

# A Separation-Based Design to Data-Driven Control for Large-Scale Partially Observed Systems

Dan Yu<sup>1</sup>, Mohammadhussein Rafieisakhaei<sup>2</sup> and Suman Chakravorty<sup>1</sup>

**Abstract**—This paper studies the partially observed stochastic optimal control problem for systems with state dynamics governed by Partial Differential Equations (PDEs) that leads to an extremely large problem. First, an open-loop deterministic trajectory optimization problem is solved using a black box simulation model of the dynamical system. Next, a Linear Quadratic Gaussian (LQG) controller is designed for the nominal trajectory-dependent linearized system, which is identified using input-output experimental data consisting of the impulse responses of the optimized nominal system. A computational nonlinear heat example is used to illustrate the performance of the approach.

## I. INTRODUCTION

In this paper, we consider the stochastic control of partially observed nonlinear dynamical systems that are governed by Partial Differential Equations (PDEs). In particular, we propose a novel data-based approach to the solution of very large Partially Observed Markov Decision Processes (POMDPs) wherein the underlying state space is obtained from the discretization of a PDE; problems whose solution has never been hitherto attempted using approximate MDP-based techniques.

It is well-known that the global optimal solution for MDPs can be found by solving the Hamilton-Jacobi-Bellman (HJB) equation. The solution techniques can be further divided into model-based and model-free techniques, according to whether the solution methodology uses an analytical model of the system or it uses a black box simulation model or actual experiments. The Reinforcement Learning (RL) techniques [1, 2] that are based on the Differential Dynamic Programming (DDP) or iLQG approach [3, 4] have shown the potential for RL algorithms to scale to higher dimensional continuous state and control space problems, such as high dimensional robotic task planning and learning problems.

Fundamentally, rather than solving the derived “Dynamic Programming” problem as in the majority of the approaches above that requires the optimization of the feedback law, our approach is to directly solve the original stochastic optimization problem in a “separated open-loop/closed-loop” fashion wherein: 1) we solve an open-loop deterministic optimization

problem to obtain an optimal nominal trajectory in a model-free fashion, and then 2) we design a closed-loop controller for the resulting linearized time-varying system around the optimal nominal trajectory, again in a model-free fashion. Nonetheless, the above “divide and conquer” strategy can be shown to be near-optimal [5, 6].

The primary contributions of the proposed approach are:

1) We specify a detailed set of experiments to accomplish the closed-loop controller design for any unknown nonlinear system, no matter how high dimensional. This series of experiments consists of a sequence of input perturbations to collect the impulse responses of the system, first to find an optimized nominal trajectory, and then to recover the Linear Time-Varying (LTV) system corresponding to the perturbations of the nominal system in order to design an LQG controller.

2) In general, for large-scale systems with partially observed states, the system identification algorithms such as time-varying Eigensystem Realization Algorithm (ERA) [7] automatically construct reduced order model of the LTV system, which results in a reduced order estimator and controller. Therefore, even for large-scale systems, such as partially observed systems with dynamics governed by PDEs, the computation of the feedback policy is computationally tractable. For instance, in the partially observed nonlinear heat control problem considered in this paper, the complexity is reduced by  $O(10^5)$  when compared to DDP-based RL techniques.

3) We provide a unification of traditional linear and nonlinear optimal control techniques with Adaptive Dynamic Programming (ADP) [8] and RL techniques in the context of Stochastic Dynamic Programming problems.

## II. PROBLEM SETUP

Consider a discrete-time nonlinear dynamical system:

$$x_{k+1} = f(x_k, u_k, w_k), \quad y_k = h(x_k, v_k), \quad (1)$$

where  $x \in \mathbb{R}^{n_x}$ ,  $y \in \mathbb{R}^{n_y}$ ,  $u \in \mathbb{R}^{n_u}$  are the state, measurement and control vectors, respectively, the system and measurement functions,  $f(\cdot)$  and  $h(\cdot)$ , are nonlinear, and  $\{w_k, v_k, k \geq 0\}$  are zero-mean, uncorrelated Gaussian white noises with covariances  $W$  and  $V$ , respectively. In considering PDEs, the dynamics are discretized using Finite Difference (FD) or Finite Element (FE), which can lead to a state space problem consisting of, e.g., millions of states.

The belief  $b(x_k)$  is the conditional distribution of the state  $x_k$  given all past data. In this paper, we consider Gaussian beliefs denoted by  $b_k := (\mu_k, \Sigma_k)$ , where  $\mu_k$  and  $\Sigma_k$  are the mean and covariance (whose size is  $O(n_x^2)$ , which for a PDE with large  $n_x$  is extremely large). We denote the belief dynamics by  $b_{k+1} = \tau(b_k, u_k, y_{k+1})$ .

\*This material is based upon work partially supported by NSF under Contract Nos. CNS-1646449 and Science & Technology Center Grant CCF-0939370, the U.S. Army Research Office under Contract No. W911NF-15-1-0279, and NPRP grant NPRP 8-1531-2-651 from the Qatar National Research Fund, a member of Qatar Foundation, AFOSR contract Dynamic Data Driven Application Systems (DDAS) contract FA9550-17-1-0068 and NSF NRI project ECCS-1637889.

<sup>1</sup>D. Yu and S. Chakravorty are with the Department of Aerospace Engineering, and <sup>2</sup>M. Rafieisakhaei is with the Department of Electrical and Computer Engineering, Texas A&M University, College Station, Texas, 77840 USA. {yudan198811@hotmail.com, mrafieis, schakrav@tamu.edu}

**Stochastic Control Problem:** Given  $b_0$  and a finite time horizon of  $N > 0$ , for unknown nonlinear  $f(\cdot)$  and  $h(\cdot)$ , find the control policy  $\pi = \{\pi_0, \pi_1, \dots, \pi_{N-1}\}$ , where  $\pi_k$  is the control policy at time  $k$  and  $u_k = \pi_k(b_k)$ , such that

$$J_\pi = \mathbb{E}\left(\sum_{k=0}^{N-1} c_k(b_k, u_k) + c_N(b_N)\right), \quad (2)$$

is minimized, where  $\{c_k(\cdot, \cdot)\}_{k=0}^{N-1}$  denotes the immediate cost function, and  $c_N(\cdot)$  denotes the terminal cost.

### III. SEPARATION-BASED FEEDBACK CONTROL DESIGN

Let  $\{\bar{u}_k\}_{k=0}^{N-1}$ ,  $\{\bar{\mu}_k\}_{k=0}^N$ ,  $\{\bar{y}_k\}_{k=0}^N$ ,  $\{\bar{b}_k\}_{k=0}^N$  denote the nominal control, state, observation, and belief trajectories of the system, respectively, where given  $\bar{u}_k = \pi_k(\bar{b}_k)$ , we have:

$\bar{\mu}_{k+1} = f(\bar{\mu}_k, \bar{u}_k, 0)$ ,  $\bar{y}_k = h(\bar{\mu}_k, 0)$ ,  $\bar{b}_{k+1} = \tau(\bar{b}_k, \bar{u}_k, \bar{y}_{k+1})$ , with the initial conditions of  $\bar{b}_0 = b_0$ , and  $\bar{\mu}_0 = \mathbb{E}[b_0]$ .

The nominal cost and its first order expansion are [5, 9]:

$$\begin{aligned} \bar{J} &:= \sum_{k=0}^{N-1} c_k(\bar{b}_k, \bar{u}_k) + c_N(\bar{b}_N), \\ J &\approx \bar{J} + \underbrace{\sum_{k=0}^{N-1} (C_k^b(b_k - \bar{b}_k) + C_k^u(u_k - \bar{u}_k)) + C_K^b(b_N - \bar{b}_N)}_{=: \delta J}. \end{aligned}$$

*Theorem 1 (Cost Function Linearization Error):* The expected first-order linearization error of the cost function is zero,  $\mathbb{E}(\delta J) = 0$ .

Theorem 1 shows that the first order approximation of the stochastic cost function is dominated by the nominal cost and depends only on the nominal trajectories of the system, independent of the feedback gain. Therefore, the design of the optimal feedback gain can be separated from the design of the optimal nominal trajectory of the system. As a result, the stochastic optimal control problem can be divided into two separate problems: the first is a deterministic problem to design the open-loop optimal control sequence, and hence, the optimal nominal trajectory of the system. The second problem is the design of an optimal linear feedback law to track the nominal trajectory (which is the optimal belief state trajectory unlike typical trajectory optimization based RL methods designed for fully observed problems such as [1, 2]).

We propose a three-step framework to solve the stochastic feedback control problem as follows.

**Step 1. Open-Loop Trajectory Optimization in Belief Space.** Solve the open-loop optimization problem given  $b_0$ :

$$\begin{aligned} \{u_k^*\}_{k=0}^{N-1} &= \arg \min_{\{u_k\}_{k=0}^{N-1}} \bar{J}(\{b_k\}_{k=0}^N, \{u_k\}_{k=0}^{N-1}), \\ b_{k+1} &= \tau(b_k, u_k, \bar{y}_{k+1}), \end{aligned} \quad (3)$$

where the nominal observations  $\bar{y}_k$  are generated as follows:  $x_{k+1} = f(x_k, u_k, 0)$ ,  $\bar{y}_k = h(x_k, 0)$  with  $x_0 = \mu_0$ . Given the nominal observations  $\bar{y}_k$ , the belief evolution is deterministic and the above is a deterministic optimization problem [10].

The open-loop optimization problem is solved using the gradient descent approach [11, 12] utilizing an Ensemble Kalman Filter (EnKF) [13]. Denote the initial guess of the control sequence by  $U^{(0)} = \{u_k^{(0)}\}_{k=0}^{N-1}$ , and the corresponding

belief state estimated using EnKF by  $\mathcal{B}^{(0)} = \{b_k^{(0)}\}_{k=0}^N$ . The control policy is updated iteratively via

$$U^{(n+1)} = U^{(n)} - \alpha \nabla_U \bar{J}(\mathcal{B}^{(n)}, U^{(n)}), \quad (4)$$

until a convergence criterion is met, where  $U^{(n)} = \{u_k^{(n)}\}_{k=0}^{N-1}$  is the control sequence in the  $n^{\text{th}}$  iteration,  $\mathcal{B}^{(n)} = \{b_k^{(n)}\}_{k=0}^N$  denotes the corresponding belief, and  $\alpha$  is the step-size parameter. Finally, denote the nominal belief by  $\{\bar{\mu}_k, \bar{\Sigma}_k\}_{k=0}^N$ .

**Step 2. Linear Time-Varying System Identification.** We linearize the system (1) around the nominal mean trajectory  $\{\bar{\mu}_k\}$ . For simplicity, assume that the control and disturbance enter through same channels and the noise is purely additive:

$\delta x_{k+1} = A_k \delta x_k + B_k(\delta u_k + w_k)$ ,  $\delta y_k = C_k \delta x_k + v_k$ , (5) where  $\delta x_k = x_k - \bar{\mu}_k$ ,  $\delta u_k = u_k - \bar{u}_k$ ,  $\delta y_k = y_k - h(\bar{\mu}_k, 0)$  describe the state, control and measurement deviations from the nominal trajectory respectively, and

$$\begin{aligned} A_k &= \frac{\partial f(x, u, w)}{\partial x} \Big|_{\bar{\mu}_k, \bar{u}_k, 0}, B_k = \frac{\partial f(x, u, w)}{\partial u} \Big|_{\bar{\mu}_k, \bar{u}_k, 0}, \\ C_k &= \frac{\partial h(x, v)}{\partial x} \Big|_{\bar{\mu}_k, 0}. \end{aligned} \quad (6)$$

We identify the system (5) using impulse responses of the system via the time-varying ERA [7]. Denote the identified system's deviations by

$\delta a_{k+1} = \hat{A}_k \delta a_k + \hat{B}_k(\delta u_k + w_k)$ ,  $\delta y_k = \hat{C}_k \delta a_k + v_k$ , (7)

where  $\delta a_k \in \mathfrak{R}^{n_r}$  denotes the reduced order model (ROM) deviation states, and  $n_r \ll n_x$ , thereby automatically providing a compact parametrization of the problem.

**Step 3. Closed-Loop Controller Design.** Given system (7), we design the closed-loop controller to follow the optimal nominal trajectory, which is to minimize the cost function

$$J_f = \sum_{k=0}^{N-1} (\delta \hat{a}'_k Q_k \delta \hat{a}_k + \delta u'_k R_k \delta u_k) + \delta \hat{a}'_N Q_N \delta \hat{a}_N, \quad (8)$$

where  $\delta \hat{a}_k$  denotes the estimates of the deviation state  $\delta a_k$ ,  $Q_k, Q_N$  are positive definite, and  $R_k$  is positive semi-definite. For the linear system (7), the ‘‘separation principle’’ of linear control theory can be used [14], and the design of the optimal linear stochastic controller can be separated into the decoupled design of a KF and a fully observed optimal LQR controller.

A flow chart for the Separation-based Nonlinear Stochastic Control Design is shown in Fig. 2.

## IV. EXPERIMENTS

We test the method on a one-dimensional nonlinear heat transfer problem. Let  $T(x, t)$  be the temperature distribution at location  $x$  and time  $t$ ,  $K(x, T)$  be the thermal diffusivity,  $\eta$  be the convective heat transfer coefficient,  $u(t)$  be the external heat sources and  $L$  be the length of the slab. The heat transfer PDE along the slab along with its boundary conditions is:

$$\frac{\partial T}{\partial t} = K(x, T) \frac{\partial^2 T}{\partial x^2} - \eta T + u(t), \quad (9)$$

$$T(x, 0) = 100^\circ F, \quad \frac{\partial T}{\partial x} \Big|_{x=0} = 0, \quad T(L, t) = 150^\circ F. \quad (10)$$

The system is discretized using finite difference method with a 100 equally-spaced grid points. There are five point sources evenly located between  $[0.1L, 0.9L]$ , where the sensors are placed, as well. The total simulation time is 62.5s with

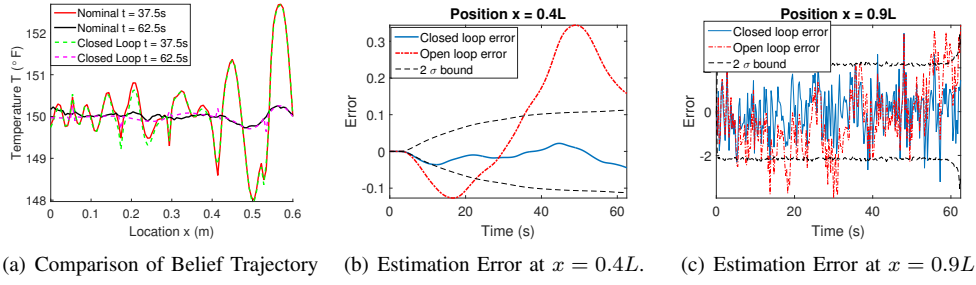


Fig. 1. Performance of the Proposed Approach

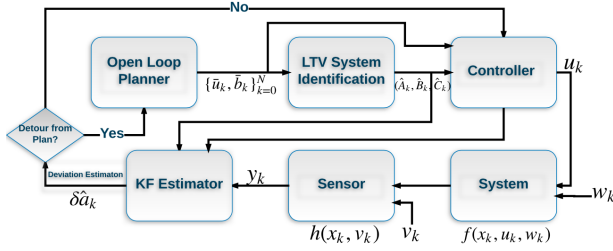


Fig. 2. Separation-Based Stochastic Feedback Control Algorithm a time-step of  $0.25s$ . The control objective is to reach the target temperature  $T_f = (150 \pm 3)^\circ F$  for the entire field within  $37.5s$ , and keep the temperature at  $(150 \pm 3)^\circ F$  between  $[37.5, 62.5]s$ .

The open-loop optimal nominal (belief mean) trajectory and optimal control are shown in Fig. 3. For the identified reduced order system, we have  $\hat{A}_k \in \mathbb{R}^{20 \times 20}$ . The feedback design decouples into the solution of two  $20 \times 20$  Riccati equations, one for the controller and one for the Kalman filter. Note that if we were to use an iLQG-based design, the size of the state space would be 10100, and the policy evaluation step would require the solution of a  $10100 \times 10100$  Riccati equation.

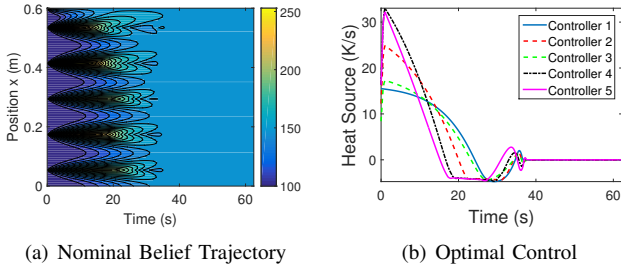


Fig. 3. Open Loop Optimization Solution

With the identified linearized system, we design the closed-loop controller. We run 1000 individual simulations with process noise  $w_k \sim N(0, I)$  and measurement noise  $v_k \sim N(0, I)$ . In Fig. 1(a), we compare the averaged closed-loop trajectory with the nominal trajectory at time  $t = 37.5s, t = 62.5s$ . In Figs. 1(b)-(c), we randomly choose two positions, and show the errors between the actual trajectory and optimal trajectory with  $2\sigma$  bounds in one simulation. For comparison, the open-loop error is also shown in the figure.

It is observed that the averaged state estimates over 1000 Monte-Carlo simulations runs are close to the open-loop optimal trajectory, which implies that the control objective to minimize the expected cost function could be achieved using the proposed approach. In this partially observed problem, the

computational complexity of designing the online estimator and controller using the identified ROM model is reduced by the order of  $O(\frac{n_4}{n_2}) = O(10^5)$ , and for a general three dimensional problem this reduction could be even more significant.

## V. CONCLUSION

In this paper, we proposed a separation-based design of the stochastic optimal control problem for systems with unknown nonlinear dynamics and partially observed states. The open-loop optimization and system identification are efficiently implemented offline using the impulse responses of the system, and an LQG controller based on the ROM is implemented online, which is computationally fast. We showed the performance of the proposed approach on a one-dimensional nonlinear heat transfer problem.

## BIBLIOGRAPHY

- [1] R. Akrou, A. Abdolmaleki, H. Abdulsamad, and G. Neumann, "Model Free Trajectory Optimization for Reinforcement Learning," in *Proceedings of the International Conference on Machine Learning*, 2016.
- [2] E. Todorov and Y. Tassa, "Iterative Local Dynamic Programming," in *Proc. of the IEEE Int. Symposium on ADP and RL*, 2009.
- [3] S. Levine and P. Abbeel, "Learning Neural Network Policies with Guided Search under Unknown Dynamics," in *Advances in Neural Information Processing Systems*, 2014.
- [4] S. Levine and K. Vladlen, "Learning Complex Neural Network Policies with Trajectory Optimization," in *Proceedings of the International Conference on Machine Learning*, 2014.
- [5] M. Rafieisakhaei, S. Chakravorty, and P. Kumar, "Near-Optimal Belief Space Planning via T-LQG," *arXiv preprint arXiv:1705.09415*, 2017.
- [6] M. Rafieisakhaei, S. Chakravorty, and P. Kumar, "A near-optimal separation principle for nonlinear stochastic systems arising in robotic path planning and control," *arXiv preprint arXiv:1705.08566*, 2017.
- [7] M. Majji, J.-N. Juang, and J. L. Junkins, "Time-varying Eigensystem Realization Algorithm," *Journal of Guidance, Control, and Dynamics*, vol. 33, no. 1, pp. 13–28, 2010.
- [8] R. P. Bithmead, V. Wertz, and M. Gerers, *Adaptive Optimal Control: The Thinking Man's G.P.C.* Prentice Hall Professional Technical Reference, 1991.
- [9] M. Rafieisakhaei, S. Chakravorty, and P. Kumar, "Belief Space Planning Simplified: Trajectory-Optimized LQG (T-LQG)," *arXiv preprint arXiv:1608.03013*, 2016.
- [10] R. Platt, R. Tedrake, L. Kaelbling, and T. Lozano-Perez, "Belief space planning assuming maximum likelihood observations," in *Proceedings of Robotics: Science and Systems (RSS)*, June 2010.
- [11] A.E. Bryson and W. Denham, "A steepest-ascent method for solving optimum programming problems," *Journal of Applied Mechanics*, vol. 29, no. 2, 1962.
- [12] A. Gosavi, *Simulation-based optimization: Parametric optimization techniques and reinforcement learning*. Norwell, MA, USA: Kluwer Academic Publishers, 2003.
- [13] S. Gillijns *et al.*, "What Is the Ensemble Kalman Filter and How Well Does it Work?" in *Proceedings of the 2006 American Control Conference*, 2006, pp. 4448–4453.
- [14] D. P. Bertsekas, *Dynamic Programming and Optimal Control, Two Volume Set*, 2nd ed. Athena Scientific, 1995.