightarrow 0^+} \frac{p^\epsilon(t) - p(t)}{\epsilon} \quad (5)$$

$$\Psi_b(t) = \left. \frac{\partial_+}{\partial \epsilon} b^\epsilon(t) \right|_{\epsilon=0} \triangleq \lim_{\epsilon \rightarrow 0^+} \frac{b^\epsilon(t) - b(t)}{\epsilon}. \quad (6)$$

The notation $\frac{\partial_+}{\partial \epsilon}$ represents the right derivative with respect to ϵ . The state variations at perturbation time τ satisfy

$$\begin{cases} \Psi_p(\tau) = f(p(\tau), v) - f(p(\tau), u(\tau)) \\ \Psi_b(\tau) = 0, \end{cases} \quad (7)$$

assuming that τ does not exactly correspond to any of the switching times t_k . For $t \geq \tau$, the physical state variation Ψ_p evolves according to the following first-order ODE:

$$\dot{\Psi}_p(t) = \frac{d}{dt} \left(\left. \frac{\partial_+}{\partial \epsilon} p^\epsilon(t) \right|_{\epsilon=0} \right) \quad (8)$$

$$= \left. \frac{\partial_+}{\partial \epsilon} \dot{p}^\epsilon(t) \right|_{\epsilon=0} \quad (9)$$

$$= \left. \frac{\partial_+}{\partial \epsilon} f(p^\epsilon(t), u(t)) \right|_{\epsilon=0} \quad (10)$$

$$= \frac{\partial}{\partial p} f(p(t), u(t)) \Psi_p(t), \quad (11)$$

where the chain rule of differentiation and $p^\epsilon(t)|_{\epsilon=0} = p(t)$ are used in (11). For a more rigorous analysis, see Appendix A. The dynamics of the belief state variation Ψ_b in the continuous region $t \in [t_k, t_{k+1})$ satisfy $\dot{\Psi}_b(t) = 0$ since the belief vector $b(t)$ is constant according to (2). However, across the jumps the belief state variation Ψ_b changes discontinuously and satisfies

$$\Psi_b(t_k) = \left. \frac{\partial_+}{\partial \epsilon} b^\epsilon(t_k) \right|_{\epsilon=0} \quad (12)$$

$$= \left. \frac{\partial_+}{\partial \epsilon} g(p^\epsilon(t_k^-), b^\epsilon(t_k^-), y_k) \right|_{\epsilon=0} \quad (13)$$

$$\begin{aligned} &= \frac{\partial}{\partial p} g(p(t_k^-), b(t_k^-), y_k) \Psi_p(t_k^-) \\ &+ \frac{\partial}{\partial b} g(p(t_k^-), b(t_k^-), y_k) \Psi_b(t_k^-). \end{aligned} \quad (14)$$

2.1.3 Perturbed Cost Functional Let us consider a total cost of the form

$$\int_{t_0}^{t_f} c(p(t), b(t), u(t)) dt + h(p(t_f), b(t_f)), \quad (15)$$

where c is the running cost and h is the terminal cost. Following the discussion above on the perturbed dynamics, let J denote the total cost of the nominal trajectory conditioned on the given observation sequence (y_1, \dots, y_T) . Under the fixed (τ, v) , we can represent the perturbed cost J^ϵ in terms of J as

$$J^\epsilon \triangleq J + \epsilon \nu(t_f) + o(\epsilon), \quad (16)$$

where $\nu(t_f) \triangleq \frac{\partial_+}{\partial \epsilon} J^\epsilon|_{\epsilon=0}$ is the variation of the total cost with respect to the perturbation. For further analysis it is

convenient to express the running cost in the Mayer form [Liberzon \(2011\)](#). Let $\hat{s}(t)$ be a new state variable defined by $\dot{\hat{s}}(t) = c(p(t), b(t), u(t))$ and $\hat{s}(t_0) = 0$. Then the total cost is a function of the appended augmented state $\bar{s} \triangleq (\hat{s}, s^T)^T \in \mathbb{R}^{1+n_s}$ at time t_f , which is given by

$$J = \hat{s}(t_f) + h(s(t_f)). \quad (17)$$

Using this form of the total cost J , the perturbed cost (16) becomes

$$J^\epsilon = J + \epsilon \left[\begin{array}{c} 1 \\ \frac{\partial}{\partial p} h(p(t_f), b(t_f)) \\ \frac{\partial}{\partial b} h(p(t_f), b(t_f)) \end{array} \right]^T \bar{\Psi}(t_f) + o(\epsilon), \quad (18)$$

where $\bar{\Psi}(t_f) \triangleq (\hat{\Psi}(t_f), \Psi_p(t_f)^T, \Psi_b(t_f)^T)^T$. Note that the dot product in (18) corresponds to $\nu(t_f)$ in (16). The variation $\hat{\Psi}$ of the appended augmented state follows the variational equation for $t \geq \tau$:

$$\dot{\hat{\Psi}}(t) = \frac{d}{dt} \left(\left. \frac{\partial_+}{\partial \epsilon} \hat{s}^\epsilon(t) \right|_{\epsilon=0} \right) \quad (19)$$

$$= \left. \frac{\partial_+}{\partial \epsilon} \dot{\hat{s}}^\epsilon(t) \right|_{\epsilon=0} \quad (20)$$

$$= \left. \frac{\partial_+}{\partial \epsilon} c(p^\epsilon(t), b^\epsilon(t), u(t)) \right|_{\epsilon=0} \quad (21)$$

$$\begin{aligned} &= \frac{\partial}{\partial p} c(p(t), b(t), u(t))^T \Psi_p(t) \\ &+ \frac{\partial}{\partial b} c(p(t), b(t), u(t))^T \Psi_b(t) \end{aligned} \quad (22)$$

where the initial condition is given by $\hat{\Psi}(\tau) = c(p(\tau), b(\tau), v) - c(p(\tau), b(\tau), u(\tau))$.

The perturbed cost equation (18), especially the dot product expressing $\nu(t_f)$, is consequential; it tells us how the total cost changes due to the perturbation applied at some time τ , up to the first order with respect to the perturbation duration ϵ . At this point, one could compute the value of $\nu(t_f)$ for a control perturbation with a specific value of (τ, v) by simulating the nominal dynamics and integrating the variational equations (11)(14)(22) from τ up to t_f .

2.1.4 Adjoint Equations Unfortunately, this forward integration of $\nu(t_f)$ is not so useful by itself since we are interested in finding the value of (τ, v) that achieves the smallest possible $\nu(t_f)$, if it exists; it would be computationally intensive to apply control perturbation at different application times τ with different values of v and re-simulate the state variation $\bar{\Psi}$. To avoid this computationally expensive search, [Ansari and Murphey \(2016\)](#) has introduced the adjoint system $\bar{\rho}$ with which the dot product remains invariant:

$$\frac{d}{dt} (\bar{\rho}(t)^T \bar{\Psi}(t)) = 0 \quad \forall t \in [t_0, t_f]. \quad (23)$$

If we let

$$\bar{\rho}(t_f) \triangleq \left(1, \frac{\partial}{\partial p} h(p(t_f), b(t_f))^T, \frac{\partial}{\partial b} h(p(t_f), b(t_f))^T \right)^T \quad (24)$$

so that its dot product with $\bar{\Psi}(t_f)$ equals $\nu(t_f)$ as in (18), the time invariance gives

$$\nu(t_f) = \bar{\rho}(t_f)^T \bar{\Psi}(t_f) \quad (25)$$

$$= \bar{\rho}(\tau)^T \bar{\Psi}(\tau) \quad (26)$$

$$= \bar{\rho}(\tau)^T \begin{bmatrix} c(p(\tau), b(\tau), v) - c(p(\tau), b(\tau), u(\tau)) \\ f(p(\tau), v) - f(p(\tau), u(\tau)) \\ 0 \end{bmatrix}. \quad (27)$$

Therefore, we can compute the first-order cost change $\nu(t_f)$ for different values of τ once the adjoint trajectory is derived. For $t \in [t_k, t_{k+1})$ the time derivative of $\bar{\Psi}$ exists, and the invariance property leads to the following equation:

$$\dot{\bar{\rho}}(t)^T \bar{\Psi}(t) + \bar{\rho}(t)^T \dot{\bar{\Psi}}(t) = 0. \quad (28)$$

It can be verified that the following system satisfies (28) with $\bar{\rho}(t) = (\hat{\rho}(t), \rho_p(t)^T, \rho_b(t)^T)^T$.

$$\begin{cases} \dot{\hat{\rho}}(t) = 0 \\ \dot{\rho}_p(t) = -\frac{\partial}{\partial p} c(p(t), b(t), u(t)) - \frac{\partial}{\partial p} f(p(t), u(t))^T \rho_p(t) \\ \dot{\rho}_b(t) = -\frac{\partial}{\partial b} c(p(t), b(t), u(t)). \end{cases} \quad (29)$$

Analogously, across discrete jumps we can still enforce the invariance by setting $\bar{\rho}(t_k)^T \bar{\Psi}(t_k) = \bar{\rho}(t_k^-)^T \bar{\Psi}(t_k^-)$, which holds for the following adjoint equations:

$$\begin{cases} \rho^0(t_k^-) = \rho^0(t_k) \\ \rho_p(t_k^-) = \rho_p(t_k) + \frac{\partial}{\partial p} g(p(t_k^-), b(t_k^-), y_k)^T \rho_b(t_k) \\ \rho_b(t_k^-) = \frac{\partial}{\partial b} g(p(t_k^-), b(t_k^-), y_k)^T \rho_b(t_k). \end{cases} \quad (30)$$

Note that the adjoint system integrates backward in time as it has the boundary condition (24) defined at t_f . More importantly, the adjoint dynamics (29)(30) only depend on the nominal trajectory of the system (p, b) and the observation sequence (y_1, \dots, y_T) . The cost variation term $\nu(t_f)$ is finally given by

$$\nu(t_f) = c(p(\tau), b(\tau), v) - c(p(\tau), b(\tau), u(\tau)) + \rho_p(\tau)^T \{f(p(\tau), v) - f(p(\tau), u(\tau))\}. \quad (31)$$

2.1.5 Control Optimization In order to efficiently optimize (31) with respect to (τ, v) , we assume that the control cost is additive quadratic $\frac{1}{2}v^T C_u v$ and the dynamics model $f(p, u)$ is control-affine with linear term $H(p)u$. Although the control-affine assumption may appear restrictive, many physical systems possess this property in engineering practice. As a result of these assumptions, (31) becomes

$$\nu(t_f) = \frac{1}{2}v^T C_u v + \rho_p(\tau)^T H(p(\tau))(v - u(\tau)) - \frac{1}{2}u(\tau)^T C_u u(\tau). \quad (32)$$

So far we have treated the observation sequence (y_1, \dots, y_T) as given and fixed. However, in practice it is a random process that we have to take into account. Fortunately, our control optimization is all based on the

nominal control u , with which we can both simulate the augmented dynamics and sample the observations. To see this, let us consider the observations as a sequence of random vectors (Y_1, \dots, Y_T) and rewrite $\nu(t_f)$ in (32) as $\nu(t_f, Y_1, \dots, Y_T)$ to clarify the dependence on it. The expected value of the first order cost variation is given by

$$\mathbb{E}[\nu(t_f)] = \int \nu(t_f, Y_1, \dots, Y_T) d\mathbb{P}, \quad (33)$$

where \mathbb{P} is the probability measure associated with the these random vectors. Although we do not know the specific values of \mathbb{P} , we have the generative model; we can simulate the augmented state trajectory using the nominal control u and sample the stochastic observations from the belief states along the trajectory.

Using the linearity of expectation for (32), we have

$$\begin{aligned} \mathbb{E}[\nu(t_f)] &= \frac{1}{2}v^T C_u v + \mathbb{E}[\rho_p(\tau)]^T H(p(\tau))(v - u(\tau)) \\ &\quad - \frac{1}{2}u(\tau)^T C_u u(\tau). \end{aligned} \quad (34)$$

Notice that only the adjoint trajectory is stochastic. We can employ Monte Carlo sampling to sample a sufficient number of observation sequences to approximate the expected adjoint trajectory. Now (34) becomes a convex quadratic in v for a positive definite C_u . Assuming that C_u is also diagonal, analytical solutions are available to the following convex optimization problem.

$$\begin{aligned} &\underset{v}{\text{minimize}} && \mathbb{E}[\nu(t_f)] \\ &\text{subject to} && a \preceq v \preceq b \end{aligned} \quad (35)$$

This optimization is solved for different values of $\tau \in (t_0 + t_{calc} + \epsilon, t_0 + \Delta t_o)$, where t_{calc} is a pre-allocated computation time budget and Δt_o is the time interval between two successive observations as well as control updates. We then search over $(\tau, v^*(\tau))$ for the optimal perturbation time τ^* to globally minimize $\mathbb{E}[\nu(t_f)]$. There is only a finite number of such τ to consider since in practice we use numerical integration such as the Euler scheme with some step size Δt_c to compute the trajectories. In Ansari and Murphey (2016) the finite perturbation duration ϵ is also optimized using line search, but in this work we set ϵ as a tunable parameter to reduce the computation time. The complete algorithm is summarized in Algorithm 1. The call to the algorithm occurs every $\Delta t_o[s]$ in a receding-horizon fashion, after the new observation is incorporated in the belief.

2.2 General Belief Space Planning Problems

If none of the state is fully observable, the same stochastic SAC framework still applies almost as is to the belief state b . In this case we consider a continuous-discrete filter Xie et al. (2007) where the prediction step follows an ODE and the update step provides an instantaneous discrete jump. The hybrid dynamics for the belief vector yields

$$\begin{cases} \dot{b}(t_k) = g(b(t_k^-), y_k) \\ \dot{b}(t) = f(b(t), u(t)) \quad \forall t \in [t_k, t_{k+1}). \end{cases} \quad (36)$$

Algorithm 1 SACBP Control Update for Problems with Mixed Observability

INPUT: Current augmented state $s_0 = (p(t_0)^T, b(t_0)^T)^T$, nominal control u , perturbation duration ϵ

OUTPUT: Optimally perturbed control schedule u^ϵ

- 1: **for** $i = 1:N$ **do**
- 2: Forward-simulate nominal augmented state trajectory (1)(2) and sample observation sequence (y_1^i, \dots, y_T^i) along the augmented state trajectory.
- 3: Backward-simulate nominal adjoint trajectory ρ_p^i, ρ_b^i (29)(30) with sampled observations.
- 4: **end for**
- 5: Compute Monte Carlo estimate: $\mathbb{E}[\rho_p] \approx \frac{1}{N} \sum_{i=1}^N \rho_p^i$.
- 6: **for** $(\tau = t_0 + t_{calc} + \epsilon; \tau \leq t_0 + \Delta t_o; \tau \leftarrow \tau + \Delta t_c)$ **do**
- 7: Solve quadratic minimization (35) with (34). Store optimal value $\nu^*(\tau)$ and optimizer $v^*(\tau)$.
- 8: **end for**
- 9: $\tau^* \leftarrow \arg \min \nu^*(\tau), v^* \leftarrow v^*(\tau^*)$
- 10: $u^\epsilon \leftarrow \text{PerturbControlTrajectory}(u, v^*, \tau^*, \epsilon)$ (3)
- 11: **return** u^ϵ

Letting $\Psi(t) = \frac{\partial}{\partial \epsilon} b^\epsilon(t) \big|_{\epsilon=0}$, the variational equation is given by

$$\begin{cases} \Psi(t_k) = \frac{\partial}{\partial b} g(b(t_k^-), y_k) \Psi(t_k^-) \\ \dot{\Psi}(t) = \frac{\partial}{\partial b} f(b(t), u(t)) \Psi(t) \quad \forall t \in [t_k, t_{k+1}) \end{cases} \quad (37)$$

with initial condition $\Psi(\tau) = f(b(\tau), v) - f(b(\tau), u(\tau))$.

Let the total cost be of the form:

$$\int_{t_0}^{t_f} c(b(t), u(t)) dt + h(b(t_f)). \quad (38)$$

Under the given (τ, v) and (y_1, \dots, y_T) , the variation $\nu(t_f)$ of the total cost can be computed as

$$\begin{aligned} \nu(t_f) &= c(b(\tau), v) - c(b(\tau), u(\tau)) \\ &+ \int_{\tau}^{t_f} \frac{\partial}{\partial b} c(b(t), u(t))^T \Psi(t) dt \\ &+ \frac{\partial}{\partial b} h(b(t_f))^T \Psi(t_f). \end{aligned} \quad (39)$$

This is equivalent to

$$\begin{aligned} \nu(t_f) &= c(b(\tau), v) - c(b(\tau), u(\tau)) \\ &+ \rho(\tau)^T \{f(b(\tau), v) - f(b(\tau), u(\tau))\}, \end{aligned} \quad (40)$$

where ρ is the adjoint system that follows the dynamics:

$$\begin{cases} \rho(t_k^-) = \frac{\partial}{\partial b} g(b(t_k^-), y_k)^T \rho(t_k) \\ \dot{\rho}(t) = -\frac{\partial}{\partial b} c(b(t), u(t)) - \frac{\partial}{\partial b} f(b(t), u(t))^T \rho(t) \end{cases} \quad (41)$$

with the boundary condition $\rho(t_f) = \frac{\partial}{\partial b} h(b(t_f))$. Under the control-affine assumption for f and the additive quadratic control cost, the expected first order cost variation (40) yields

$$\begin{aligned} \mathbb{E}[\nu(t_f)] &= \frac{1}{2} v^T C_u v + \mathbb{E}[\rho(\tau)]^T H(b(\tau))(v - u(\tau)) \\ &- \frac{1}{2} u(\tau)^T C_u u(\tau), \end{aligned} \quad (42)$$

Algorithm 2 SACBP Control Update for General Belief Space Planning Problems

INPUT: Current belief state $b_0 = b(t_0)$, nominal control u , perturbation duration ϵ

OUTPUT: Optimally perturbed control schedule u^ϵ

- 1: **for** $i = 1:N$ **do**
- 2: Forward-simulate nominal belief state trajectory (36) and sample observation sequence (y_1^i, \dots, y_T^i) along the belief trajectory.
- 3: Backward-simulate nominal adjoint trajectory ρ^i (41) with sampled observations.
- 4: **end for**
- 5: Compute Monte Carlo estimate: $\mathbb{E}[\rho] \approx \frac{1}{N} \sum_{i=1}^N \rho^i$.
- 6: **for** $(\tau = t_0 + t_{calc} + \epsilon; \tau \leq t_0 + \Delta t_o; \tau \leftarrow \tau + \Delta t_c)$ **do**
- 7: Solve quadratic minimization (35) with (42). Store optimal value $\nu^*(\tau)$ and optimizer $v^*(\tau)$.
- 8: **end for**
- 9: $\tau^* \leftarrow \arg \min \nu^*(\tau), v^* \leftarrow v^*(\tau^*)$
- 10: $u^\epsilon \leftarrow \text{PerturbControlTrajectory}(u, v^*, \tau^*, \epsilon)$ (3)
- 11: **return** u^ϵ

where $H(b(\tau))$ is the control coefficient term in f .

Although it is difficult to state the general conditions under which this control-affine assumption holds, one can verify that the continuous-discrete EKF Xie et al. (2007) satisfies this property if the underlying system dynamics f_{sys} is control-affine.

$$\begin{cases} \dot{\mu}(t) = f_{sys}(\mu(t), u(t)) \\ \dot{\Sigma}(t) = A\Sigma + \Sigma A^T + Q \end{cases} \quad (43)$$

In the above continuous-time prediction equations, A is the Jacobian of the dynamics function $f_{sys}(x(t), u(t))$ evaluated at the mean $\mu(t)$ and Q is the process noise covariance. If f_{sys} is control-affine, so is A and therefore so is $\dot{\Sigma}$. Obviously $\dot{\mu}$ is control affine as well. As a result the dynamics for the belief vector $b = (\mu^T, \text{vec}(\Sigma)^T)^T$ satisfy the control-affine assumption.

Mirroring the approach in Section 2.1, we can use Monte Carlo sampling to estimate the expected value in (42). The resulting algorithm is presented in Algorithm 2.

2.3 Closed-loop Nominal Policy

In Sections 2.1 and 2.2 we assumed that the nominal control u was an open-loop control trajectory. However, one can think of a scenario where a nominal control is a closed-loop policy computed off-line, such as a discrete POMDP policy. Indeed, SACBP can also handle closed-loop nominal policies. Let π be a closed-loop nominal policy, which is a mapping from either an augmented state $s(t)$ or a belief state $b(t)$ to a control value $u(t)$. Due to the stochastic belief dynamics, the control values returned by π in the future are also stochastic for $t \geq t_1$. This is reflected when we forward-propagate the nominal dynamics. Specifically, each sampled trajectory has a different control trajectory in addition to a different observation sequence. However, the equations are still convex quadratic in v as shown below. For problems with

mixed observability, we have

$$\mathbb{E}[\nu(t_f)] = \frac{1}{2}v^T C_u v + \mathbb{E}[\rho_p(\tau)]^T H(p(\tau)) \{v - \pi(s(\tau))\} - \frac{1}{2}\pi(s(\tau))^T C_u \pi(s(\tau)). \quad (44)$$

The general belief space planning case also yields a similar equation:

$$\mathbb{E}[\nu(t_f)] = \frac{1}{2}v^T C_u v + \mathbb{E}[\rho(\tau)]^T H(b(\tau)) \{v - \pi(b(\tau))\} - \frac{1}{2}\pi(b(\tau))^T C_u \pi(b(\tau)). \quad (45)$$

Note that $s(\tau)$ and $b(\tau)$ are both deterministic since the first observation y_1 is not yet taken at $\tau < t_1$. The expectations in (44) and (45) can be estimated using Monte Carlo sampling. The forward-simulation of the nominal trajectory in Line 2 of Algorithms 1 and 2 is now with the closed loop policy π , and the equations in Line 7 need to be replaced with (44) and (45), respectively. However, the rest remains unchanged.

2.4 Computation Time Complexity

Let us analyze the time complexity of the SACBP algorithm. The bottleneck of the computation is when the forward-backward simulation is performed multiple times (lines 1–5 of Algorithms 1 and 2). The asymptotic complexity of this part is given by $O(N(\frac{t_f - t_0}{\Delta t_o})(M_{\text{forward}} + M_{\text{backward}}))$, where M_{forward} and M_{backward} are the times to respectively integrate the forward and backward dynamics between two successive observations. For a more concrete analysis let us use the Gaussian belief dynamics given by EKF as an example. For simplicity we assume the same dimension n for the state, the control, and the observation. The belief state has dimension $O(n^2)$. Using the Euler scheme, the forward integration takes $M_{\text{forward}} = O((\frac{\Delta t_o}{\Delta t_c} + 1)n^3)$ since evaluating continuous and discrete EKF equations are both $O(n^3)$. Computation of the continuous part of the costate dynamics (41) is dominated by the evaluation of the Jacobian $\frac{\partial f}{\partial b}$, which is $O(n^5)$ because $O(n^3)$ operations to evaluate f are carried out $O(n^2)$ times. The discrete part is also $O(n^5)$. Therefore, $M_{\text{backward}} = O((\frac{\Delta t_o}{\Delta t_c} + 1)n^5)$. Overall, the time complexity is $O(N(\frac{t_f - t_0}{\Delta t_o})(\frac{\Delta t_o}{\Delta t_c} + 1)n^5)$. This is asymptotically smaller in n than belief iLQG, which is $O(n^6)$. See Rafieisakhaei et al. (2017) for a comparison of time complexity among different belief space planning algorithms. We also remind the readers that SACBP is an online method and a naive implementation already achieves near real-time performance, computing control in less than 0.4[s]. By near real-time we mean that a naive implementation of SACBP requires approximately $3 \times t_{\text{calc}}$ to $7 \times t_{\text{calc}}$ time to compute an action that must be applied t_{calc} [s] in the future. We expect that parallelization in a GPU and a more efficient implementation will result in real-time computation for SACBP.

3 Analysis of Mode Insertion Gradient for Stochastic Hybrid Systems

The SACBP algorithm presented in Section 2 as well as the original SAC algorithm Ansari and Murphey (2016) both

rely on the local sensitivity analysis of the cost functional with respect to the control perturbation. This first-order sensitivity term (i.e. $\nu(t_f)$ in our notation) is known as the mode insertion gradient in the mode scheduling literature Egerstedt et al. (2006); Wardi and Egerstedt (2012). In Ansari and Murphey (2016) the notion of the mode insertion gradient has been generalized to handle a broader class of hybrid systems than discussed before. What remains to be seen is a further generalization of the mode insertion gradient to stochastic hybrid systems, such as the belief dynamics discussed in this paper. Indeed, the quantity we can optimize in (35) is essentially the expected value of the first-order sensitivity of the total cost. This is not to be confused with the first-order sensitivity of the expected total cost, which would be a natural generalization of the mode insertion gradient to stochastic systems. In general, those two quantities can be different, since the order of expectation and differentiation may not be swapped arbitrarily. In this section, we provide a set of sufficient conditions under which the order can be exchanged. By doing so we show that 1) the notion of mode insertion gradient can be generalized to stochastic hybrid systems, and 2) the SACBP algorithm optimizes this generalized mode insertion gradient. Through this analysis we will see that the SACBP algorithm has a guarantee that, in expectation it performs at least as good as the nominal policy for an appropriate choice of ϵ .

3.1 Assumptions

Let us begin with a set of underlying assumptions for the system dynamics, the control, and the cost functions. Without loss of generality, we assume that the system starts at time $t = 0$ and ends at $t = T$, with a sequence of T observations (y_1, \dots, y_T) made every unit time. For generality, we use notation x to represent the state variable of the system in this section, in place of b or s that respectively represented the belief state or the augmented state in Section 2. This means that the analysis presented here is not restricted to belief systems where the dynamics is governed by Bayesian filters, but rather applies to a broader class of systems.

Assumption 1. Control Model. *The controls are in $\tilde{C}^{0,m}[0, T]$, the space of piecewise continuous functions from $[0, T]$ into \mathbb{R}^m . We further assume that there exists some $\rho_{\max} < \infty$ such that for all $t \in [0, T]$, we have $u(t) \in B(0, \rho_{\max})$ where $B(0, \rho_{\max})$ is the closed Euclidean ball of radius ρ_{\max} centered at 0, i.e. $\|u(t)\|_2 \leq \rho_{\max}$. We denote this admissible control set by $U \triangleq \{u \in \tilde{C}^{0,m}[0, T] \mid \forall t \in [0, T] \ u(t) \in B(0, \rho_{\max})\}$.*

Remark 1. *The control model described above takes the form of an open-loop control, where time t determines the control signal $u(t)$. The generalization to closed-loop nominal policies are discussed in Appendix A. (See Remark 5.)*

Assumption 2. Dynamics Model. *Let $x_0 \in \mathbb{R}^{n_x}$ be the given initial state value at $t = 0$. Given a control $u \in U$ and a sequence of observations $(y_1, \dots, y_T) \in \mathbb{R}^{n_y} \times \dots \times \mathbb{R}^{n_y}$, the dynamics model is the following hybrid system with time-driven switching:*

$$x(t) \triangleq x_i(t) \ \forall t \in [i-1, i) \ \forall i \in \{1, 2, \dots, T\}, \quad (46)$$

where x_i is the i -th "mode" of the system state defined on $[i-1, i]$ as:

$$x_i(i-1) = g(x_{i-1}(i-1), y_{i-1}) \quad (47)$$

$$\dot{x}_i(t) = f(x_i(t), u(t)) \quad \forall t \in [i-1, i], \quad (48)$$

with $x(0) = x_1(0) = x_0$. We also define the final state as $x(T) \triangleq g(x_T(T), y_T)$.

For the transition functions f and g we assume the following:

(2a) the function $f: \mathbb{R}^{n_x} \times \mathbb{R}^m \rightarrow \mathbb{R}^{n_x}$ is continuously differentiable;

(2b) the function $g: \mathbb{R}^{n_x} \times \mathbb{R}^{n_y} \rightarrow \mathbb{R}^{n_x}$ is continuous. It is also differentiable in x ;

(2c) for function f , there exist constants $K_1 \in [1, \infty)$ and $K_2 \in (0, \infty)$ such that $\forall x', x'' \in \mathbb{R}^{n_x}$ and $\forall u', u'' \in B(0, \rho_{\max})$, the following relations hold:

$$\begin{aligned} \|f(x', u') - f(x'', u'')\|_2 \\ \leq K_1 (\|x' - x''\|_2 + \|u' - u''\|_2) \end{aligned} \quad (49)$$

$$\left\| \frac{\partial}{\partial x} f(x', u') \right\|_2 \leq K_2 \quad (50)$$

(2d) for function g , there exist finite non-negative constants K_3, K_4, K_5, K_6 and positive integers L_1, L_2 such that $\forall x \in \mathbb{R}^{n_x}$ and $\forall y \in \mathbb{R}^{n_y}$, the following relations hold:

$$\begin{aligned} \|g(x, y)\|_2 \leq K_3 + K_4 \|x\|_2^{L_1} + K_5 \|y\|_2^{L_2} \\ + K_6 \|x\|_2^{L_1} \|y\|_2^{L_2} \end{aligned} \quad (51)$$

$$\begin{aligned} \left\| \frac{\partial}{\partial x} g(x, y) \right\|_2 \leq K_3 + K_4 \|x\|_2^{L_1} + K_5 \|y\|_2^{L_2} \\ + K_6 \|x\|_2^{L_1} \|y\|_2^{L_2} \end{aligned} \quad (52)$$

Remark 2. Assumptions (2a) and (2c) are related to the existence and uniqueness of the solution to the differential equation (48) as well as the variational equation under control perturbation. (See Propositions 3 and 15 in Appendix A.) These assumptions are similar to Assumption 5.6.2 in *Elijah (1997)*. Assumptions (2b) and (2d) are the growth conditions on x across adjacent modes. Recall that in belief space planning where the system state x is the belief state b , the jump function g corresponds to the observation update of the Bayesian filter. The form of the bound in (51) and (52) allows a broad class of continuous functions to be considered as g , and is inspired by a few examples of the Bayesian update equations as presented below.

Proposition 1. Bounded Jump for Univariate Gaussian Distribution. Let $b = (\mu, s)^T \in \mathbb{R}^2$ be the belief state, where μ is the mean parameter and $s > 0$ is the variance. Suppose that the observation y is the underlying state $x \in \mathbb{R}$ corrupted by additive Gaussian white noise $v \sim \mathcal{N}(0, 1)$. Then, the Bayesian update function g for this belief system satisfies Assumption (2d).

Proof. The Bayesian update formula for this system is given by $g(b, y) = \hat{b} \triangleq (\hat{\mu}, \hat{s})^T$, where

$$\hat{\mu} = \mu + \frac{s}{s+1}(y - \mu) \quad (53)$$

$$\hat{s} = s - \frac{s^2}{s+1} \quad (54)$$

is the update step of the Kalman filter. Rearranging the terms, we have

$$g(b, y) = \frac{1}{s+1} \begin{pmatrix} \mu + sy \\ s \end{pmatrix} \quad (55)$$

and consequently,

$$\frac{\partial}{\partial b} g(b, y) = \frac{1}{(s+1)^2} \begin{pmatrix} s+1 & y \\ 0 & 1 \end{pmatrix}. \quad (56)$$

We will show that the function g satisfies Assumption (2d). For the bound on $g(b, y)$,

$$\|g(b, y)\|_2^2 = \frac{1}{(s+1)^2} \{(\mu + sy)^2 + s^2\} \quad (57)$$

$$\leq (\mu + sy)^2 + s^2 \quad (58)$$

$$\leq \|b\|_2^2 + \left(b^T \begin{pmatrix} y \\ 1 \end{pmatrix}\right)^2 \quad (59)$$

$$\leq \|b\|_2^2 (2 + \|y\|_2^2), \quad (60)$$

where we have used $(s+1)^2 \geq 1$ and the Cauchy-Schwarz inequality. Thus,

$$\|g(b, y)\|_2 \leq \|b\|_2 \sqrt{2 + \|y\|_2^2} \quad (61)$$

$$\leq \sqrt{2} \|b\|_2 + \|b\|_2 \|y\|_2 \quad (62)$$

Similarly, the bound on the Jacobian yields

$$\left\| \frac{\partial}{\partial b} g(b, y) \right\|_2^2 \leq \left\| \frac{\partial}{\partial b} g(b, y) \right\|_F^2 \quad (63)$$

$$= \frac{1}{(s+1)^4} \{(s+1)^2 + y^2 + 1\} \quad (64)$$

$$= \frac{1}{(s+1)^2} + \frac{1}{(s+1)^4} (y^2 + 1) \quad (65)$$

$$\leq 2 + \|y\|_2^2. \quad (66)$$

Therefore,

$$\left\| \frac{\partial}{\partial b} g(b, y) \right\|_2 \leq \sqrt{2} + \|y\|_2 \quad (67)$$

This shows that the jump function g for the above univariate Gaussian model satisfies Assumption (2d) with $(K_3, K_4, K_5, K_6) = (\sqrt{2}, \sqrt{2}, 1, 1)$ and $(L_1, L_2) = (1, 1)$.

Proposition 2. Bounded Jump for Categorical Distribution. Let $b = (b_1, \dots, b_n)^T \in \mathbb{R}^n$ be the n -dimensional belief state representing the categorical distribution over the underlying state $x \in \{1, \dots, n\}$. We choose the unnormalized form where the probability of $x = i$ is given by $b_i / \sum_{i=1}^n b_i$. Let the observation $y \in \{1, \dots, m\}$ be modeled by a conditional probability mass function $p(y | x) \in [0, 1]$. Then, the Bayesian update function g for this belief system satisfies Assumption (2d).

Proof. The Bayes rule gives $g(b, y) = \hat{b} \triangleq (\hat{b}_1, \dots, \hat{b}_n)$, where

$$\begin{pmatrix} \hat{b}_1 \\ \hat{b}_2 \\ \vdots \\ \hat{b}_n \end{pmatrix} = \begin{pmatrix} p(y | 1)b_1 \\ p(y | 2)b_2 \\ \vdots \\ p(y | n)b_n \end{pmatrix}. \quad (68)$$

Therefore, we can easily bound the norm of the posterior belief \hat{b} by

$$\|g(b, y)\|_2 = \|\hat{b}\|_2 \leq \|b\|_2, \quad (69)$$

as $p(y | x) \leq 1$. The Jacobian is simply the diagonal matrix $\text{diag}(p(y | 1), \dots, p(y | n))$, and hence

$$\left\| \frac{\partial}{\partial b} g(b, y) \right\|_2 \leq 1. \quad (70)$$

This shows that the jump function g for the categorical belief model above satisfies Assumption (2d) with $(K_3, K_4, K_5, K_6) = (1, 1, 0, 0)$ and $(L_1, L_2) = (1, 1)$.

Assumption 3. Cost Model. *The instantaneous cost $c: \mathbb{R}^{n_x} \times \mathbb{R}^m \rightarrow \mathbb{R}$ is continuous. It is also continuously differentiable in x . The terminal cost $h: \mathbb{R}^{n_x} \rightarrow \mathbb{R}$ is differentiable. Furthermore, we assume that there exist finite non-negative constants K_7, K_8 and a positive integer L_3 such that for all $x \in \mathbb{R}^{n_x}$ and $u \in B(0, \rho_{\max})$, the following relations hold:*

$$|c(x, u)| \leq K_7 + K_8 \|x\|_2^{L_3} \quad (71)$$

$$\left\| \frac{\partial}{\partial x} c(x, u) \right\|_2 \leq K_7 + K_8 \|x\|_2^{L_3} \quad (72)$$

$$|h(x)| \leq K_7 + K_8 \|x\|_2^{L_3} \quad (73)$$

$$\left\| \frac{\partial}{\partial x} h(x) \right\|_2 \leq K_7 + K_8 \|x\|_2^{L_3}. \quad (74)$$

Remark 3. Assumption 3 is to guarantee that the cost function is integrable with respect to stochastic observations, which are introduced in Assumption 4. Note that even though the above bound is not general enough to apply to all analytic functions, it does include all finite order polynomials of $\|x(t)\|_2$ and $\|u(t)\|_2$, for example, since $\|u(t)\|_2$ is bounded by Assumption 1.

Assumption 4. Stochastic Observations. *Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space. Let (Y_1, \dots, Y_T) be a sequence of random vectors in \mathbb{R}^{n_y} defined on this space, representing the sequence of observations. Assume that for each Y_i all the moments of the ℓ^2 norm is finite. That is,*

$$\forall i \in \{1, \dots, T\} \quad \forall k \in \mathbb{N} \quad \mathbb{E} [\|Y_i\|_2^k] < \infty. \quad (75)$$

Definition 1. Perturbed Control. *Let $u \in U$ be a control. For $\tau \in (0, 1)$ and $v \in B(0, \rho_{\max})$, define the perturbed control u^ϵ by*

$$u^\epsilon(t) \triangleq \begin{cases} v & \text{if } t \in (\tau - \epsilon, \tau] \\ u(t) & \text{otherwise,} \end{cases} \quad (76)$$

where $\epsilon \in [0, \tau]$. By definition if $\epsilon = 0$ then u^ϵ is the same as u . We assume that the nominal control $u(t)$ is left continuous in t at $t = \tau$.

3.2 Main Results

The main result of the analysis is the following theorem.

Theorem 1. Mode Insertion Gradient. *Suppose that Assumptions 1 – 4 are satisfied. For a given (τ, v) , let u^ϵ denote the perturbed control of the form (76). The perturbed control u^ϵ and the stochastic observations (Y_1, \dots, Y_T) result in the stochastic perturbed state trajectory x^ϵ . For such u^ϵ and x^ϵ , let us define the mode insertion gradient of the expected total cost as*

$$\frac{\partial_+}{\partial \epsilon} \mathbb{E} \left[\int_0^T c(x^\epsilon(t), u^\epsilon(t)) dt + h(x^\epsilon(T)) \right] \Big|_{\epsilon=0}. \quad (77)$$

Then, this right derivative exists and we have

$$\begin{aligned} \frac{\partial_+}{\partial \epsilon} \mathbb{E} \left[\int_0^T c(x^\epsilon(t), u^\epsilon(t)) dt + h(x^\epsilon(T)) \right] \Big|_{\epsilon=0} &= c(x(\tau), v) - c(x(\tau), u(\tau)) \\ &+ \mathbb{E} \left[\int_\tau^T \frac{\partial}{\partial x} c(x(t), u(t))^T \Psi(t) dt \right. \\ &\quad \left. + \frac{\partial}{\partial x} h(x(T))^T \Psi(T) \right], \end{aligned} \quad (78)$$

where $\Psi(t) = \frac{\partial}{\partial \epsilon} x^\epsilon(t) \Big|_{\epsilon=0}$ is the state variation.

The proof of the theorem is deferred to Appendix A. One can see that the mode insertion gradient (77) is a natural generalization of the ones discussed in Egerstedt et al. (2006); Wardi and Egerstedt (2012); Ansari and Murphey (2016) to stochastic hybrid systems. Furthermore, by comparing (78) with (39) it is apparent that the right hand side of (78) is mathematically equivalent to $\mathbb{E}[\nu(t_f)]$, the quantity to be optimized with the SACBP algorithm in Section 2.

The fact that SACBP optimizes (77) leads to a certain performance guarantee of the algorithm. In the open-loop nominal control case, the term $\mathbb{E}[\nu(t_f)]$ as in (34) or (42) becomes 0 if the control perturbation v is equal to the nominal control $u(\tau)$. Therefore, as long as $u(\tau)$ is a feasible solution to (35) the optimal value is guaranteed to be less than or equal to zero. Furthermore, in expectation the actual value of $\mathbb{E}[\nu(t_f)]$ matches the one approximated with samples, since the Monte Carlo estimate is unbiased. In other words, the perturbation (τ^*, v^*) computed by the algorithm is expected to result in a non-negative mode insertion gradient. If the mode insertion gradient is negative, there always exists a sufficiently small $\epsilon > 0$ such that the expected total cost is decreased by the control perturbation. In the corner case that the mode insertion gradient is zero, one can set $\epsilon = 0$ to not perturb the control at all. Therefore, for an appropriate choice of ϵ the expected performance of the SACBP algorithm is at least as good as that of the nominal control.

The same discussion holds for the case of closed-loop nominal control policies, when the expression for $\mathbb{E}[\nu(t_f)]$ is given by (44) or (45). Therefore, the expected worst-case performance of the algorithm is lower-bounded by that of the nominal policy.

If we have a reasonable nominal policy, this turns into a rather strong guarantee compared to other belief space planning algorithms that can exhibit extremely large regret. For example, the UCT algorithm is known to have a very poor expected performance in the worst case due to its over-optimistic behavior [Coquelin and Munos \(2007\)](#). Therefore, the Monte Carlo Tree Search methods such as POMCP [Silver and Veness \(2010\)](#) or MCTS-DPW that are reliant on the UCT algorithm inherit this issue.

4 Simulation Results

We evaluated the performance of SACBP in the following simulation studies: (i) active multi-target tracking with range-only observations; (ii) object manipulation under model uncertainty. All the computation was performed on a desktop computer with Intel Core i7-6800K CPU and 62.8GB RAM. The Monte Carlo sampling of SACBP was parallelized on the CPU.

4.1 Active Multi-Target Tracking with Range-only Observations

This problem focuses on pure information gathering, namely identifying where the moving targets are in the environment. In doing so, the surveillance robot modeled as a single integrator can only use relative distance observations. The robot's position p is fully observable and the transitions are deterministic. Assuming perfect data association, the observation for target i is $d_i = \|q_i - p + v_i\|_2$, where q_i is the true target position and v_i is zero-mean Gaussian white noise with state-dependent covariance $R(p, q_i) = R_0 + \|q_i - p\|_2 R_1$. We used $0.01I_{2 \times 2}$ for the nominal noise R_0 . The range-dependent noise $R_1 = 0.001I_{2 \times 2}$ degrades the observation quality as the robot gets farther from the target. The discrete-time UKF was employed for state estimation in tracking 20 independent targets. The target dynamics are modeled by a 2D Brownian motion with covariance $Q = 0.1I_{2 \times 2}$. Similarly to [Spinello and Stilwell \(2010\)](#), an approximated observation covariance $R(p, \mu_i)$ was used in the filter to obtain tractable estimation results, where μ_i is the most recent mean estimate of q_i .

The SACBP algorithm generated the continuous robot trajectory over 200[s] with planning horizon $t_f - t_0 = 2[s]$, update interval $\Delta t_o = 0.2[s]$, perturbation duration $\epsilon = 0.16[s]$, and $N = 10$ Monte Carlo samples. The Euler scheme was used for integration with $\Delta t_c = 0.01[s]$. The Jacobians and the gradients were computed either analytically or using an automatic differentiation tool [Revels et al. \(2016\)](#) to retain both speed and precision. In this simulation $t_{\text{calc}} = 0.05[s]$ was assumed no matter how long the actual control update took. We used $c(p, b, u) = 0.05u^T u$ for the running cost and $h(p, b) = \sum_{i=1}^{20} \exp(\text{entropy}(b_i))$ for the terminal cost, with an intention to reduce the worst-case uncertainty among the targets. This expression for $h(p, b)$ is equivalent to:

$$h(p, b) = \sum_{i=1}^{20} \sqrt{\det(2\pi e \Sigma_i)}, \quad (79)$$

where Σ_i is the covariance for the i -th target. The nominal control was constantly zero.

We compared SACBP against three benchmarks: (i) a greedy algorithm based on the gradient descent of terminal cost h , similar to [Schwager et al. \(2017\)](#); (ii) MCTS-DPW [Couëtoux et al. \(2011\)](#); [Egorov et al. \(2017\)](#) in the Gaussian belief space; (iii) projection-based trajectory optimization for ergodic exploration [Miller and Murphey \(2013\)](#); [Miller et al. \(2016\)](#); [Dressel and Kochenderfer \(2018\)](#). We also implemented the belief iLQG algorithm, but the policy did not converge for this problem. We suspect that the non-convex terminal cost h contributed to this behavior, which in fact violates one of the underlying assumptions made in the paper [van den Berg et al. \(2012\)](#).

MCTS-DPW used the same planning horizon as SACBP, however it drew $N = 15$ samples from the belief tree so the computation time of the two algorithms matched approximately. Ergodic trajectory optimization is not a belief space planning approach but has been used in the active sensing literature. Beginning with the nominal control of zero, it locally optimized the ergodicity of the trajectory with respect to the spatial information distribution based on Fisher information. This optimization was open-loop since the future observations were not considered. As a new observation became available, the distribution and the trajectory were recomputed. All the controllers were saturated at the same limit. The results presented in Figure 1 clearly indicates a significant performance improvement of SACBP while achieving near real-time computation. More notably, SACBP generated a trajectory that periodically revisited the two groups whereas other methods failed to do so (Figure 2). With SACBP the robot was moving into one of the four diagonal directions for most of the time. This is plausible, as SACBP solves the quadratic program with a box input constraint (35), which tends to find optimal solutions at the corners. MCTS-DPW resulted in a highly non-smooth trajectory and failed to explore the environment. The greedy approach improved the smoothness, but the robot eventually followed a cyclic trajectory in a small region of the environment. To our surprise, the ergodic method did not generate a trajectory that covers the two groups of the targets. This is likely due to the use of a projection-based trajectory optimization method, which has been recently found to perform rather poorly with rapid replanning [Dressel \(2018\)](#).

4.2 Object Manipulation under Model Uncertainty

This problem is identical to the model-based Bayesian reinforcement learning problem studied in [Slade et al. \(2017\)](#), therefore a detailed description of the nonlinear dynamics and the observation models are omitted. See Figure 3 for the illustration of the environment. A 2D robot attached to a rigid body object applies forces and torques to move the object to the origin. The object's mass, moment of inertia, moment arm lengths, and linear friction coefficient are unknown. These parameters as well as the object's 2D state need to be estimated using EKF, with noisy sensors which measure the robot's position, velocity, and acceleration in the global frame. The same values for $t_f - t_0$, Δt_o , Δt_c , t_{calc} as in the previous problem were assumed. SACBP used $\epsilon = 0.04[s]$ and $N = 10$.

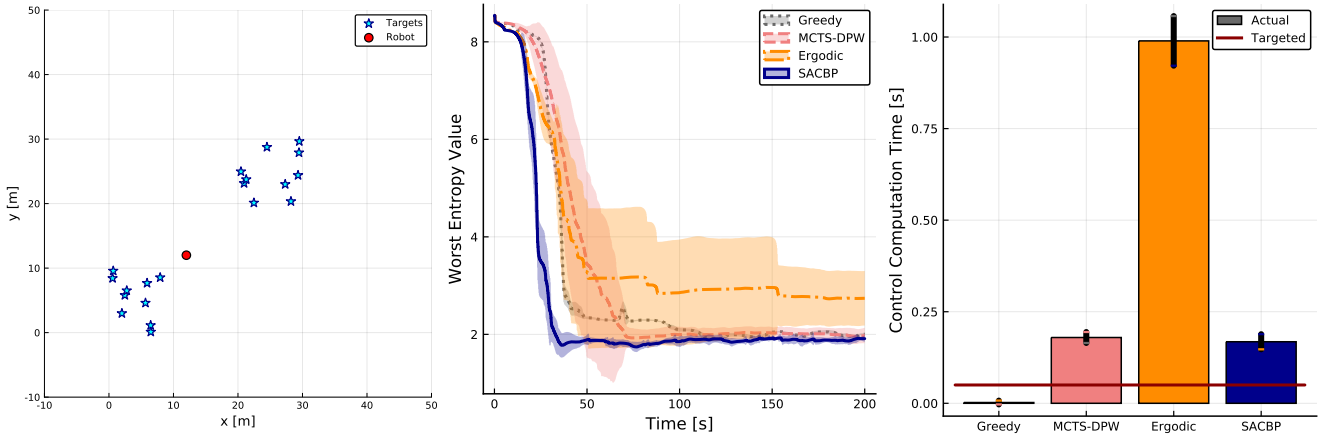


Figure 1. (Left) Simulation environment with 20 targets and a surveillance robot. (Middle) The history of the worst entropy value among the targets averaged over 20 runs with the standard deviation. With the budget of 10 Monte Carlo samples, SACBP had small variance and consistently outperformed the other benchmarks on average. (Right) Computation time of SACBP achieved a reasonable value compared with the benchmarks, only 0.11[s] slower than the targeted value, i.e., simulated t_{calc} .

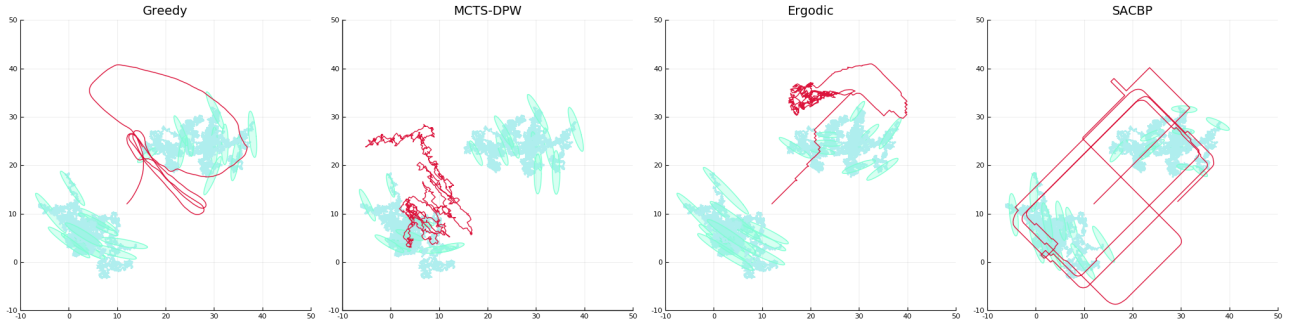


Figure 2. Sample robot trajectories (depicted in red) generated by each algorithm. Greedy, MCTS-DPW, and Ergodic did not result in a trajectory that fully covers the two groups of the targets, whereas SACBP periodically revisited both of them. With SACBP, the robot traveled into one of the four diagonal directions for most of the time. This is due to the fact that SACBP optimizes a convex quadratic under a box saturation constraint, which tends to find optimal solutions at the corners. In all the figures, the blue lines represent the target trajectories and the ellipses are 99% error ellipses.

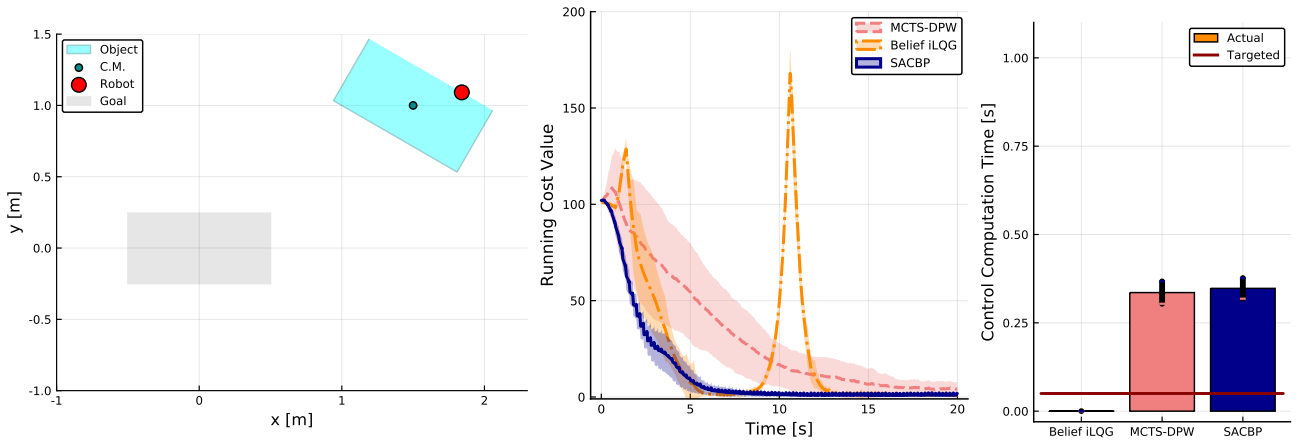


Figure 3. (Left) The robot is attached to the rectangular object. (Middle) The history of the true running cost $\frac{1}{2}x^T C_x x + \frac{1}{2}u^T C_u u$ averaged over 20 cases. SACBP with $N = 10$ samples successfully brought the cost to almost 0, meaning that the goal was reached. MCTS-DPW with $N = 190$ was not as successful. Belief iLQG resulted in large overshoots around 2[s] and 11[s]. (Right) Computation time of SACBP increased from the multi-target tracking problem due to increased complexity related to the continuous-discrete belief dynamics. Note that belief iLQG took 5945[s] to derive the policy off-line, although the average online execution time was only 4.0×10^{-5} [s] per iteration.

The nominal control was a closed-loop position controller whose input was the mean x-y position and the rotation estimates of the object. The cost function was quadratic in the true state x and control u , given by $\frac{1}{2}x^T C_x x + \frac{1}{2}u^T C_u u$. Taking expectations yielded the equivalent cost in the Gaussian belief space $c(b, u) = \frac{1}{2}\mu^T C_x \mu + \frac{1}{2}\text{tr}(C_x \Sigma) + \frac{1}{2}u^T C_u u$, where Σ is the covariance matrix. We let terminal cost h be the same as c except without the control term.

We compared SACBP against (i) MCTS-DPW in the Gaussian belief space and (ii) belief iLQG. We allowed MCTS-DPW to draw $N = 190$ samples to set the computation time comparable to SACBP. As suggested in [Slade et al. \(2017\)](#), MCTS-DPW used the position controller mentioned above as the rollout policy. Similarly, belief iLQG was initialized with a nominal trajectory generated by the same position controller. Note that both MCTS-DPW and belief iLQG computed controls for the discrete-time models whereas SACBP directly used the continuous-time model. However the simulation was all performed in continuous time, meaning that the control for MCTS-DPW and belief iLQG remained constant over each $\Delta t_o[s]$ interval. This could explain the overshoot of the belief iLQG trajectory around 2[s] in Figure 3. Another large overshoot around 11[s] is likely due to the locally optimal behavior of the iLQG solver. Overall, the results presented in Figure 3 demonstrate that SACBP succeeded in this task with only 10 Monte Carlo samples, reducing the running cost to almost 0 within 10[s]. Although the computation time increased from the previous problem due to the continuous-discrete filtering, it still achieved near real-time performance and much shorter than belief iLQG, which took 5945[s] until convergence in our implementation.

5 Conclusions and Future Work

In this paper we have presented SACBP, a novel belief space planning algorithm for continuous-time dynamical systems. We have viewed the stochastic belief dynamics as a hybrid system with time-driven switching and derived the optimal control perturbation based on the perturbation theory of differential equations. The resulting algorithm extends the framework of SAC to stochastic belief dynamics and is highly parallelizable to run in near real-time. The rigorous mathematical analysis showed that the notion of mode insertion gradient can be generalized to stochastic hybrid systems, which leads to the property of SACBP that the algorithm is expected to perform at least as good as the nominal policy for an appropriate choice of the perturbation duration. Through an extensive simulation study we have confirmed that SACBP outperforms other algorithms including a greedy algorithm, a local trajectory optimization method, and an approximate dynamic programming approach. In future work we are interested to consider a distributed multi-robot version of SACBP as well as problems with hard state constraints. We also plan to provide additional case studies for more complex belief distributions with efficient implementation.

Funding

Toyota Research Institute ("TRI") provided funds to assist the authors with their research but this article solely reflects the

opinions and conclusions of its authors and not TRI or any other Toyota entity. This work was also supported in part by NSF grant CMMI1562335, and a JASSO fellowship. The authors are grateful for this support.

References

- Ansari AR and Murphey TD (2016) Sequential Action Control: Closed-Form Optimal Control for Nonlinear and Nonsmooth Systems. *IEEE Transactions on Robotics* 32(5): 1196–1214.
- Bajcsy R (1988) Active perception. In: *Proceedings of the IEEE*, volume 76. pp. 966–1005.
- Bourbaki N and Spain P (2004) *Elements of Mathematics Functions of a Real Variable: Elementary Theory*. Berlin, Heidelberg: Springer Berlin Heidelberg.
- Bourgault F, Makarenko A, Williams S, Grocholsky B and Durrant-Whyte H (2002) Information based adaptive robotic exploration. In: *IEEE/RSJ International Conference on Intelligent Robots and System*, volume 1. IEEE, pp. 540–545.
- Browne CB, Powley E, Whitehouse D, Lucas SM, Cowling PI, Rohlfshagen P, Tavener S, Perez D, Samothrakis S and Colton S (2012) A survey of monte carlo tree search methods. *IEEE Transactions on Computational Intelligence and AI in Games* 4(1): 1–43.
- Coquelin PA and Munos R (2007) Bandit algorithms for tree search. In: *Proceedings of the Twenty-Third Conference on Uncertainty in Artificial Intelligence*, UAI'07. Arlington, Virginia, United States: AUAI Press, pp. 67–74.
- Couëtoux A, Hoock JB, Sokolovska N, Teytaud O and Bonnard N (2011) Continuous Upper Confidence Trees. In: *2011 International Conference on Learning and Intelligent Optimization*. Springer, Berlin, Heidelberg, pp. 433–445.
- Diestel J and Uhl J (1977) *Vector Measures*. American Mathematical Society.
- Dressel L and Kochenderfer MJ (2018) Tutorial on the generation of ergodic trajectories with projection-based gradient descent. *IET Cyber-Physical Systems: Theory & Applications*.
- Dressel LK (2018) *Efficient and Low-cost Localization of Radio Sources with an Autonomous Drone*. PhD Thesis, Stanford University.
- Egerstedt M, Wardi Y and Axelsson H (2006) Transition-time optimization for switched-mode dynamical systems. *IEEE Transactions on Automatic Control* 51(1): 110–115.
- Egorov M, Sunberg ZN, Balaban E, Wheeler TA, Gupta JK and Kochenderfer MJ (2017) Pomdps. jl: A framework for sequential decision making under uncertainty. *Journal of Machine Learning Research* 18(26): 1–5.
- Elijah P (1997) *Optimization: Algorithms and Consistent Approximations*. Springer Verlage Publications.
- Erez T and Smart WD (2010) A scalable method for solving high-dimensional continuous pomdps using local approximation. In: *Proceedings of the Twenty-Sixth Conference on Uncertainty in Artificial Intelligence*, UAI'10. Arlington, Virginia, United States: AUAI Press, pp. 160–167.
- Gowrisankaran K (1972) Measurability of functions in product spaces. *Proceedings of the American Mathematical Society* 31(2): 485–488.
- Heemels W, Lehmann D, Lunze J and Schutter BD (2009) Introduction to hybrid systems. In: Lunze J and Lamnabhi-Lagarigue F (eds.) *Handbook of Hybrid Systems Control* –

- Theory, Tools, Applications*, chapter 1. Cambridge University Press, pp. 3–30.
- Hollinger GA and Sukhatme GS (2014) Sampling-based robotic information gathering algorithms. *The International Journal of Robotics Research* 33(9): 1271–1287.
- Kaelbling LP, Littman ML and Cassandra AR (1998) Planning and acting in partially observable stochastic domains. *Artif. Intell.* 101(1-2): 99–134.
- Le Ny J and Pappas GJ (2009) On trajectory optimization for active sensing in Gaussian process models. In: *Proceedings of the 48th IEEE Conference on Decision and Control (CDC) held jointly with 2009 28th Chinese Control Conference*. IEEE, pp. 6286–6292.
- Liberzon D (2011) *Calculus of variations and optimal control theory: A concise introduction*. Princeton University Press.
- Mavrommati A, Tzorakoleftherakis E, Abraham I and Murphey TD (2018) Real-time area coverage and target localization using receding-horizon ergodic exploration. *IEEE Transactions on Robotics* : 62–80.
- Mihaylova L, Lefebvre T, Bruyninckx H, Gadeyne K and De Schutter J (2002) Active Sensing for Robotics - A Survey. In: *Proceedings of 5th International Conference on Numerical Methods and Applications*. pp. 316–324.
- Miller LM and Murphey TD (2013) Trajectory optimization for continuous ergodic exploration. In: *American Control Conference (ACC), 2013*. IEEE, pp. 4196–4201.
- Miller LM, Silverman Y, MacIver MA and Murphey TD (2016) Ergodic exploration of distributed information. *IEEE Transactions on Robotics* 32(1): 36–52.
- Nishimura H and Schwager M (2018a) Active motion-based communication for robots with monocular vision. In: *2018 IEEE International Conference on Robotics and Automation (ICRA)*. pp. 2948–2955.
- Nishimura H and Schwager M (2018b) Sacbp: Belief space planning for continuous-time dynamical systems via stochastic sequential action control. In: *The 13th International Workshop on the Algorithmic Foundations of Robotics (WAFR)*. Mérida, México.
- Patil S, Kahn G, Laskey M, Schulman J, Goldberg K and Abbeel P (2014) Scaling up gaussian belief space planning through covariance-free trajectory optimization and automatic differentiation. In: *WAFR, Springer Tracts in Advanced Robotics*, volume 107. Springer, pp. 515–533.
- Platt R, Tedrake R, Kaelbling L and Lozano-Perez T (2010) Belief space planning assuming maximum likelihood observations. In: *Robotics Science and Systems Conference (RSS)*.
- Popovi M, Hitz G, Nieto J, Sa I, Siegwart R and Galceran E (2017) Online informative path planning for active classification using uavs. In: *2017 IEEE International Conference on Robotics and Automation (ICRA)*. pp. 5753–5758.
- Rafieisakhaei M, Chakravorty S and Kumar PR (2017) T-lqg: Closed-loop belief space planning via trajectory-optimized lqg. In: *2017 IEEE International Conference on Robotics and Automation (ICRA)*. pp. 649–656.
- Revels J, Lubin M and Papamarkou T (2016) Forward-mode automatic differentiation in julia. *arXiv:1607.07892 [cs.MS]*
- Schwager M, Dames P, Rus D and Kumar V (2017) A multi-robot control policy for information gathering in the presence of unknown hazards. In: *Robotics Research : The 15th International Symposium ISRR*. Springer International Publishing, pp. 455–472.
- Seekircher A, Laue T and Röfer T (2011) Entropy-based active vision for a humanoid soccer robot. In: *RoboCup 2010: Robot Soccer World Cup XIV*, volume 6556 LNCS. Springer Berlin Heidelberg, pp. 1–12.
- Silver D and Veness J (2010) Monte-carlo planning in large pomdps. In: Lafferty JD, Williams CKI, Shawe-Taylor J, Zemel RS and Culotta A (eds.) *Advances in Neural Information Processing Systems 23*. Curran Associates, Inc., pp. 2164–2172.
- Slade P, Culbertson P, Sunberg Z and Kochenderfer M (2017) Simultaneous active parameter estimation and control using sampling-based bayesian reinforcement learning. In: *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. pp. 804–810.
- Somani A, Ye N, Hsu D and Lee WS (2013) Despot: Online pomdp planning with regularization. In: *Proceedings of the 26th International Conference on Neural Information Processing Systems - Volume 2, NIPS'13*. USA: Curran Associates Inc., pp. 1772–1780.
- Spinello D and Stilwell DJ (2010) Nonlinear estimation with state-dependent gaussian observation noise. *IEEE Transactions on Automatic Control* 55(6): 1358–1366.
- Sunberg Z and Kochenderfer MJ (2017) POMCPOW: an online algorithm for pomdps with continuous state, action, and observation spaces. *CoRR* abs/1709.06196.
- van den Berg J, Patil S and Alterovitz R (2012) Motion planning under uncertainty using iterative local optimization in belief space. *The International Journal of Robotics Research* 31(11): 1263–1278.
- Wardi Y and Egerstedt M (2012) Algorithm for optimal mode scheduling in switched systems. In: *2012 American Control Conference (ACC)*. pp. 4546–4551.
- Xie L, Popa D and Lewis FL (2007) *Optimal and robust estimation: with an introduction to stochastic control theory*. CRC press.

A Detailed Analysis of Mode Insertion Gradient for Stochastic Hybrid Systems

In this appendix, we provide a thorough analysis of the stochastic hybrid systems with time-driven switching that satisfy Assumptions 1 – 4. Our goal is to prove Theorem 1.

A.1 Nominal Trajectory under Specific Observations

First, we analyze the properties of the system $x(t)$ for $t \in [0, T]$ under a given initial condition x_0 , control $u \in U$, and a specific sequence of observations (y_1, \dots, y_T) sampled from (Y_1, \dots, Y_T) .

Proposition 3. Existence and Uniqueness of Solutions. *Given a control $u \in U$ and a sequence of observations (y_1, \dots, y_T) , the system $x(t)$ starting at x_0 has a unique solution for $t \in [0, T]$.*

Proof. We will show that each x_i for $i \in \{1, \dots, T\}$ has a unique solution, and thus x is uniquely determined as a whole. First, by Assumption (2a), (2c) and the Picard Lemma

(Lemma 5.6.3 in [Elijah \(1997\)](#)), the differential equation

$$\dot{x}_1(t) = f(x_1(t), u(t)) \quad (80)$$

with initial condition $x_1(0) = x_0$ has a solution for $t \in [0, 1]$. Furthermore, Proposition 5.6.5 in [Elijah \(1997\)](#) assures that the solution x_1 is unique under Assumption (2a) and (2c). This guarantees that the initial condition for x_2 defined by $x_2(1) = g(x_1(1), y_1)$ is unique. Therefore, proceeding by induction each x_1, \dots, x_T has a unique solution, which completes the proof.

Corollary 4. Right Continuity. *Given a control $u \in U$ and a sequence of observations (y_1, \dots, y_T) , the system $x(t)$ starting at x_0 is right continuous in t on $[0, T]$.*

Proof. By Proposition 3 each x_i has a unique solution that follows $\dot{x}_i = f(x_i, u)$. Clearly each x_i is continuous on $[i-1, i]$, which proves that $x(t) \triangleq x_i(t) \forall t \in [i-1, i] \forall i \in \{1, 2, \dots, T\}$ with $x(T) \triangleq g(x_T(T), y_T)$ is right continuous on $[0, T]$.

Lemma 5. Let ξ_i denote the initial condition for x_i . Then, there exists a constant $K_9 < \infty$ such that for all $i \in \{1, \dots, T\}$,

$$\forall t \in [i-1, i] \quad \|x_i(t)\|_2 \leq (1 + \|\xi_i\|_2) e^{K_9} \quad (81)$$

Proof. Using Assumption (2a) and (2c), the claim follows directly from Proposition 5.6.5 in [Elijah \(1997\)](#).

Proposition 6. Lipschitz Continuity. *For each $i \in \{1, \dots, T\}$, let ξ'_i and ξ''_i be two distinct initial conditions for x_i . Furthermore, let u' and u'' be two controls from U . The pairs (ξ', u') and (ξ'', u'') respectively define two solutions x'_i and x''_i to ODE $\dot{x}_i = f(x_i, u)$ over $[i-1, i]$. Then, there exists an $L < \infty$, independent of ξ'_i, ξ''_i, u' and u'' , such that*

$$\forall i \in \{1, \dots, T\} \quad \forall t \in [i-1, i] \quad \|x'_i(t) - x''_i(t)\|_2 \leq L \left(\|\xi'_i - \xi''_i\|_2 + \int_{i-1}^i \|u'(t) - u''(t)\|_2 dt \right). \quad (82)$$

Proof. The proof is similar to that of Lemma 5.6.7 in [Elijah \(1997\)](#). Making use of the Picard Lemma (Lemma 5.6.3 in [Elijah \(1997\)](#)) and Assumption (2c), we obtain

$$\|x'_i(t) - x''_i(t)\|_2 \leq e^{K_1} \left(\|\xi'_i - \xi''_i\|_2 + K_1 \int_{i-1}^i \|u'(t) - u''(t)\|_2 dt \right). \quad (83)$$

As $K_1 \geq 1$,

$$\|x'_i(t) - x''_i(t)\|_2 \leq K_1 e^{K_1} \left(\|\xi'_i - \xi''_i\|_2 + \int_{i-1}^i \|u'(t) - u''(t)\|_2 dt \right). \quad (84)$$

Defining $L \triangleq K_1 e^{K_1} < \infty$ completes the proof.

Corollary 7. Uniform Continuity in Initial Conditions. *Let $u \in U$ be a given control. For each $i \in \{1, \dots, T\}$, let ξ'_i and ξ''_i be two distinct initial conditions for x_i . The pairs (ξ', u) and (ξ'', u) respectively define two solutions x'_i and x''_i to ODE $\dot{x}_i = f(x_i, u)$ over $[i-1, i]$. Note that they share*

the same control but have different initial conditions, unlike Proposition 6. Then, for any $\epsilon > 0$ there exists $\delta > 0$ such that

$$\begin{aligned} \forall i \in \{1, \dots, T\} \quad \forall t \in [i-1, i] \quad & \|\xi'_i - \xi''_i\|_2 < \delta \\ \Rightarrow \|x'_i(t) - x''_i(t)\|_2 & < \epsilon. \end{aligned} \quad (85)$$

Proof. By Proposition 6, we find

$$\|x'_i(t) - x''_i(t)\|_2 \leq L \|\xi'_i - \xi''_i\|_2 \quad (86)$$

for all $i \in \{1, \dots, T\}$ and $t \in [i-1, i]$, where $L < \infty$. Take any $\delta < \frac{\epsilon}{L}$ to prove the claim.

Proposition 8. Continuity in Observations. *Given a control $u \in U$, the map $(y_1, \dots, y_T) \mapsto x(t)$ is continuous for all $t \in [0, T]$, where x represents the solution to the system under Assumption 2 starting at $x(0) = x_0$.*

Proof. We will show the continuity of $(y_1, \dots, y_T) \mapsto x_i(t)$ for each $i \in \{1, \dots, T\}$ by mathematical induction. First, by Assumption 2 the value of $x_1(t)$ is solely determined by x_0 and u . Therefore, for any $t \in [0, 1]$ the function $(y_1, \dots, y_T) \mapsto x_1(t)$ is a constant map, which is continuous. Next, suppose that $\forall t \in [i-1, i] \quad (y_1, \dots, y_T) \mapsto x_i(t)$ is continuous for some $i \in \{1, \dots, T-1\}$. Now consider x_{i+1} . Let $F_{i+1}(\xi_{i+1}, t)$ be the map from an initial condition $x_{i+1}(i) = \xi_{i+1}$ to the solution x_{i+1} at $t \in [i, i+1]$ under the given u . In other words, $F_{i+1}(\xi_{i+1}, t)$ is equivalent to the integral equation

$$F_{i+1}(\xi_{i+1}, t) \triangleq \xi_{i+1} + \int_i^t f(x_{i+1}(a), u(a)) da \quad (87)$$

and takes initial condition ξ_{i+1} as well as time t as its arguments. Substituting $\xi_{i+1} = g(x_i(i), y_i)$, $F_{i+1}(g(x_i(i), y_i), t)$ gives the actual value of $x_{i+1}(t)$. We will prove the continuity of $F_{i+1}(g(x_i(i), y_i), t)$ as follows. First, note that the map from (y_1, \dots, y_T) to $F_{i+1}(g(x_i(i), y_i), t)$ is the result of the composition:

$$\begin{pmatrix} y_1 \\ \vdots \\ y_T \end{pmatrix} \mapsto \begin{pmatrix} x_i(i) \\ y_i \end{pmatrix} \mapsto g(x_i(i), y_i) \mapsto F_{i+1}(g(x_i(i), y_i), t). \quad (88)$$

The first map is continuous since $x_i(i)$ is continuous in (y_1, \dots, y_T) by the induction hypothesis. The second map is also continuous by Assumption (2b). Lastly, the Corollary 7 shows that $F_{i+1}(\xi_{i+1}, t)$ is (uniformly) continuous in ξ_{i+1} for $t \in [i, i+1]$. Therefore, $(y_1, \dots, y_T) \mapsto x_{i+1}(t)$ is continuous for all $t \in [i, i+1]$.

Proposition 9. Bounded State Trajectory. *Given a control $u \in U$ and a sequence of observations (y_1, \dots, y_T) , the system $x(t)$ starting at $x(0) = x_1(0) = x_0$ has the following bound:*

$$\begin{aligned} \forall i \in \{2, \dots, T\} \quad \forall t \in [i-1, i] \\ \|x_i(t)\|_2 \leq \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{K}_i} \alpha_i^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m}, \end{aligned} \quad (89)$$

where \mathcal{K}_i is a finite set of sequences of non-negative integers of length $i - 1$, and $\alpha_i^{(j_1, \dots, j_{i-1})}(x_0)$ is a finite positive constant that depends on x_0 and (j_1, \dots, j_{i-1}) but not on any of the observations or the control.

For $i = 1$ the bound is given by $\forall t \in [0, 1] \quad \|x_1(t)\|_2 \leq \alpha_1(x_0)$ for some finite positive constant $\alpha_1(x_0)$.

Proof. For $i = 1$, Lemma 5 gives $\forall t \in [0, 1] \quad \|x_1(t)\|_2 \leq (1 + \|x_0\|_2) e^{K_9} \triangleq \alpha_1(x_0)$. For $i = 2$, by Assumption (2d) and the case for $i = 1$, we have

$$\|x_2(1)\|_2 = \|g(x_1(1), y_1)\|_2 \quad (90)$$

$$\leq K_3 + K_4 \|x_1(1)\|_2^{L_1} + K_5 \|y_1\|_2^{L_2} + K_6 \|x_1(1)\|_2^{L_1} \|y_1\|_2^{L_2} \quad (91)$$

$$\leq K_3 + K_4 \alpha_1(x_0)^{L_1} + K_5 \|y_1\|_2^{L_2} + K_6 \alpha_1(x_0)^{L_1} \|y_1\|_2^{L_2}. \quad (92)$$

Then, by Lemma 5, $\forall t \in [1, 2]$

$$\|x_2(t)\|_2 \leq (1 + \|x_2(1)\|_2) e^{K_9} \quad (93)$$

$$\leq e^{K_9} \left(1 + K_3 + K_4 \alpha_1(x_0)^{L_1} + K_5 \|y_1\|_2^{L_2} + K_6 \alpha_1(x_0)^{L_1} \|y_1\|_2^{L_2} \right) \quad (94)$$

$$\triangleq \sum_{(j_1) \in \mathcal{K}_2} \alpha_2^{(j_1)}(x_0) \|y_1\|_2^{j_1}, \quad (95)$$

where $\mathcal{K}_2 = \{(0), (L_2)\}$, and

$$\alpha_2^{(0)}(x_0) = e^{K_9} (1 + K_3 + K_4 \alpha_1(x_0)^{L_1}) \quad (96)$$

$$\alpha_2^{(L_2)}(x_0) = e^{K_9} (K_5 + K_6 \alpha_1(x_0)^{L_1}) \quad (97)$$

are both finite positive constants that depend on x_0 but not on any of the observations or the control.

Next, suppose that the claim holds for some $i \geq 2$. That is,

$$\forall t \in [i - 1, i]$$

$$\|x_i(t)\|_2 \leq \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{K}_i} \alpha_i^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m}, \quad (98)$$

where \mathcal{K}_i and $\alpha_i^{(j_1, \dots, j_{i-1})}(x_0)$ are as defined in the statement of the proposition. Making use of this assumption, Assumption (2d) and Lemma 5, we find that for all $t \in [i, i + 1]$,

$$\|x_{i+1}(t)\|_2 \leq (1 + \|g(x_i(i), y_i)\|_2) e^{K_9} \quad (99)$$

$$\leq e^{K_9} (1 + K_3) + e^{K_9} K_5 \|y_i\|_2^{L_2} + e^{K_9} \|x_i(i)\|_2^{L_1} (K_4 + K_6 \|y_i\|_2^{L_2}). \quad (100)$$

The first two terms in the above sum can be rewritten as

$$e^{K_9} (1 + K_3) = e^{K_9} (1 + K_3) \|y_1\|_2^0 \times \dots \times \|y_i\|_2^0 \quad (101)$$

$$e^{K_9} K_5 \|y_i\|_2^{L_2} = e^{K_9} K_5 \|y_1\|_2^0 \times \dots \times \|y_{i-1}\|_2^0 \times \|y_i\|_2^{L_2} \quad (102)$$

For the last term, we can use (98) and the multinomial theorem to write

$$\begin{aligned} & \|x_i(i)\|_2^{L_1} \\ & \leq \left(\sum_{(j_1, \dots, j_{i-1}) \in \mathcal{K}_i} \alpha_i^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m} \right)^{L_1} \\ & = \sum_{k_1 + \dots + k_{|\mathcal{K}_i|} = L_1} \binom{L_1}{k_1, \dots, k_{|\mathcal{K}_i|}} \\ & \quad \times \left\{ \prod_{l=1}^{|\mathcal{K}_i|} \alpha_i^{(j_1^{(l)}, \dots, j_{i-1}^{(l)})}(x_0)^{k_l} \right\} \times \prod_{m=1}^{i-1} \|y_m\|_2^{\sum_{l=1}^{|\mathcal{K}_i|} k_l j_m^{(l)}}, \end{aligned} \quad (103)$$

$$\begin{aligned} & = \sum_{k_1 + \dots + k_{|\mathcal{K}_i|} = L_1} \binom{L_1}{k_1, \dots, k_{|\mathcal{K}_i|}} \\ & \quad \times \left\{ \prod_{l=1}^{|\mathcal{K}_i|} \alpha_i^{(j_1^{(l)}, \dots, j_{i-1}^{(l)})}(x_0)^{k_l} \right\} \times \prod_{m=1}^{i-1} \|y_m\|_2^{\sum_{l=1}^{|\mathcal{K}_i|} k_l j_m^{(l)}}, \end{aligned} \quad (104)$$

where $(j_1^{(l)}, \dots, j_{i-1}^{(l)})$ is the l -th element in \mathcal{K}_i . Note that k_l is non-negative for all $l \in \{1, \dots, |\mathcal{K}_i|\}$. By the induction hypothesis (98), exponent $\sum_{l=1}^{|\mathcal{K}_i|} k_l j_m^{(l)}$ is also non-negative for all $m \in \{1, \dots, i - 1\}$.

Thus, substituting (101), (102), and (104) into (100), rearranging the sums and re-labeling the sequences of integer exponents, we can write

$$\|x_{i+1}(t)\|_2 \leq \sum_{(j_1, \dots, j_i) \in \mathcal{K}_{i+1}} \alpha_{i+1}^{(j_1, \dots, j_i)}(x_0) \prod_{m=1}^i \|y_m\|_2^{j_m} \quad (105)$$

for all $t \in [i, i + 1]$, where \mathcal{K}_{i+1} is a set of sequences of non-negative integers of length i , and each $\alpha_{i+1}^{(j_1, \dots, j_i)}(x_0)$ does not depend on any of the observations or the control. Here the cardinality of \mathcal{K}_{i+1} is at most finite, since

$$|\mathcal{K}_{i+1}| \leq 2 + 2 \binom{L_1 + |\mathcal{K}_i| - 1}{|\mathcal{K}_i| - 1} \quad (106)$$

by (100) and the multinomial theorem.

Finally, proceeding by mathematical induction over $i \in \{2, \dots, T\}$ completes the proof.

Proposition 10. Bounded Cost Functions. *Given a control $u \in U$ and a sequence of observations (y_1, \dots, y_T) , the instantaneous cost $c(x(t), u(t))$ induced by the state trajectory $x(t)$ has the following bound.*

$$\begin{aligned} & \forall i \in \{2, \dots, T\} \quad \forall t \in [i - 1, i] \quad |c(x_i(t), u(t))| \\ & \leq \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{K}'_i} \alpha'_i{}^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m}, \end{aligned} \quad (107)$$

where \mathcal{K}'_i is a finite set of sequences of non-negative integers of length $i - 1$, and $\alpha'_i{}^{(j_1, \dots, j_{i-1})}(x_0)$ is a finite positive constant that depends on x_0 and (j_1, \dots, j_{i-1}) but not on any of the observations or the control.

For $i = 1$ the bound is given by

$$\forall t \in [0, 1] \quad |c(x_1(t), u(t))| \leq \alpha'_1(x_0) \quad (108)$$

for some finite positive constant $\alpha'_1(x_0)$.

Similarly, the terminal cost $h(x(T))$ is bounded by

$$|h(x(T))| \leq \sum_{(j_1, \dots, j_T) \in \mathcal{K}'_{T+1}} \alpha'_{T+1}{}^{(j_1, \dots, j_T)}(x_0) \prod_{m=1}^T \|y_m\|_2^{j_m}. \quad (109)$$

Proof. For the instantaneous cost function $c(x_i(t), u(t))$, Assumption 3 along with Proposition 9 yields the following bound:

$$\forall i \in \{2, \dots, T\} \forall t \in [i-1, i] \quad |c(x_i(t), u(t))| \leq K_7 + K_8 \left(\sum_{(j_1, \dots, j_{i-1}) \in \mathcal{K}_i} \alpha_i^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m} \right)^{L_3}. \quad (110)$$

Making use of the multinomial expansion formula in the same manner as in the proof of Proposition 9, we conclude that

$$\forall i \in \{2, \dots, T\} \forall t \in [i-1, i] \quad |c(x_i(t), u(t))| \leq \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{K}'_i} \alpha_i'^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m} \quad (111)$$

for some finite set \mathcal{K}'_i of sequences of non-negative integers and finite positive constants $\alpha_i'^{(j_1, \dots, j_{i-1})}(x_0)$. Similarly, for $i = 1$ we obtain

$$|c(x_1(t), u(t))| \leq K_7 + K_8 \alpha_1(x_0)^{L_3} \triangleq \alpha_1'(x_0) \quad (112)$$

for all $t \in [0, 1]$.

To bound the terminal cost, note that

$$|h(x(T))| \leq K_7 + K_8 \|x(T)\|_2^{L_3} \quad (113)$$

$$= K_7 + K_8 \|g(x_T(T), y_T)\|_2^{L_3} \quad (114)$$

$$\leq K_7 + K_8 \left\{ K_3 + K_5 \|y_T\|_2^{L_2} + \|x_T(T)\|_2^{L_1} \left(K_4 + K_6 \|y_T\|_2^{L_2} \right) \right\}^{L_3} \quad (115)$$

by Assumptions (2d) and 3. Since $L_3 < \infty$, (115) yields a polynomial of $\|x_T(T)\|_2$ and $\|y_T\|_2$ of finite terms, for each of which we can use Proposition 9 and apply the multinomial expansion formula to show

$$|h(x(T))| \leq \sum_{(j_1, \dots, j_T) \in \mathcal{K}'_{T+1}} \alpha_{T+1}'^{(j_1, \dots, j_T)}(x_0) \prod_{m=1}^T \|y_m\|_2^{j_m} \quad (116)$$

for some finite set \mathcal{K}'_{T+1} of sequence of non-negative integers and finite positive constants $\alpha_{T+1}'^{(j_1, \dots, j_T)}(x_0)$.

A.2 Perturbed Trajectory under Specific Observations

Next, we will perturb the nominal state x of the system while assuming the same initial condition x_0 and the specific observations (y_1, \dots, y_T) as in Section A.1. The perturbed control is an open-loop perturbation as defined below.

Definition 1. Perturbed Control. Let $u \in U$ be a control. For $\tau \in (0, 1)$ and $v \in B(0, \rho_{\max})$, define the perturbed control u^ϵ by

$$u^\epsilon(t) \triangleq \begin{cases} v & \text{if } t \in (\tau - \epsilon, \tau] \\ u(t) & \text{otherwise,} \end{cases} \quad (117)$$

where $\epsilon \in [0, \tau]$. By definition if $\epsilon = 0$ then u^ϵ is the same as u . We assume that the nominal control $u(t)$ is left continuous in t at $t = \tau$.

Remark 4. It is obvious that $u^\epsilon(t)$ is piecewise continuous on $[0, T]$. Therefore, for $v \in B(0, \rho_{\max})$ we have $u^\epsilon \in U$, i.e. u^ϵ is an admissible control. Thus, Proposition 3 assures that there exists a unique solution x^ϵ for the trajectory of the system under the control perturbation. In the remainder of the analysis, we assume that (τ, v) is given and fixed.

Lemma 11. Let $\epsilon, \epsilon' \in [0, \tau]$, and let $u^\epsilon, u^{\epsilon'} \in U$ be two perturbed controls of the form (117). Let $x_1^\epsilon, x_1^{\epsilon'}$ be the solutions of x_1 for $t \in [0, 1]$ by applying u^ϵ and $u^{\epsilon'}$ respectively to the initial condition x_0 . Then, there exists an $L' < \infty$, independent of $\epsilon, \epsilon', x_1^\epsilon$ and $x_1^{\epsilon'}$, such that

$$\forall \epsilon, \epsilon' \in [0, \tau] \forall t \in [0, 1] \quad \|x_1^\epsilon(t) - x_1^{\epsilon'}(t)\|_2 \leq L' |\epsilon - \epsilon'|. \quad (118)$$

Proof. By Proposition 6, we find that

$$\forall t \in [0, 1] \quad \|x_1^\epsilon(t) - x_1^{\epsilon'}(t)\|_2 \leq L \int_0^1 \|u^\epsilon(t) - u^{\epsilon'}(t)\|_2 dt \quad (119)$$

for some $L < \infty$. Let us derive an upper-bound on the integral on the right hand side. If $\epsilon \geq \epsilon'$,

$$\int_0^1 \|u^\epsilon(t) - u^{\epsilon'}(t)\|_2 dt = \int_{\tau-\epsilon}^{\tau-\epsilon'} \|v - u(t)\|_2 dt, \quad (120)$$

where $u(t)$ is the nominal control that both u^ϵ and $u^{\epsilon'}$ are based on. Since $u(t) \in B(0, \rho_{\max})$, we obtain

$$\int_0^1 \|u^\epsilon(t) - u^{\epsilon'}(t)\|_2 dt \leq \sup_{u \in B(0, \rho_{\max})} \|v - u\|_2 (\epsilon - \epsilon'). \quad (121)$$

Similarly, if $\epsilon < \epsilon'$ we have

$$\int_0^1 \|u^\epsilon(t) - u^{\epsilon'}(t)\|_2 dt \leq \sup_{u \in B(0, \rho_{\max})} \|v - u\|_2 (\epsilon' - \epsilon). \quad (122)$$

Put these two cases together and substitute into (119) to get

$$\forall t \in [0, 1] \quad \|x_1^\epsilon(t) - x_1^{\epsilon'}(t)\|_2 \leq L' |\epsilon - \epsilon'|, \quad (123)$$

where $L' \triangleq L \sup_{u \in B(0, \rho_{\max})} \|v - u\|_2 \leq 2L\rho_{\max} < \infty$.

Lemma 12. Let u^ϵ and x_1^ϵ be as in Lemma 11. Then,

$$\lim_{\epsilon \rightarrow 0^+} \frac{1}{\epsilon} \int_{\tau-\epsilon}^{\tau} \{f(x_1^\epsilon(t), u^\epsilon(t)) - f(x_1(t), u(t))\} dt = f(x_1(\tau), v) - f(x_1(\tau), u(\tau)), \quad (124)$$

where x_1 denotes the solution under the nominal control $u \in U$.

Proof. We will show that the difference norm

$$\left\| \frac{1}{\epsilon} \int_{\tau-\epsilon}^{\tau} \{f(x_1^\epsilon(t), u^\epsilon(t)) - f(x_1(t), u(t))\} dt - f(x_1(\tau), v) + f(x_1(\tau), u(\tau)) \right\|_2 \quad (125)$$

converges to 0 as $\epsilon \rightarrow 0^+$. Indeed, (125) becomes

$$\left\| \frac{1}{\epsilon} \int_{\tau-\epsilon}^{\tau} \{f(x_1^\epsilon(t), u^\epsilon(t)) - f(x_1(t), u(t)) - f(x_1(\tau), v) + f(x_1(\tau), u(\tau))\} dt \right\|_2 \quad (126)$$

$$\leq \frac{1}{\epsilon} \int_{\tau-\epsilon}^{\tau} \|f(x_1^\epsilon(t), u^\epsilon(t)) - f(x_1(t), u(t)) - f(x_1(\tau), v) + f(x_1(\tau), u(\tau))\|_2 dt \quad (127)$$

$$\leq \frac{1}{\epsilon} \int_{\tau-\epsilon}^{\tau} \{\|f(x_1^\epsilon(t), u^\epsilon(t)) - f(x_1(\tau), v)\|_2 + \|f(x_1(t), u(t)) - f(x_1(\tau), u(\tau))\|_2\} dt. \quad (128)$$

We used the triangle inequality in (128). Now, making use of the fact that $\forall t \in (\tau - \epsilon, \tau]$ $u^\epsilon(t) = v$, Assumption (2c) yields

$$\|f(x_1^\epsilon(t), u^\epsilon(t)) - f(x_1(\tau), v)\|_2 \leq K_1 \|x_1^\epsilon(t) - x_1(\tau)\|_2 \quad (129)$$

$$= K_1 \|x_1^\epsilon(t) - x_1(t) + x_1(t) - x_1(\tau)\|_2 \quad (130)$$

$$\leq K_1 \|x_1^\epsilon(t) - x_1(t)\|_2 + K_1 \|x_1(t) - x_1(\tau)\|_2 \quad (131)$$

$$\leq K_1 L' \epsilon + K_1 \|x_1(t) - x_1(\tau)\|_2 \quad (132)$$

for any $t \in (\tau - \epsilon, \tau]$, where we applied Lemma 11 in (132) with $\epsilon' = 0$. Similarly,

$$\|f(x_1(t), u(t)) - f(x_1(\tau), u(\tau))\|_2 \leq K_1 \|x_1(t) - x_1(\tau)\|_2 + K_1 \|u(t) - u(\tau)\|_2. \quad (133)$$

Therefore, (125) is upper-bounded by

$$\frac{1}{\epsilon} \int_{\tau-\epsilon}^{\tau} \{K_1 L' \epsilon + 2K_1 \|x_1(t) - x_1(\tau)\|_2 + K_1 \|u(t) - u(\tau)\|_2\} dt \quad (134)$$

$$\leq K_1 L' \epsilon + 2K_1 \sup_{t \in [\tau-\epsilon, \tau]} \|x_1(t) - x_1(\tau)\|_2 + K_1 \sup_{t \in [\tau-\epsilon, \tau]} \|u(t) - u(\tau)\|_2, \quad (135)$$

which converges to 0 as $\epsilon \rightarrow 0^+$, since

$$0 \leq \sup_{t \in [\tau-\epsilon, \tau]} \|x_1(t) - x_1(\tau)\|_2 \rightarrow 0 \quad (136)$$

and

$$0 \leq \sup_{t \in [\tau-\epsilon, \tau]} \|u(t) - u(\tau)\|_2 \rightarrow 0 \quad (137)$$

as $\epsilon \rightarrow 0^+$.

Lemma 13. Let x_1^ϵ and x_1 be as in Lemma 12. Let $\Psi_1(t)$ be the right derivative of $x_1^\epsilon(t)$ with respect to ϵ evaluated at $\epsilon = 0$. That is,

$$\Psi_1(\tau) = \frac{\partial_+ x_1^\epsilon(\tau)}{\partial \epsilon} \Big|_{\epsilon=0} \triangleq \lim_{\epsilon \rightarrow 0^+} \frac{x_1^\epsilon(\tau) - x_1(\tau)}{\epsilon}. \quad (138)$$

$$(139)$$

Then, we have

$$\Psi_1(\tau) = f(x_1(\tau), v) - f(x_1(\tau), u(\tau)). \quad (140)$$

Proof. Let us express both $x_1^\epsilon(\tau)$ and $x_1(\tau)$ in the integral form:

$$x_1^\epsilon(\tau) = x_0 + \int_0^\tau f(x_1^\epsilon(t), u^\epsilon(t)) dt \quad (141)$$

$$x_1(\tau) = x_0 + \int_0^\tau f(x_1(t), u(t)) dt. \quad (142)$$

Note that $u^\epsilon(t) = u(t)$ and $x_1^\epsilon(t) = x_1(t)$ for $t \in [0, \tau - \epsilon]$, since no perturbation is applied to the system until $t > \tau - \epsilon$. Therefore,

$$\begin{aligned} x_1^\epsilon(\tau) - x_1(\tau) &= \int_{\tau-\epsilon}^\tau \{f(x_1^\epsilon(t), u^\epsilon(t)) - f(x_1(t), u(t))\} dt. \end{aligned} \quad (143)$$

Making use of Lemma 12, we conclude that

$$\lim_{\epsilon \rightarrow 0^+} \frac{x_1^\epsilon(\tau) - x_1(\tau)}{\epsilon} = f(x_1(\tau), v) - f(x_1(\tau), u(\tau)). \quad (144)$$

Lemma 14. Let x_1^ϵ and x_1 be as in Lemma 12. Then, $\Psi_1(t) = \frac{\partial_+ x_1^\epsilon(t)}{\partial \epsilon} \Big|_{\epsilon=0}$ uniquely exists for $t \in [\tau, 1]$ and follows the ODE:

$$\dot{\Psi}_1(t) = \frac{\partial}{\partial x_1} f(x_1(t), u(t)) \Psi_1(t), \quad (145)$$

with the initial condition $\Psi_1(\tau)$ given by Lemma 13.

Proof. Taking some $a \in (\tau, 1]$, let us express $x_1^\epsilon(a)$ and $x_1(a)$ in the integral form:

$$x_1^\epsilon(a) = x_1^\epsilon(\tau) + \int_\tau^a f(x_1^\epsilon(t), u^\epsilon(t)) dt \quad (146)$$

$$x_1(a) = x_1(\tau) + \int_\tau^a f(x_1(t), u(t)) dt. \quad (147)$$

Thus, we have

$$\begin{aligned} \Psi_1(a) &= \Psi_1(\tau) \\ &+ \lim_{\epsilon \rightarrow 0^+} \int_\tau^a \frac{1}{\epsilon} \{f(x_1^\epsilon(t), u(t)) - f(x_1(t), u(t))\} dt, \end{aligned} \quad (148)$$

where we used $\forall t \in [\tau, a]$ $u^\epsilon(t) = u(t)$. We will take a measure-theoretic approach to prove that the order of the limit and the integration can be switched in (148). First, think of the integral as a Lebesgue integral on the measure space $([\tau, a], \mathcal{B}([\tau, a]), \lambda)$, where $\mathcal{B}([\tau, a])$ is the Borel σ -algebra on $[\tau, a]$ and λ is the Lebesgue measure. Furthermore, consider the integrand as a function from $[\tau, a]$ into the Banach space $(\mathbb{R}^{n_x}, \|\cdot\|_2)$, i.e. the Euclidean space endowed with the ℓ^2 norm. By the piecewise continuity of u^ϵ, u and the continuity of x_1^ϵ, x_1 , and f , the integrand is a piecewise continuous function with respect to t , which is λ -measurable. In fact it is also Bochner-integrable, since for $t \in [\tau, a]$ we have the constant bound:

$$\frac{1}{\epsilon} \|f(x_1^\epsilon(t), u(t)) - f(x_1(t), u(t))\|_2 \quad (149)$$

$$\leq \frac{1}{\epsilon} K_1 \|x_1^\epsilon(t) - x_1(t)\|_2 \quad (150)$$

$$\leq K_1 L' \quad (151)$$

by Assumption (2c) and Lemma 11. Furthermore, the chain rule gives

$$\lim_{\epsilon \rightarrow 0^+} \frac{1}{\epsilon} (f(x_1^\epsilon(t), u(t)) - f(x_1(t), u(t))) \quad (152)$$

$$= \frac{\partial}{\partial x_1} f(x_1(t), u(t)) \Psi_1(t), \quad (153)$$

assuming that $\Psi_1(t)$ exists for $t \in [\tau, a]$. Therefore, by the Bochner-integral version of the dominated convergence theorem (Theorem 3 in Diestel and Uhl (1977), Chapter II), we obtain

$$\Psi_1(a) = \Psi_1(t) + \int_\tau^a \frac{\partial}{\partial x_1} f(x_1(t), u(t)) \Psi_1(t) dt. \quad (154)$$

This is equivalent to the ordinary differential equation:

$$\dot{\Psi}_1(t) = \frac{\partial}{\partial x_1} f(x_1(t), u(t)) \Psi_1(t). \quad (155)$$

It remains to show that the solution $\Psi_1(t)$ that satisfies (155) does exist and is unique. First, let Ψ'_1 and Ψ''_1 be two systems that follow (155) and share the same initial condition $\Psi_1(\tau)$. Then, by Assumption (2c) we have

$$\begin{aligned} \|\dot{\Psi}'_1(t) - \dot{\Psi}''_1(t)\|_2 &= \left\| \frac{\partial}{\partial x_1} f(x_1(t), u(t)) (\Psi'_1(t) - \Psi''_1(t)) \right\|_2 \\ &\leq \left\| \frac{\partial}{\partial x_1} f(x_1(t), u(t)) \right\|_2 \cdot \|\Psi'_1(t) - \Psi''_1(t)\|_2 \end{aligned} \quad (156)$$

$$(157)$$

$$\leq K_2 \|\Psi'_1(t) - \Psi''_1(t)\|_2. \quad (158)$$

Existence follows from this inequality in conjunction with the Picard Lemma (Lemma 5.6.3 in Elijah (1997)). To show the uniqueness, apply the Bellman-Gronwall Lemma (Lemma 5.6.4 in Elijah (1997)) to the following integral inequality:

$$\forall a \in [\tau, 1]$$

$$\|\Psi'_1(a) - \Psi''_1(a)\|_2 \leq K_2 \int_\tau^a \|\Psi'_1(t) - \Psi''_1(t)\|_2 dt. \quad (159)$$

Proposition 15. Variational Equation. *Let $u \in U$ and x_1 be the nominal control and the resulting state trajectory. Let x_1^ϵ be the perturbed state induced by the perturbed control u^ϵ of the form (117). Propagating $x_1^\epsilon(t)$ through the hybrid dynamics, we get a series of modes $x_2^\epsilon, \dots, x_T^\epsilon$ that constitutes the entire trajectory $x^\epsilon(t)$ for $t \in [\tau, T]$. Define the state variation $\Psi(t)$ for $t \in [\tau, T]$ by*

$$\Psi(t) = \frac{\partial_+}{\partial \epsilon} x^\epsilon(t) \Big|_{\epsilon=0} \triangleq \lim_{\epsilon \rightarrow 0^+} \frac{x^\epsilon(t) - x(t)}{\epsilon}. \quad (160)$$

Then, $\Psi(t)$ exists for $t \in [\tau, T]$ and follows the hybrid system with time-driven switching:

$$\Psi(t) = \begin{cases} \Psi_1(t) & \forall t \in [\tau, 1) \\ \Psi_i(t) & \forall t \in [i-1, i) \forall i \in \{2, \dots, T\} \end{cases} \quad (161)$$

with $\Psi(T) = \frac{\partial}{\partial x_T} g(x_T(T), y_T) \Psi_T(T)$, where Ψ_1 is defined on $[\tau, 1]$ as in Lemma 14, and Ψ_i for $i \geq 2$ is defined on

$[i-1, i]$ with

$$\Psi_i(i-1) = \frac{\partial}{\partial x_{i-1}} g(x_{i-1}(i-1), y_{i-1}) \Psi_{i-1}(i-1) \quad (162)$$

$$\dot{\Psi}_i(t) = \frac{\partial}{\partial x_i} f(x_i(t), u(t)) \Psi_i(t) \quad \forall t \in [i-1, i]. \quad (163)$$

Proof. The case for $t \in [\tau, 1)$ follows from Lemma 14, since in this case we have

$$\Psi(t) \triangleq \lim_{\epsilon \rightarrow 0^+} \frac{x^\epsilon(t) - x(t)}{\epsilon} \quad (164)$$

$$= \lim_{\epsilon \rightarrow 0^+} \frac{x_1^\epsilon(t) - x_1(t)}{\epsilon} \quad (165)$$

$$= \Psi_1(t). \quad (166)$$

At $t = 1$, we obtain

$$\Psi(1) \triangleq \lim_{\epsilon \rightarrow 0^+} \frac{x^\epsilon(1) - x(1)}{\epsilon} \quad (167)$$

$$= \lim_{\epsilon \rightarrow 0^+} \frac{x_2^\epsilon(1) - x_2(1)}{\epsilon} \quad (168)$$

$$= \lim_{\epsilon \rightarrow 0^+} \frac{g(x_1^\epsilon(1), y_1) - g(x_1(1), y_1)}{\epsilon} \quad (169)$$

$$= \frac{\partial}{\partial x} g(x_1(1), y_1) \Psi_1(1) \quad (170)$$

by (47) and the chain rule. Let us define Ψ_2 by $\Psi_2(t) \triangleq \lim_{\epsilon \rightarrow 0^+} (x_2^\epsilon(t) - x_2(t))/\epsilon$. Similarly to the proof of Lemma 14, one can show that Ψ_2 follows the integral equation:

$$\Psi_2(a) = \Psi_2(1) + \int_1^a \frac{\partial}{\partial x_2} f(x_2(t), u(t)) \Psi_2(t) dt \quad (171)$$

with $\Psi_2(1) = \Psi(1)$, and that $\Psi_2(t)$ that satisfies (162) and (163) uniquely exists. This proves the case for $t \in [1, 2)$. Proceeding by induction completes the proof.

So far we have focused entirely on the right derivative $\frac{\partial_+}{\partial \epsilon} x^\epsilon(t)$ evaluated at $\epsilon = 0$. The next proposition shows that x_1^ϵ is in fact right differentiable with respect to ϵ at all $\epsilon \in [0, \tau)$.

Proposition 16. Right Differentiability of State Perturbation. *Let $x^\epsilon(t)$ be the perturbed state trajectory under the perturbed control u^ϵ defined by (117). Let $\Psi^\epsilon(t)$ denote the right derivative $\frac{\partial_+}{\partial \epsilon} x^\epsilon(t)$ evaluated at a particular $\epsilon \in [0, \tau)$. (Note that when $\epsilon = 0$ we have $\Psi^\epsilon = \Psi$.) Then, $\Psi^\epsilon(t)$ exists for $t \in [\tau - \epsilon, T]$ and follows the hybrid system with time-driven switching:*

$$\Psi^\epsilon(t) = \begin{cases} \Psi_1^\epsilon(t) & \forall t \in [\tau - \epsilon, 1) \\ \Psi_i^\epsilon(t) & \forall t \in [i-1, i) \forall i \in \{2, \dots, T\} \end{cases} \quad (172)$$

with $\Psi^\epsilon(T) = \frac{\partial}{\partial x_T} g(x_T^\epsilon(T), y_T) \Psi_T^\epsilon(T)$. Ψ_1^ϵ is defined on $[\tau - \epsilon, 1]$, where

$$\Psi_1^\epsilon(\tau - \epsilon) = f(x_1^\epsilon(\tau - \epsilon), v) - f(x_1^\epsilon(\tau - \epsilon), u^\epsilon(\tau - \epsilon)) \quad (173)$$

and

$$\dot{\Psi}_1^\epsilon(t) = \frac{\partial}{\partial x_1^\epsilon} f(x_1^\epsilon(t), u^\epsilon(t)) \Psi_1^\epsilon(t) \quad \forall t \in [\tau - \epsilon, 1]. \quad (174)$$

Ψ_i^ϵ for $i \geq 2$ is defined on $[i - 1, i]$ with

$$\Psi_i^\epsilon(i - 1) = \frac{\partial}{\partial x_{i-1}^\epsilon} g(x_{i-1}^\epsilon(i - 1), y_{i-1}) \Psi_{i-1}^\epsilon(i - 1) \quad (175)$$

$$\dot{\Psi}_i^\epsilon(t) = \frac{\partial}{\partial x_i^\epsilon} f(x_i^\epsilon(t), u^\epsilon(t)) \Psi_i^\epsilon(t) \quad \forall t \in [i - 1, i]. \quad (176)$$

Proof. The proposition can be proven by considering the perturbed control u^ϵ as the new nominal control, and defining a further perturbation based on this nominal control. Formally, define $\tilde{u}^{\epsilon'}$ by

$$\tilde{u}^{\epsilon'}(t) \triangleq \begin{cases} v & \text{if } t \in (\tau - \epsilon - \epsilon', \tau - \epsilon] \\ u^\epsilon(t) & \text{otherwise} \end{cases} \quad (177)$$

with $\epsilon' \in [0, \tau - \epsilon]$, where (τ, v) is the same pair of values as for $u^\epsilon(t)$. Since $u^\epsilon(t)$ is left continuous in t at $t = \tau - \epsilon$, $\tilde{u}^{\epsilon'}$ is a valid perturbed control of the form (117) with u^ϵ being the nominal control. Here ϵ is considered as fixed and the parameters defining this new perturbation are v , $\tau - \epsilon$, and ϵ' . This perturbation yields the new perturbed state trajectory $\tilde{x}^{\epsilon'}$ based on the nominal trajectory x^ϵ . Note that when $\epsilon' = 0$ we have $\tilde{u}^{\epsilon'} = u^\epsilon$ and $\tilde{x}^{\epsilon'} = x^\epsilon$. We can define the new state variation:

$$\tilde{\Psi}(t) = \frac{\partial}{\partial \epsilon'} \tilde{x}^{\epsilon'}(t) \Big|_{\epsilon'=0} \triangleq \lim_{\epsilon' \rightarrow 0^+} \frac{\tilde{x}^{\epsilon'}(t) - x^\epsilon(t)}{\epsilon'}. \quad (178)$$

Applying Proposition 15 to this new setting, we find that $\tilde{\Psi}(t)$ exists for $t \in [\tau - \epsilon, T]$ and follows the hybrid system with time-driven switching:

$$\tilde{\Psi}(t) = \begin{cases} \tilde{\Psi}_1(t) & \forall t \in [\tau - \epsilon, 1] \\ \tilde{\Psi}_i(t) & \forall t \in [i - 1, i] \quad \forall i \in \{2, \dots, T\} \end{cases} \quad (179)$$

with $\tilde{\Psi}(T) = \frac{\partial}{\partial x_T^\epsilon} g(x_T^\epsilon(T), y_T) \tilde{\Psi}_T(T)$. $\tilde{\Psi}_1$ is defined on $[\tau - \epsilon, 1]$, where

$$\tilde{\Psi}_1(\tau - \epsilon) = f(x_1^\epsilon(\tau - \epsilon), v) - f(x_1^\epsilon(\tau - \epsilon), u^\epsilon(\tau - \epsilon)) \quad (180)$$

and

$$\dot{\tilde{\Psi}}_1(t) = \frac{\partial}{\partial x_1^\epsilon} f(x_1^\epsilon(t), u^\epsilon(t)) \tilde{\Psi}_1(t) \quad \forall t \in [\tau - \epsilon, 1]. \quad (181)$$

$\tilde{\Psi}_i$ for $i \geq 2$ is defined on $[i - 1, i]$ with

$$\tilde{\Psi}_i(i - 1) = \frac{\partial}{\partial x_{i-1}^\epsilon} g(x_{i-1}^\epsilon(i - 1), y_{i-1}) \tilde{\Psi}_{i-1}(i - 1) \quad (182)$$

$$\dot{\tilde{\Psi}}_i(t) = \frac{\partial}{\partial x_i^\epsilon} f(x_i^\epsilon(t), u^\epsilon(t)) \tilde{\Psi}_i(t) \quad \forall t \in [i - 1, i]. \quad (183)$$

On the other hand, notice that the new perturbed control $\tilde{u}^{\epsilon'}$ is actually equivalent to the perturbed control $u^{\epsilon+\epsilon'}$ that

is based on the original nominal control u . Namely,

$$\tilde{u}^{\epsilon'}(t) = u^{\epsilon+\epsilon'}(t) \triangleq \begin{cases} v & \text{if } t \in (\tau - (\epsilon + \epsilon'), \tau] \\ u(t) & \text{otherwise.} \end{cases} \quad (184)$$

Consequently, the new perturbed state $\tilde{x}^{\epsilon'}$ is equal to $x^{\epsilon+\epsilon'}$, and thus

$$\tilde{\Psi}(t) = \lim_{\epsilon' \rightarrow 0^+} \frac{\tilde{x}^{\epsilon'}(t) - x^\epsilon(t)}{\epsilon'} \quad (185)$$

$$= \lim_{\epsilon' \rightarrow 0^+} \frac{x^{\epsilon+\epsilon'}(t) - x^\epsilon(t)}{\epsilon'} \quad (186)$$

$$= \Psi^\epsilon(t). \quad (187)$$

This completes the proof.

Lemma 17. Let $\Psi_1^\epsilon, \dots, \Psi_T^\epsilon$ be as given by Proposition 16. Then, for Ψ_1^ϵ the following holds.

$$\|\Psi_1^\epsilon(t)\|_2 \leq 2K_1 \rho_{\max} e^{K_2} \quad \forall t \in [\tau - \epsilon, 1] \quad (188)$$

Similarly, for all $i \in \{2, \dots, T\}$ we have

$$\|\Psi_i^\epsilon(t)\|_2 \leq \|\Psi_i^\epsilon(i - 1)\|_2 e^{K_2} \quad \forall t \in [i - 1, i]. \quad (189)$$

Proof. We begin with the integral equation:

$$\begin{aligned} \Psi_1^\epsilon(a) &= \Psi_1^\epsilon(\tau - \epsilon) \\ &\quad + \int_{\tau - \epsilon}^a \frac{\partial}{\partial x_1^\epsilon} f(x_1^\epsilon(t), u^\epsilon(t)) \Psi_1^\epsilon(t) dt. \end{aligned} \quad (190)$$

Therefore,

$$\begin{aligned} \|\Psi_1^\epsilon(a)\|_2 &\leq \|\Psi_1^\epsilon(\tau - \epsilon)\|_2 \\ &\quad + \int_{\tau - \epsilon}^a \left\| \frac{\partial}{\partial x_1^\epsilon} f(x_1^\epsilon(t), u^\epsilon(t)) \right\|_2 \cdot \|\Psi_1^\epsilon(t)\|_2 dt. \end{aligned} \quad (191)$$

Using

$$\Psi_1^\epsilon(\tau - \epsilon) = f(x_1^\epsilon(\tau - \epsilon), v) - f(x_1^\epsilon(\tau - \epsilon), u^\epsilon(\tau - \epsilon)) \quad (192)$$

and Assumption (2c), we get

$$\begin{aligned} \|\Psi_1^\epsilon(a)\|_2 &\leq K_1 \|v - u^\epsilon(\tau - \epsilon)\|_2 + K_2 \int_{\tau - \epsilon}^a \|\Psi_1^\epsilon(t)\|_2 dt \\ &\leq K_1 \sup_{u \in B(0, \rho_{\max})} \|v - u\|_2 \\ &\quad + K_2 \int_{\tau - \epsilon}^a \|\Psi_1^\epsilon(t)\|_2 dt \end{aligned} \quad (193)$$

$$\leq K_1 \sup_{u \in B(0, \rho_{\max})} \|v - u\|_2 + K_2 \int_{\tau - \epsilon}^a \|\Psi_1^\epsilon(t)\|_2 dt \quad (194)$$

$$\leq 2K_1 \rho_{\max} + K_2 \int_{\tau - \epsilon}^a \|\Psi_1^\epsilon(t)\|_2 dt \quad (195)$$

for all $a \in [\tau - \epsilon, 1]$. Thus, by the Bellman-Gronwall Lemma (Lemma 5.6.4 in Eljah (1997)) it follows that

$$\|\Psi_1^\epsilon(t)\|_2 \leq 2K_1 \rho_{\max} e^{K_2} \quad \forall t \in [\tau - \epsilon, 1]. \quad (196)$$

For general $i \geq 2$, apply the Bellman-Gronwall Lemma to the similar integral inequality:

$$\forall a \in [i-1, i] \quad (197)$$

$$\|\Psi_i^\epsilon(a)\|_2 \leq \|\Psi_i^\epsilon(i-1)\|_2 + K_2 \int_{i-1}^a \|\Psi_i^\epsilon(t)\|_2 dt \quad (198)$$

to get the result.

Proposition 18. Bounded State Variation. *Given u^ϵ and (y_1, \dots, y_T) , Ψ^ϵ defined in Proposition 16 has the following bound:*

$$\forall \epsilon \in [0, \tau) \quad \forall i \in \{2, \dots, T\} \quad \forall t \in [i-1, i] \quad (199)$$

$$\|\Psi_i^\epsilon(t)\|_2 \leq \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{L}_i} \beta_i^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m}, \quad (200)$$

where \mathcal{L}_i is a finite set of sequences of non-negative integers of length $i-1$, and $\beta_i^{(j_1, \dots, j_{i-1})}(x_0)$ is a finite positive constant that depends on x_0 and (j_1, \dots, j_{i-1}) but not on ϵ , u^ϵ , or (y_1, \dots, y_T) .

Proof. The proof of this proposition is similar to that of Proposition 9. Take any $\epsilon \in [0, \tau)$. For $i = 2$, we have $\forall t \in [1, 2]$

$$\|\Psi_2^\epsilon(t)\|_2 \leq \|\Psi_2^\epsilon(1)\|_2 e^{K_2} \quad (201)$$

$$\leq \left\| \frac{\partial}{\partial x_1^\epsilon} g(x_1^\epsilon(1), y_1) \Psi_1^\epsilon(1) \right\|_2 e^{K_2} \quad (202)$$

$$\leq \left\| \frac{\partial}{\partial x_1^\epsilon} g(x_1^\epsilon(1), y_1) \right\|_2 \cdot \|\Psi_1^\epsilon(1)\|_2 e^{K_2} \quad (203)$$

$$\leq 2K_1 \rho_{\max} e^{2K_2} \left\{ K_3 + K_5 \|y_1\|_2^{L_2} + (K_4 + K_6 \|y_1\|_2^{L_2}) \|x_1^\epsilon(1)\|_2^{L_1} \right\} \quad (204)$$

by Assumption (2d), Proposition 16, and Lemma 17. Using Proposition 9, we can bound $x_1^\epsilon(1)$ by

$$\|x_1^\epsilon(1)\|_2 \leq \sum_{(j_1) \in \mathcal{K}_2} \alpha_2^{(j_1)}(x_0) \|y_1\|_2^{j_1}. \quad (205)$$

Substituting (205) into (204) and using the multinomial theorem, one can verify that

$$\forall t \in [1, 2] \quad \|\Psi_2^\epsilon(t)\|_2 \leq \sum_{(j_1) \in \mathcal{L}_1} \beta_1^{(j_1)}(x_0) \|y_1\|_2^{j_1} \quad (206)$$

for some finite set \mathcal{L}_1 and finite $\beta_1^{(j_1)}(x_0)$.

Next, suppose that the claim holds for some $i \leq 2$. That is,

$$\forall t \in [i-1, i] \quad \|\Psi_i^\epsilon(t)\|_2 \leq \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{L}_i} \beta_i^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m}, \quad (207)$$

where \mathcal{L}_i and $\beta_i^{(j_1, \dots, j_{i-1})}(x_0)$ are as defined in the statement of the proposition. Similar to the case for $i = 2$, we have

$$\forall t \in [i, i+1]$$

$$\begin{aligned} & \|\Psi_{i+1}^\epsilon(t)\|_2 \\ & \leq e^{K_2} \left(\sum_{(j_1, \dots, j_{i-1}) \in \mathcal{L}_i} \beta_i^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m} \right) \\ & \times \left\{ K_3 + K_5 \|y_i\|_2^{L_2} + (K_4 + K_6 \|y_i\|_2^{L_2}) \|x_i^\epsilon(i)\|_2^{L_1} \right\} \end{aligned} \quad (208)$$

Proposition 9 gives the following bound:

$$\|x_i^\epsilon(i)\|_2 \leq \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{K}_i} \alpha_i^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m}. \quad (209)$$

Substituting (209) into (208) and using the multinomial theorem, we conclude that

$$\forall t \in [i, i+1] \quad \|\Psi_{i+1}^\epsilon(t)\|_2 \leq \sum_{(j_1, \dots, j_i) \in \mathcal{L}_{i+1}} \beta_{i+1}^{(j_1, \dots, j_i)}(x_0) \prod_{m=1}^i \|y_m\|_2^{j_m} \quad (210)$$

for some finite set \mathcal{L}_{i+1} and finite $\beta_{i+1}^{(j_1, \dots, j_i)}(x_0)$.

Finally, proceeding by mathematical induction over $i \in \{2, \dots, T\}$ completes the proof.

Lemma 19. *Let u^ϵ and x_1^ϵ be as in Lemma 11. Then, similarly to Lemma 12 we have*

$$\begin{aligned} \lim_{\epsilon \rightarrow 0^+} \frac{1}{\epsilon} \int_{\tau-\epsilon}^{\tau} \{c(x_1^\epsilon(t), u^\epsilon(t)) - c(x_1(t), u(t))\} dt \\ = c(x_1(\tau), v) - c(x_1(\tau), u(\tau)). \end{aligned} \quad (211)$$

Proof. As $u^\epsilon(t) = v \quad \forall t \in (\tau - \epsilon, \tau]$, we will show that

$$\begin{aligned} \lim_{\epsilon \rightarrow 0^+} \frac{1}{\epsilon} \int_{\tau-\epsilon}^{\tau} \{c(x_1^\epsilon(t), v) - c(x_1(t), u(t))\} dt \\ = c(x_1(\tau), v) - c(x_1(\tau), u(\tau)). \end{aligned} \quad (212)$$

By Assumption 3 and the continuity of $x_1^\epsilon(t)$, $c(x_1^\epsilon(t), v)$ is continuous with respect to t on $[\tau - \epsilon, \tau]$. Thus, the mean value theorem yields

$$\frac{1}{\epsilon} \int_{\tau-\epsilon}^{\tau} c(x_1^\epsilon(t), v) dt = c(x_1^\epsilon(\tilde{t}), v) \quad (213)$$

for some $\tilde{t} \in [\tau - \epsilon, \tau]$. From the triangle inequality and Lemma 11 it follows that

$$\|x_1^\epsilon(\tilde{t}) - x_1(\tau)\|_2 \leq \|x_1^\epsilon(\tilde{t}) - x_1(\tilde{t})\|_2 + \|x_1(\tilde{t}) - x_1(\tau)\|_2 \quad (214)$$

$$\leq L' \epsilon + \|x_1(\tilde{t}) - x_1(\tau)\|_2. \quad (215)$$

Therefore, $\lim_{\epsilon \rightarrow 0^+} x_1^\epsilon(\tilde{t}) = x_1(\tau)$ and

$$\lim_{\epsilon \rightarrow 0^+} \frac{1}{\epsilon} \int_{\tau-\epsilon}^{\tau} c(x_1^\epsilon(t), v) dt = c(x_1(\tau), v). \quad (216)$$

On the other hand, $u(t)$ is continuous on $[\tau - \epsilon, \tau]$ for all sufficiently small ϵ , since u is left continuous at τ by Definition 1. Therefore, $c(x_1(t), u(t))$ is continuous with respect to t on $[\tau - \epsilon, \tau]$ for ϵ small. The mean value theorem gives

$$\frac{1}{\epsilon} \int_{\tau-\epsilon}^{\tau} c(x_1(t), u(t)) dt = c(x_1(\tilde{t}), u(\tilde{t})) \quad (217)$$

for some $\tilde{t} \in [\tau - \epsilon, \tau]$. Taking the limit $\epsilon \rightarrow 0^+$, the right hand side converges to $c(x_1(\tau), u(\tau))$. Combining this result with (216), we conclude that

$$\lim_{\epsilon \rightarrow 0^+} \frac{1}{\epsilon} \int_{\tau-\epsilon}^{\tau} \{c(x_1^\epsilon(t), v) - c(x_1(t), u(t))\} dt = c(x_1(\tau), v) - c(x_1(\tau), u(\tau)). \quad (218)$$

Lemma 20. Given $\epsilon \in [0, \tau]$, u^ϵ and (y_1, \dots, y_T) , the right derivative of the instantaneous cost function with respect to ϵ is given by

$$\forall t \in [i-1, i] \quad (219)$$

$$\frac{\partial_+}{\partial \epsilon} c(x_i^\epsilon(t), u^\epsilon(t)) = \frac{\partial}{\partial x_i^\epsilon} c(x_i^\epsilon(t), u(t))^T \Psi_i^\epsilon(t). \quad (220)$$

for each $i \in \{2, \dots, T\}$.

For $i = 1$ we have

$$\forall t \in (\tau, 1] \quad (221)$$

$$\frac{\partial_+}{\partial \epsilon} c(x_1^\epsilon(t), u^\epsilon(t)) = \frac{\partial}{\partial x_1^\epsilon} c(x_1^\epsilon(t), u(t))^T \Psi_1^\epsilon(t), \quad (222)$$

Similarly, the right derivative of the terminal cost function with respect to ϵ is given by

$$\frac{\partial_+}{\partial \epsilon} h(x^\epsilon(T)) = \frac{\partial}{\partial x^\epsilon} h(x^\epsilon(T))^T \Psi^\epsilon(T). \quad (223)$$

Proof. To prove the claim for the instantaneous cost, note that $u^\epsilon(t) = u(t)$ for all $t \in (\tau, T]$ and use the chain rule. The case for the terminal cost also follows from the chain rule.

Proposition 21. Bounded Cost Variations. Given $\epsilon \in [0, \tau]$, u^ϵ and (y_1, \dots, y_T) , the right derivative of the instantaneous cost function with respect to ϵ has the following uniform bound:

$$\forall t \in [i-1, i] \quad (224)$$

$$\left| \frac{\partial_+}{\partial \epsilon} c(x_i^\epsilon(t), u^\epsilon(t)) \right| \leq \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{L}'_i} \beta_i'^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m} \quad (225)$$

for each $i \in \{2, \dots, T\}$, where \mathcal{L}'_i is a finite set of sequences of non-negative integers of length $i-1$, and $\beta_i'^{(j_1, \dots, j_{i-1})}(x_0)$ is a finite positive constant that depends on x_0 and (j_1, \dots, j_{i-1}) but not on ϵ , u^ϵ , or (y_1, \dots, y_T) .

For $i = 1$ the bound is given by

$$\forall t \in (\tau, 1] \quad \left| \frac{\partial_+}{\partial \epsilon} c(x_1^\epsilon(t), u^\epsilon(t)) \right| \leq \beta_1'(x_0) \quad (226)$$

for some finite positive constant $\beta_1'(x_0)$.

Similarly, the right derivative of the terminal cost function with respect to ϵ has the following bound:

$$\left| \frac{\partial_+}{\partial \epsilon} h(x^\epsilon(T)) \right| \leq \sum_{(j_1, \dots, j_T) \in \mathcal{L}'_{T+1}} \beta_{T+1}'^{(j_1, \dots, j_T)}(x_0) \prod_{m=1}^T \|y_m\|_2^{j_m} \quad (227)$$

for some finite set \mathcal{L}'_{T+1} of sequence of non-negative integers and finite positive constants $\beta_{T+1}'^{(j_1, \dots, j_T)}(x_0)$.

Proof. The proof of this proposition is similar to that of Proposition 10. Take any $\epsilon \in [0, \tau]$. For $i \in \{2, \dots, T\}$, Assumption 3 along with Lemma 20 yields

$$\left| \frac{\partial_+}{\partial \epsilon} c(x_i^\epsilon(t), u^\epsilon(t)) \right| = \left| \frac{\partial}{\partial x_i^\epsilon} c(x_i^\epsilon(t), u(t))^T \Psi_i^\epsilon(t) \right| \quad (228)$$

$$\leq \left\| \frac{\partial}{\partial x_i^\epsilon} c(x_i^\epsilon(t), u(t)) \right\|_2 \cdot \|\Psi_i^\epsilon(t)\|_2 \quad (229)$$

$$\leq (K_7 + K_8 \|x_i^\epsilon(t)\|_2^{L_3}) \|\Psi_i^\epsilon(t)\|_2. \quad (230)$$

By Propositions 9 and 18, we have

$$\|x_i^\epsilon(t)\|_2 \leq \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{K}_i} \alpha_i^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m} \quad (231)$$

$$\|\Psi_i^\epsilon(t)\|_2 \leq \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{L}_i} \beta_i^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m} \quad (232)$$

for all $t \in [i-1, i]$. Substituting these into (230) and using the multinomial expansion formula, we conclude that

$$\forall i \in \{2, \dots, T\} \quad \forall t \in [i-1, i] \quad \left| \frac{\partial_+}{\partial \epsilon} c(x_i^\epsilon(t), u^\epsilon(t)) \right| \leq \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{L}'_i} \beta_i'^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m} \quad (233)$$

for some finite set \mathcal{L}'_i of sequences of non-negative integers and finite positive constants $\beta_i'^{(j_1, \dots, j_{i-1})}(x_0)$. Similarly, for $i = 1$ we have $\forall t \in (\tau, 1]$

$$\left| \frac{\partial_+}{\partial \epsilon} c(x_1^\epsilon(t), u^\epsilon(t)) \right| \leq (K_7 + K_8 \|x_1^\epsilon(t)\|_2^{L_3}) \|\Psi_1^\epsilon(t)\|_2 \quad (234)$$

$$\leq 2(K_7 + K_8 \alpha_1(x_0)^{L_3}) K_1 \rho_{\max}^{K_2} \quad (235)$$

$$\triangleq \beta_1'(x_0), \quad (236)$$

by Proposition 9 and Lemma 17.

To bound the right derivative of the terminal cost, note that

$$\left| \frac{\partial_+}{\partial \epsilon} h(x^\epsilon(T)) \right| = \left| \frac{\partial}{\partial x^\epsilon} h(x^\epsilon(T))^T \Psi^\epsilon(T) \right| \quad (237)$$

$$\leq (K_7 + K_8 \|x^\epsilon(T)\|_2^{L_3}) \|\Psi^\epsilon(T)\|_2 \quad (238)$$

$$= \left(K_7 + K_8 \|g(x_T^\epsilon(T), y_T)\|_2^{L_3} \right) \|\Psi^\epsilon(T)\|_2 \quad (239)$$

$$\leq \left(K_7 + K_8 \|g(x_T^\epsilon(T), y_T)\|_2^{L_3} \right) \times \left\| \frac{\partial}{\partial x_T^\epsilon} g(x_T^\epsilon(T), y_T) \right\|_2 \cdot \|\Psi_T^\epsilon(T)\|_2 \quad (240)$$

by Assumption 3, Proposition 15, and Lemma 20. One can apply Assumption (2d) to bound the norms of g and its Jacobian in terms of $\|x_T^\epsilon(T)\|_2$ and $\|y_T\|$. Then, (240) becomes a polynomial of $\|x_T^\epsilon(T)\|_2$ and $\|y_T\|$, multiplied by $\|\Psi_T^\epsilon(T)\|_2$. Finally, using (231) and (232) with $i = T$ to replace $\|x_T^\epsilon(T)\|_2$ and $\|\Psi_T^\epsilon(T)\|_2$, one can verify that

$$\left| \frac{\partial_+}{\partial \epsilon} h(x^\epsilon(T)) \right| \leq \sum_{(j_1, \dots, j_T) \in \mathcal{L}'_{T+1}} \beta_{T+1}'^{(j_1, \dots, j_T)}(x_0) \prod_{m=1}^T \|y_m\|_2^{j_m} \quad (241)$$

for some finite set \mathcal{L}'_{T+1} of sequence of non-negative integers and finite positive constants $\beta_{T+1}'^{(j_1, \dots, j_T)}(x_0)$.

Lemma 22. *Let x^ϵ be the perturbed state induced by the perturbed control u^ϵ , and let (y_1, \dots, y_T) be the given observations. Then, the function $\epsilon \mapsto c(x^\epsilon(t), u^\epsilon(t))$ is continuous with respect to $\epsilon \in [0, \tau]$ for all $t \in (\tau, T]$.*

Proof. Note that for $t \in (\tau, T]$ we have $u^\epsilon(t) = u(t)$. Thus, $c(x^\epsilon(t), u^\epsilon(t)) = c(x^\epsilon(t), u(t))$. The continuity of $x_1^\epsilon(t)$ with respect to $\epsilon \in [0, \tau]$ follows from Lemma 11. In particular, $x_1^\epsilon(1)$ is continuous with respect to ϵ . Next, suppose that $\epsilon \mapsto x_i^\epsilon(i)$ is continuous for some $i \in \{1, \dots, T\}$. Then, by Assumption (2b) and Corollary 7 it follows that $\epsilon \mapsto x_{i+1}^\epsilon(t)$ is continuous for all $t \in [i, i+1]$. Proceeding by mathematical induction, we conclude that $x^\epsilon(t)$ is continuous with respect to $\epsilon \in [0, \tau]$ for all $t \in (\tau, T]$. Therefore, $\epsilon \mapsto c(x^\epsilon(t), u(t))$ is continuous by Assumption 3.

Proposition 23. *Let $u \in U$ be a control, which yields the nominal state x . Let x^ϵ be the perturbed state induced by the perturbed control u^ϵ , and let (y_1, \dots, y_T) be the given observations. Then, the following bounds hold for all $\epsilon \in [0, \tau]$:*

$$\forall t \in (\tau, 1] \quad |c(x_1^\epsilon(t), u^\epsilon(t)) - c(x_1(t), u(t))| \leq \epsilon \beta_1'(x_0) \quad (242)$$

and

$$\begin{aligned} & \forall i \in \{2, \dots, T\} \quad \forall t \in [i-1, i] \\ & |c(x_i^\epsilon(t), u^\epsilon(t)) - c(x_i(t), u(t))| \\ & \leq \epsilon \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{L}'_i} \beta_i'^{(j_1, \dots, j_{i-1})}(x_0) \prod_{m=1}^{i-1} \|y_m\|_2^{j_m}, \end{aligned} \quad (243)$$

where $\beta_1'(x_0)$, $\beta_i'^{(j_1, \dots, j_{i-1})}(x_0)$ and \mathcal{L}' are as defined in Proposition 21.

Proof. For $\epsilon \in [0, \tau]$ and $t \in (\tau, T]$, the function $\epsilon \mapsto c(x^\epsilon(t), u^\epsilon(t))$ is continuous by Lemma 22 and finite by Proposition 10. It is also right differentiable with respect to ϵ for $\epsilon \in [0, \tau]$ and $t \in (\tau, T]$ by Lemma 20. Therefore, the mean value theorem (Corollary in Bourbaki and Spain (2004), p.15) along with Proposition 21 proves the claim.

Lemma 24. *Let x^ϵ be the perturbed state induced by the perturbed control u^ϵ , and let (y_1, \dots, y_T) be the given observations. Then, the function $\epsilon \mapsto h(x^\epsilon(T))$ is continuous with respect to $\epsilon \in [0, \tau]$.*

Proof. By the proof of Lemma 22 it follows that the function $\epsilon \mapsto x^\epsilon(T)$ is continuous with respect to $\epsilon \in [0, \tau]$. The continuity of h by Assumption 3 completes the proof.

Proposition 25. *Let $u \in U$ be a control, which yields the nominal state x . Let x^ϵ be the perturbed state induced by the perturbed control u^ϵ , and let (y_1, \dots, y_T) be the given observations. Then, the following bound holds for all $\epsilon \in [0, \tau]$:*

$$\begin{aligned} & |h(x^\epsilon(T)) - h(x(T))| \\ & \leq \epsilon \sum_{(j_1, \dots, j_T) \in \mathcal{L}'_{T+1}} \beta_{T+1}'^{(j_1, \dots, j_T)}(x_0) \prod_{m=1}^T \|y_m\|_2^{j_m}, \end{aligned} \quad (244)$$

where $\beta_{T+1}'^{(j_1, \dots, j_T)}(x_0)$ and \mathcal{L}'_{T+1} are as defined in Proposition 21.

Proof. The proof is very similar to that of Proposition 23. Use Proposition 10, Lemma 24, and Lemma 20 to show finiteness, continuity, and right differentiability of $\epsilon \mapsto h(x^\epsilon(T))$. Then use the same mean value theorem with Proposition 21 to prove the claim.

A.3 Expected Total Cost under Stochastic Observations

In this last part of the analysis, we finally let the observations (y_1, \dots, y_T) take random values; more formally, we treat them as a sequence of random variables $(Y_1(\omega), \dots, Y_T(\omega))$ for $\omega \in \Omega$, where $(\Omega, \mathcal{F}, \mathbb{P})$ is the probability space and each Y_i satisfies Assumption 4. With (τ, v) given and fixed, (ω, t) and ϵ uniquely determine the perturbed control u^ϵ and the observations, hence the resulting state trajectory x^ϵ and the costs $c(x^\epsilon(t), u^\epsilon(t))$, $h(x^\epsilon(T))$.

Lemma 26. *Let $([\tau, T], \mathcal{B}([\tau, T]), \lambda)$ be a measure space, where $\mathcal{B}([\tau, T])$ is the Borel σ -algebra on $[\tau, T]$ and λ is the Lebesgue measure. Let $\mu \triangleq \lambda \times \mathbb{P}$ be the product measure defined on the product space $(\Omega \times [\tau, T], \mathcal{F} \otimes \mathcal{B}([\tau, T]))$, where $\mathcal{F} \otimes \mathcal{B}([\tau, T])$ is the product σ -algebra. Then, the function $(\omega, t) \mapsto c(x^\epsilon(t), u^\epsilon(t))$ is $\mathcal{F} \otimes \mathcal{B}([\tau, T])$ -measurable for every $\epsilon \in [0, \tau]$.*

Proof. Take any $\epsilon \in [0, \tau]$. Then, u^ϵ is in U and thus the function $(Y_1(\omega), \dots, Y_T(\omega)) \mapsto x^\epsilon(t)$ is continuous for every $t \in [\tau, T]$ by Proposition 8. Therefore, the map $(\omega, t) \mapsto x^\epsilon(t)$ as a function of ω is \mathcal{F} -measurable for every $t \in [\tau, T]$. By Assumption 2, $x^\epsilon(t)$ is also right continuous

with respect to t for every $\omega \in \Omega$. Therefore, from Theorem 3 in [Gowrisankaran \(1972\)](#) it follows that $(\omega, t) \mapsto x^\epsilon(t)$ is measurable with respect to the product σ -algebra $\mathcal{F} \otimes \mathcal{B}([\tau, T])$.

On the other hand, $u^\epsilon(t)$ is piecewise continuous in t and is constant with respect to $(Y_1(\omega), \dots, Y_T(\omega))$. Therefore, $(\omega, t) \mapsto u^\epsilon(t)$ is also measurable with respect to $\mathcal{F} \otimes \mathcal{B}([\tau, T])$.

Finally, the continuity of the instantaneous cost c by Assumption 3 proves the claim.

Proposition 27. *For the perturbed control u^ϵ and the perturbed state x^ϵ , we have*

$$\begin{aligned} \frac{\partial_+}{\partial \epsilon} \mathbb{E} \left[\int_\tau^T c(x^\epsilon(t), u^\epsilon(t)) dt \right] \Big|_{\epsilon=0} \\ = \mathbb{E} \left[\int_\tau^T \frac{\partial}{\partial x} c(x(t), u(t))^T \Psi(t) dt \right], \end{aligned} \quad (245)$$

where Ψ is the state variation defined in Proposition 15.

Proof. By definition, the left hand side of (245) is

$$\begin{aligned} \frac{\partial_+}{\partial \epsilon} \mathbb{E} \left[\int_\tau^T c(x^\epsilon(t), u^\epsilon(t)) dt \right] \Big|_{\epsilon=0} \\ = \lim_{\epsilon \rightarrow 0^+} \mathbb{E} \left[\int_\tau^T \frac{1}{\epsilon} \{c(x^\epsilon(t), u^\epsilon(t)) - c(x(t), u(t))\} dt \right]. \end{aligned} \quad (246)$$

Consider the expected value above as the equivalent Lebesgue integral:

$$\int_\Omega \left(\int_{[\tau, T]} \frac{1}{\epsilon} \{c(x^\epsilon(t), u^\epsilon(t)) - c(x(t), u(t))\} d\lambda(t) \right) \times d\mathbb{P}(\omega). \quad (247)$$

By Lemma 26 the integrand is measurable in the product space. In addition, Proposition 23 shows that the absolute value:

$$\frac{1}{\epsilon} |c(x^\epsilon(t), u^\epsilon(t)) - c(x(t), u(t))| \quad (248)$$

is bounded by an integrable function for μ -a.e. (ω, t) . Indeed, if we let $\hat{c}(\omega, t)$ to be a function defined by

$$\hat{c}(\omega, t) = \begin{cases} \beta'_1(x_0) & \forall t \in [\tau, 1) \\ \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{L}'_i} \beta'_i(j_1, \dots, j_{i-1})(x_0) \prod_{m=1}^{i-1} \|Y_m\|_2^{j_m} & \forall t \in [i-1, i) \forall i \in \{1, \dots, T\}, \end{cases} \quad (249)$$

then we have

$$\frac{1}{\epsilon} |c(x^\epsilon(t), u^\epsilon(t)) - c(x(t), u(t))| \leq \hat{c}(\omega, t) \quad (250)$$

for every non-zero ϵ and μ -a.e. (ω, t) , and

$$\begin{aligned} \int_\Omega \left(\int_{[\tau, T]} \hat{c}(\omega, t) d\lambda(t) \right) d\mathbb{P}(\omega) &= \beta'_1(x_0)(1 - \tau) \\ &+ \sum_{i=2}^T \sum_{(j_1, \dots, j_{i-1}) \in \mathcal{L}'_i} \beta'_i(j_1, \dots, j_{i-1})(x_0) \\ &\times \mathbb{E} \left[\prod_{m=1}^{i-1} \|Y_m\|_2^{j_m} \right], \end{aligned} \quad (251)$$

where

$$\mathbb{E} \left[\prod_{m=1}^{i-1} \|Y_m\|_2^{j_m} \right] \leq \sqrt{\mathbb{E} [\|Y_1\|_2^{j_1}]} \sqrt{\mathbb{E} \left[\prod_{m=2}^{i-1} \|Y_m\|_2^{j_m} \right]} \quad (252)$$

$$\leq \vdots$$

$$\leq \prod_{m=1}^{i-1} \left(\mathbb{E} [\|Y_m\|_2^{j_m}] \right)^{\frac{1}{2^m}} < \infty \quad (253)$$

by the Cauchy-Schwarz inequality and Assumption 4.

Furthermore, Lemma 20 proves that

$$\begin{aligned} \lim_{\epsilon \rightarrow 0^+} \frac{1}{\epsilon} \{c(x^\epsilon(t), u^\epsilon(t)) - c(x(t), u(t))\} \\ = \frac{\partial}{\partial x} c(x(t), u(t))^T \Psi(t) \end{aligned} \quad (254)$$

for μ -a.e. (ω, t) . Therefore, the dominated convergence theorem yields

$$\begin{aligned} \frac{\partial_+}{\partial \epsilon} \int_\Omega \left(\int_\tau^T c(x^\epsilon(t), u^\epsilon(t)) d\lambda(t) \right) d\mathbb{P}(\omega) \Big|_{\epsilon=0} \\ = \int_\Omega \left(\int_\tau^T \frac{\partial}{\partial x} c(x(t), u(t))^T \Psi(t) d\lambda(t) \right) d\mathbb{P}(\omega). \end{aligned} \quad (255)$$

Proposition 28. *For the perturbed control u^ϵ and the perturbed state x^ϵ , we have*

$$\frac{\partial_+}{\partial \epsilon} \mathbb{E} [h(x^\epsilon(T))] \Big|_{\epsilon=0} = \mathbb{E} \left[\frac{\partial}{\partial x} h(x(T))^T \Psi(T) \right], \quad (256)$$

where Ψ is the state variation defined in Proposition 15.

Proof. The proof is similar to that of Proposition 27. By Assumption 3 and Proposition 8, $(Y_1(\omega), \dots, Y_T(\omega)) \mapsto h(x^\epsilon(T))$ is a continuous map for every $\epsilon \in [0, \tau]$. Therefore, the function $\omega \mapsto h(x^\epsilon(T))$ is \mathcal{F} -measurable.

By definition, the left hand side of (256) is

$$\begin{aligned} \frac{\partial_+}{\partial \epsilon} \mathbb{E} [h(x^\epsilon(T))] \Big|_{\epsilon=0} \\ = \lim_{\epsilon \rightarrow 0^+} \int_\Omega \frac{1}{\epsilon} \{h(x^\epsilon(T)) - h(x(T))\} d\mathbb{P}(\omega) \end{aligned} \quad (257)$$

The analysis above implies that the integrand is \mathcal{F} -measurable. In addition, Proposition 25 shows that the

absolute value:

$$\frac{1}{\epsilon} |h(x^\epsilon(T)) - h(x(T))| \quad (258)$$

is bounded by an integrable function for every non-zero ϵ and every $\omega \in \Omega$. Furthermore, Lemma 20 proves that

$$\lim_{\epsilon \rightarrow 0^+} \frac{1}{\epsilon} \{h(x^\epsilon(T)) - h(x(T))\} = \frac{\partial}{\partial x} h(x(T))^T \Psi(T) \quad (259)$$

for every $\omega \in \Omega$. Therefore, the dominated convergence theorem yields

$$\left. \frac{\partial_+}{\partial \epsilon} \int_{\Omega} h(x^\epsilon(T)) d\mathbb{P}(\omega) \right|_{\epsilon=0} = \int_{\Omega} \frac{\partial}{\partial x} h(x(T))^T \Psi(T) d\mathbb{P}(\omega). \quad (260)$$

Theorem 1. Mode Insertion Gradient. *Suppose that Assumptions 1 – 4 are satisfied. For a given (τ, v) , let u^ϵ denote the perturbed control of the form (117). The perturbed control u^ϵ and the stochastic observations (Y_1, \dots, Y_T) result in the stochastic perturbed state trajectory x^ϵ . For such u^ϵ and x^ϵ , let us define the mode insertion gradient of the expected total cost as*

$$\left. \frac{\partial_+}{\partial \epsilon} \mathbb{E} \left[\int_0^T c(x^\epsilon(t), u^\epsilon(t)) dt + h(x^\epsilon(T)) \right] \right|_{\epsilon=0}. \quad (261)$$

Then, this right derivative exists and we have

$$\begin{aligned} \left. \frac{\partial_+}{\partial \epsilon} \mathbb{E} \left[\int_0^T c(x^\epsilon(t), u^\epsilon(t)) dt + h(x^\epsilon(T)) \right] \right|_{\epsilon=0} \\ = c(x(\tau), v) - c(x(\tau), u(\tau)) \\ + \mathbb{E} \left[\int_{\tau}^T \frac{\partial}{\partial x} c(x(t), u(t))^T \Psi(t) dt \right. \\ \left. + \frac{\partial}{\partial x} h(x(T))^T \Psi(T) \right], \quad (262) \end{aligned}$$

where Ψ is the state variation defined in Proposition 15.

Proof. We first consider the instantaneous cost c . Split the integration interval to get

$$\begin{aligned} \mathbb{E} \left[\int_0^T c(x^\epsilon(t), u^\epsilon(t)) dt \right] \\ = \mathbb{E} \left[\int_0^{\tau-\epsilon} c(x^\epsilon(t), u^\epsilon(t)) dt \right] \\ + \mathbb{E} \left[\int_{\tau-\epsilon}^{\tau} c(x^\epsilon(t), u^\epsilon(t)) dt \right] \\ + \mathbb{E} \left[\int_{\tau}^T c(x^\epsilon(t), u^\epsilon(t)) dt \right] \quad (263) \end{aligned}$$

For the first two terms in the sum, recall that the evolution of the state $x^\epsilon(t)$ is not affected by any observations for all

$t \in [0, \tau]$. Thus,

$$\mathbb{E} \left[\int_0^{\tau-\epsilon} c(x^\epsilon(t), u^\epsilon(t)) dt \right] = \int_0^{\tau-\epsilon} c(x^\epsilon(t), u^\epsilon(t)) dt \quad (264)$$

$$\mathbb{E} \left[\int_{\tau-\epsilon}^{\tau} c(x^\epsilon(t), u^\epsilon(t)) dt \right] = \int_{\tau-\epsilon}^{\tau} c(x^\epsilon(t), u^\epsilon(t)) dt. \quad (265)$$

Note that (264) is constant with respect to ϵ , since for all $t \in [0, \tau - \epsilon]$ we have $u^\epsilon(t) = u(t)$ and $x^\epsilon(t) = x(t)$. On the other hand, for (265) we can apply Lemma 19 to obtain

$$\left. \frac{\partial_+}{\partial \epsilon} \mathbb{E} \left[\int_{\tau-\epsilon}^{\tau} c(x^\epsilon(t), u^\epsilon(t)) dt \right] \right|_{\epsilon=0} = c(x(\tau), v) - c(x(\tau), u(\tau)) \quad (266)$$

For the last term, Proposition 27 gives

$$\begin{aligned} \left. \frac{\partial_+}{\partial \epsilon} \mathbb{E} \left[\int_{\tau}^T c(x^\epsilon(t), u^\epsilon(t)) dt \right] \right|_{\epsilon=0} \\ = \mathbb{E} \left[\int_{\tau}^T \frac{\partial}{\partial x} c(x(t), u(t))^T \Psi(t) dt \right]. \quad (267) \end{aligned}$$

Finally, for the terminal cost h we have

$$\left. \frac{\partial_+}{\partial \epsilon} \mathbb{E} [h(x^\epsilon(T))] \right|_{\epsilon=0} = \mathbb{E} \left[\frac{\partial}{\partial x} h(x(T))^T \Psi(T) \right] \quad (268)$$

by Proposition 28.

Remark 5. Closed-loop Nominal Policy. *As far as the control is concerned, the analysis above only requires that the nominal control u is in U (as in Assumption 1) and the perturbed control u^ϵ is measurable with respect to $\mathcal{F} \otimes \mathcal{B}([\tau, T])$ (as in Lemma 26). To guarantee that these requirements are satisfied with a closed-loop nominal policy $\pi : \mathbb{R}^{n_x} \rightarrow \mathbb{R}^m$, it is sufficient that π is a measurable map and the induced nominal control trajectory $u(t) = \pi(x(t))$ for $t \in [0, T]$ belongs to U for any observations (y_1, \dots, y_T) . Note that the model of the control perturbation considered here is still open-loop:*

$$u^\epsilon(t) = \begin{cases} v & \text{if } t \in (\tau - \epsilon, \tau] \\ \pi(x(t)) & \text{otherwise.} \end{cases} \quad (269)$$

That is, the nominal state trajectory x is used in the control feedback. This is not to be confused with the closed-loop perturbation:

$$u_{\text{closed}}^\epsilon(t) = \begin{cases} v & \text{if } t \in (\tau - \epsilon, \tau] \\ \pi(x^\epsilon(t)) & \text{otherwise,} \end{cases} \quad (270)$$

where the perturbed state trajectory x^ϵ is fed back to the controller.