

---

# Control Adaptation via Meta-Learning Dynamics

---

James Harrison<sup>\*,1</sup>, Apoorva Sharma<sup>\*,1</sup>, Roberto Calandra<sup>2</sup>, Marco Pavone<sup>1</sup>

<sup>1</sup>Stanford University, <sup>2</sup>University of California, Berkeley

{jharrison, apoorva}@stanford.edu, roberto.calandra@berkeley.edu, pavone@stanford.edu

## Abstract

Model-based control relies on having an accurate dynamics model. In situations where the dynamics of the robot is not static (e.g., after a delivery quadcopter picks or drops a package), it is paramount to accurately both estimate how the dynamics is influenced and account for these changes at the control level. One way of doing so is to condition the dynamics model on some latent variable which represents the context. In this paper, we present CAMELiD, a data-driven approach to learn controllers that can quickly adapt to previously unseen dynamics. This is done through a high-capacity neural network trained in a meta-learning manner to infer the context of the dynamics and maintain approximate posteriors over the dynamics model. This model is paired with an uncertainty-aware model-based control framework based on iLQG. We demonstrate CAMELiD on a quadrotor delivery experiment, and show it results in substantial performance improvement over baselines.

## 1 Introduction

Model-based learning control is a compelling framework for robotic learning due to its adaptivity, sample efficiency and scalability [9, 6, 25, 3]. However, high capacity models such as neural networks can not be rapidly adapted online, thus preventing their use in adaptive control-based methods [14]. Gaussian Process models are attractive, as their Bayesian predictions are useful even with limited data online, but they scale poorly [5]. Furthermore, incorporating prior information beyond relatively simple characteristics such as smoothness or periodicity properties into GP regression can be challenging, limiting their sample efficiency. Recently, meta-learning approaches have been shown to be effective choices for rapid adaptation of high capacity models online [8]. However, previous approaches such as [4] naively apply sampling-based control methods in a certainty-equivalent fashion with a point estimate of the posterior dynamics model. In contrast, in this work we maintain a more expressive estimate of the posterior dynamics model. We show how a Bayesian approach to meta-learning dynamics models can be paired with an uncertainty-aware (or cautious, in the terminology of adaptive control [19]) control algorithm to provide a more robust framework for nonlinear adaptive control. In contrast to certainty-equivalent methods such as [4] that only leverage an adaptive, deterministic dynamics model, we exhibit better performance and stability characteristics.

**Problem Statement.** We aim to control an unknown system, from which we have observed trajectory data with varying latent parameters. For example, if we have a delivery quadrotor, we may observe flight data for the vehicle with different loading conditions. Formally, we consider a discrete-time dynamical system of the form

$$x_{t+1} = f(x_t, u_t; \theta) + \epsilon, \quad (1)$$

where  $f(\cdot, \cdot; \cdot)$  is a continuous nonlinear function,  $x \in \mathcal{X}$  denotes the state, and  $u \in \mathcal{U}$  is the action. The dynamics are parameterized by latent parameters  $\theta \sim p(\theta)$ , and  $\epsilon$  is a disturbance term that is assumed to be Gaussian with known variance  $\Sigma_\epsilon$ . We assume the cost function,  $\text{Cost}(x, u)$  is known.

---

\* Authors contributed equally. Order determined randomly.

Critically, however, we do not assume knowledge  $f$ ,  $\theta$ , or  $p(\theta)$  to be known. We assume interaction with the system is episodic, and  $\theta$  is sampled at the beginning of each episode and is constant over the episode.

We assume access to trajectory data  $\mathcal{D} = \{D_\tau(\theta_j)\}_{j=1}^J$ , where  $D_\tau(\theta_j) = \{(x_t, u_t), x_{t+1}\}_{t=0}^{\tau-1}$ . From this data, we wish to learn an approximate model for both the dynamics function  $f$ , as well as an approximate prior over  $\theta$ . This setting is similar to the standard system identification setting, although  $\theta$  is not explicitly known to system designer as it may be in system identification. The problem can thus be formulated as, given dataset  $\mathcal{D}$ , problem horizon  $T$ , and initial state distribution  $\rho(x_0)$ ,

$$\min_{\pi \in \Pi} \mathbb{E}_{\theta \sim p(\theta), x_0 \sim \rho(x_0)} \left[ \sum_{t=0}^T \text{Cost}(x_t, \pi(x_t)) \right], \quad (2)$$

subject to the dynamics, (1). We write  $\pi$  to denote a policy, and  $\Pi$  to denote the set of possible policies. While this notation is more commonly used in the model-free regime, we take a model-based approach in this paper. This problem can be seen as a generalization of the standard adaptive control problem [2] to the ‘‘multi-task’’ case [1] or as the meta-reinforcement learning problem [13, 20, 22] in which the cost function is known.

## 2 Meta-learning via ALPaCA models

Model-based control in uncertain environments with nonlinear dynamics requires a regression model that is *high-capacity*, able to accurately capture the nonlinear dynamics; *data-efficient*, so as to make the most of the limited transition data observed online; and *computationally efficient*, so that it may be used in real-time.

ALPaCA is an efficient Bayesian approach to function regression in a meta-learning/multi-task learning setting [15]. In this formulation, function regression is performed by using observed data to update a posterior belief over a family of possible functions. In the context of adaptive dynamics modeling, we are interested in the family of functions  $f : \mathbb{R}^{n_x} \times \mathbb{R}^{n_u} \times \Theta \rightarrow \mathbb{R}^{n_x}$  parameterized by latent variables  $\theta \in \Theta$  distributed according to  $p(\theta)$ . ALPaCA models a distribution over  $f$  through the parametric model  $\hat{f}(x) = K^T \phi(x, u; w)$ , where  $\phi(x, u; w)$  represents a feed-forward neural network with weights  $w$  and output dimension  $n_\phi$ , and  $K$  is a  $n_\phi \times n_x$  matrix which can be thought of as the linear last-layer weights of the network. ALPaCA maintains a matrix-normal distribution over the last layer weights  $K \sim \mathcal{MN}(\bar{K}_0, \Lambda_0^{-1}, \Sigma_\epsilon)$ . This particular form reduces the task of online function regression to Bayesian Linear Regression in the feature space defined by  $\phi$ . This prior over  $K$  is self-conjugate, so at time  $t$  the posterior over  $K$  is  $p(K | X_t, \Phi_{t-1}) = \mathcal{MN}(\bar{K}_t, \Lambda_t^{-1}, \Sigma_\epsilon)$ , where

$$\Lambda_t = \Phi_{t-1}^T \Phi_{t-1} + \Lambda_0 \quad \text{and} \quad \bar{K}_t = \Lambda_t^{-1} (\Phi_{t-1} X_t + \Lambda_0 \bar{K}_0),$$

where  $X_t^T = [x_1, \dots, x_t]$  and  $\Phi_{t-1}^T = [\phi(x_0, u_0), \dots, \phi(x_{t-1}, u_{t-1})]$  are constructed from the transition data observed up until time  $t$ .

How quickly and how well this Bayesian linear regression algorithm can learn online depends on how well the feature mapping  $\phi$  and the prior on  $K$  are suited to the domain. Critically, ALPaCA employs a meta-learning strategy, training offline to optimize the performance of the online learning algorithm for the domain of interest.

Specifically, leveraging an offline dataset  $\mathcal{D}$  containing collections of  $(x, u, x')$  tuples from different contexts (i.e. values of  $\theta$ ), ALPaCA optimizes the weights  $w$  of the feature mapping and the parameters  $\bar{K}_0, \Lambda_0^{-1}$  of the prior. In this offline meta-learning phase, we sample from  $\mathcal{D}$  a collection  $D_T(\theta_j) = \{(x_{j,i}, u_{j,i}, x_{j,i+1})\}_{i=0}^{T-1}$ , with  $x_{j,i+1} \sim \mathcal{N}(f(x_{j,i}, u_{j,i}; \theta_j), \Sigma_\epsilon)$ . We use the BLR update rules to compute the posterior given the first  $t$  samples from the collection, which gives the posterior predictive density for new transitions

$$\hat{x}_{t+1} \sim q(x_{t+1} | x_t, u_t, \tilde{D}_t) = \mathcal{N}(\bar{K}_t^T \phi(x_t, u_t; w), \Sigma_t), \quad (3)$$

where

$$\Sigma_t = (1 + \phi^T(x_t, u_t; w) \Lambda_{t-1}^{-1} \phi(x_t, u_t; w)) \odot \Sigma_\epsilon, \quad (4)$$

and  $\tilde{D}_t$  is an online context dataset. ALPaCA optimizes the predictive density for  $w, \bar{K}_0$ , and  $\Lambda_0^{-1}$  via stochastic gradient descent, using the negative log likelihood of the remaining samples in  $\tilde{D}_T$  under this posterior as a loss function. This is equivalent to minimizing the KL divergence between ALPaCA’s posterior distributions and the empirical posterior distributions seen in the dataset  $\mathcal{D}$ .

### 3 CAMELiD

**Forward Uncertainty Propagation.** In model-based RL and adaptive control, it is necessary to make good *long-term* predictions. However, models trained to perform one-step prediction have been shown to be susceptible to compounding error when used to predict longer horizons [7, 11]. We derive an approach to Gaussian uncertainty propagation tailored to the ALPaCA dynamics model that enables computationally tractable yet accurate cost estimation over multistep predictions. Specifically, we will replace the one-step max likelihood objective of ALPaCA with the multistep objective

$$\max_{K_0, \Lambda_0, w} \sum_{j=1}^J p(x_T^{(j)} | x_{0:t}^{(j)}, u_{0:t}^{(j)}, \bar{K}_0, \Lambda_0, w), \quad (5)$$

for arbitrary  $t$  and  $T \geq t + 1$ . Due to the nonlinear dynamics, exactly computing the density of  $x_T^{(j)}$  is intractable, and so we leverage linearization and moment matching to compute a tractable Gaussian approximation. We do not consider the updates to the dynamics model that would occur online.<sup>1</sup>

Due to space constraints, we will exclude the derivation of the forward uncertainty propagation. The result, however, is that  $p(x_{i+1}^{(j)} | x_{0:t}^{(j)}, u_{0:t}^{(j)}, \bar{K}_0, \Lambda_0, w) \approx \mathcal{N}(\mu_{i+1}, \Sigma_{i+1})$ , with

$$\mu_{i+1} = \bar{K}^T \phi(\mu_i, u_i) \quad \text{and} \quad \Sigma_{i+1} = \Sigma_\epsilon + \bar{K}^T \phi_x(\mu_i, u_i) \Sigma_i \phi_x^T(\mu_i, u_i) \bar{K} + \mathbb{E}[\phi_i^T \Lambda^{-1} \phi_i] \odot \Sigma_\epsilon, \quad (6)$$

with  $\bar{K} := \bar{K}_t$ ,  $\Lambda := \Lambda_t$ ,  $\phi_x = \nabla_x \phi$ ,  $\phi_i = \phi(x_i, u_i)$ , and

$$\mathbb{E}[\phi_i^T \Lambda^{-1} \phi_i] = \text{tr}(\phi_x^T(\mu_i, u_i) \Lambda^{-1} \phi_x(\mu_i, u_i) \Sigma_i) + \phi^T(\mu_i, u_i) \Lambda^{-1} \phi(\mu_i, u_i).$$

This objective may be used to train the ALPaCA dynamics model for extended-horizon prediction.

**Uncertainty-Aware Control.** In this section, we derive a locally-optimal stochastic control algorithm, which may be seen as a version of iLQR-style control schemes [24]. This controller alternates between forward passes, in which the system dynamics are linearized, and backwards passes in which the control is locally optimized. These local quadratic optimal control methods have been shown to be practical in numerous applications, from large-scale robot control [23] to several applications in reinforcement learning [18, 12]. The controller derivation proceeds in terms of deviations from a nominal trajectory, which are written  $\delta x_i$  and  $\delta u_i$ . The forward pass is done with nominal controls  $\bar{u}_i$  to generate nominal means  $\bar{\mu}_i$ . For the backward pass, we must linearize the ALPaCA dynamics model around the nominal trajectory. Let  $A_i = \phi_x(\bar{\mu}_i, \bar{u}_i)$  and  $B_i = \phi_u(\bar{\mu}_i, \bar{u}_i)$ . Then, we may write the linearized deviation dynamics as  $\delta \mu_{i+1} = \bar{K}^T A_i \delta \mu_i + \bar{K}^T B_i \delta u_i$ . We will leverage the approximate relationship,  $\phi(x_i, u_i) \approx \bar{\phi}_i + A_i \delta \mu_i + B_i \delta u_i$ . Note that  $\delta x_{i+1} \sim \mathcal{N}(\delta \mu_{i+1}, \Sigma_{i+1})$ .

Quadratizing, the cost-to-go may be written as  $J_i(\delta x_i) = V_i + V_{x,i}^T \delta x_i + \frac{1}{2} \delta x_i^T V_{xx,i} \delta x_i$ . By the Bellman equation,  $J_i(\delta x_i)$  may be written

$$\min_{u \in \mathcal{U}} \{ \text{Cost}(\delta x_i, \delta u_i) + \mathbb{E}[J_{i+1}(\delta x_{i+1})] \} = \min_{u \in \mathcal{U}} \frac{1}{2} \begin{bmatrix} 1 \\ \delta x_i \\ \delta u_i \end{bmatrix}^T \begin{bmatrix} \hat{Q}_i & \hat{Q}_{x,i}^T & \hat{Q}_{u,i}^T \\ \hat{Q}_{x,i} & \hat{Q}_{xx,i} & \hat{Q}_{xu,i} \\ \hat{Q}_{u,i} & \hat{Q}_{ux,i} & \hat{Q}_{uu,i} \end{bmatrix} \begin{bmatrix} 1 \\ \delta x_i \\ \delta u_i \end{bmatrix} \quad (7)$$

with

$$\begin{aligned} \hat{Q}_i &= C_i + V_{i+1} + (1 + \bar{\phi}_i^T \Lambda^{-1} \bar{\phi}_i) \text{tr}(\Sigma_\epsilon V_{xx,i+1}), \\ \hat{Q}_{x,i} &= C_{x,i} + A_i^T \bar{K} V_{x,i+1} + (A_i^T \Lambda^{-1} \bar{\phi}_i) \odot \text{tr}(\Sigma_\epsilon V_{xx,i+1}), \\ \hat{Q}_{u,i} &= C_{u,i} + B_i^T \bar{K} V_{x,i+1} + (B_i^T \Lambda^{-1} \bar{\phi}_i) \odot \text{tr}(\Sigma_\epsilon V_{xx,i+1}), \\ \hat{Q}_{xx,i} &= C_{xx,i} + A_i^T \bar{K} V_{xx,i+1} \bar{K}^T A_i + (A_i^T \Lambda^{-1} A_i) \odot \text{tr}(\Sigma_\epsilon V_{xx,i+1}), \\ \hat{Q}_{uu,i} &= C_{uu,i} + B_i^T \bar{K} V_{xx,i+1} \bar{K}^T B_i + (B_i^T \Lambda^{-1} B_i) \odot \text{tr}(\Sigma_\epsilon V_{xx,i+1}), \\ \hat{Q}_{ux,i} &= \hat{Q}_{xu,i}^T = C_{ux,i} + B_i^T \bar{K} V_{xx,i+1} \bar{K}^T A_i + S_{xu,i}^T + (B_i^T \Lambda^{-1} A_i) \odot \text{tr}(\Sigma_\epsilon V_{xx,i+1}), \end{aligned}$$

where the  $C$  terms are the quadratization of the cost function. The last term in each equation captures the model uncertainty. Indeed, note that as  $\Lambda^{-1} \rightarrow 0$ , the effect of the model uncertainty decays to zero. Combining these expressions, the optimal control law may be written  $\delta u_i = l_i + L_i \delta x_i$ , where  $l_i = -\hat{Q}_{uu,i}^{-1} \hat{Q}_{u,i}$ , and  $L_i = -\hat{Q}_{uu,i}^{-1} \hat{Q}_{ux,i}$ . Finally, the value may be computed recursively via  $V_i = \hat{Q}_i - \hat{Q}_{u,i}^T \hat{Q}_{uu,i}^{-1} \hat{Q}_{u,i}$ ,  $V_{x,i} = \hat{Q}_{x,i} - \hat{Q}_{u,i}^T \hat{Q}_{uu,i}^{-1} \hat{Q}_{ux,i}$ , and  $V_{xx,i} = \hat{Q}_{xx,i} - \hat{Q}_{u,i}^T \hat{Q}_{uu,i}^{-1} \hat{Q}_{ux,i}$ .

<sup>1</sup>Considering the impact of control to parameter uncertainty is known as *dual control*, and is known to be extremely difficult, even in the approximate case [17].

Controller	Dyn.	Cost
Baseline	1-Step	519 $\pm$ 29.8
	N-Step	830 $\pm$ 182
CE	1-Step	567 $\pm$ 55.6
	N-Step	515.5 $\pm$ 32.7
Cautious	1-Step	605 $\pm$ 110
	N-Step	<b>506.5 <math>\pm</math> 31.0</b>

Controller	Dyn.	Failure Rate
Baseline	1-Step	0.10 (0.045, 0.165)
	N-Step	0.02 (0.000, 0.057)
CE	1-Step	0.26 (0.178, 0.349)
	N-Step	0.07 (0.023, 0.127)
Cautious	1-Step	0.13 (0.067, 0.201)
	N-Step	<b>0.00 (0.000, 0.022)</b>

Table 1: Performance (average cost) and failure rate on the quadrotor delivery task.

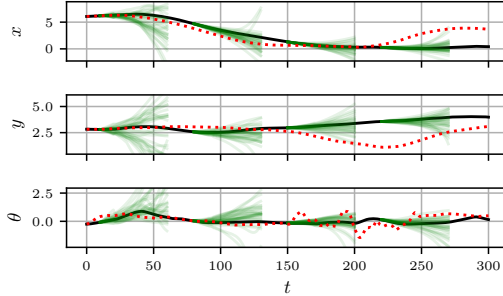


Figure 1: One sample trajectory of the quadrotor problem under the CAMELiD controller (black). The dotted red trajectory shows the baseline controller. The green trajectories are rollouts under the sequence of actions taken by the controller, under  $K$  sampled from the posterior.

**Algorithm Overview.** Using the multi-step training objective discussed previously, CAMELiD may be trained via maximum likelihood. The training objective closely follows that of [15], just with a multi-step objective. Because the approximate posterior predictive distribution is Gaussian, the feature weights  $w$  as well as the prior over the last layer may be learned via stochastic gradient descent.

During the online (control) phase, the uncertainty-aware iLQG controller may be used. In our experiments, we use the regularization approach and convergence criteria of [23], which substantially improved performance. As is standard for iLQG, CAMELiD consists of a forward pass in which the cost function is quadratized and the dynamics are linearized around the nominal trajectory. In the backward pass, the nominal control is updated via computing the cost-to-go terms and the associated control policy. To improve the efficiency of the online model update, a recursive last layer update is used [15].

## 4 Experimental Results

We demonstrate the capabilities of the CAMELiD framework on the physically intuitive domain of a planar quadrotor, a standard benchmark problem in control and reinforcement learning [21, 16, 10]. We consider a delivery task, where a quadrotor has the goal of moving to a pre-specified goal state, with a random mass attached at a random position along the quad. We encode the locomotion task through a cost function which encourages the agent to move to position  $(0, 3)$  while remaining upright.

We compare three control schemes on this task. We refer to the uncertainty-aware model we have described in the previous section as *Cautious* in our experimental results. The *Certainty Equivalent* (CE) model consists of the adaptive ALPaCA-based dynamics model without the additive  $S$  terms in the iterative LQG. As such, the CE approach matches the approach taken in [4], which used meta-learned dynamics models that provide only a point estimate of the posterior. The *Baseline* model is simply a dynamics model trained via domain randomization. It is trained to maximize log likelihood over all training data, without online adaptation. It is paired with a certainty-equivalent controller. We investigate these three approaches with both dynamics models that have been trained on one step prediction and multi-step prediction. The dynamics models were trained on a dataset of trajectories of the quadrotor maneuvering with different payloads and configurations.

The mean performance, together with 95% confidence intervals, of CAMELiD is reported in Table 1. As the high cost of crashing skewed the mean computations, we have separately reported the failure rate, and the mean cost incurred on successful runs. Note that the mean cost for the cautious controller includes trials on which the other controllers failed, and thus this is a penalty against the cautious method. In general, the certainty equivalent and cautious controllers outperform the non-adaptive controller. However, the certainty equivalent controller has a higher failure rate. By incorporating the model uncertainty into the control algorithm, the CAMELiD is able to maintain the low costs of an adaptive controller while also remaining robust. A representative rollout is shown in Figure 1.

## References

- [1] Andreas Argyriou, Theodoros Evgeniou, and Massimiliano Pontil. Multi-task feature learning. *Neural Information Processing Systems (NIPS)*, 2007.
- [2] Karl J Åström and Björn Wittenmark. *Adaptive control*. Courier Corporation, 2013.
- [3] Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. *Neural Information Processing Systems (NIPS)*, 2018.
- [4] Ignasi Clavera, Anusha Nagabandi, Ronald S Fearing, Pieter Abbeel, Sergey Levine, and Chelsea Finn. Learning to adapt: Meta-learning for model-based control. *arXiv:1803.11347*, 2018.
- [5] Marc Peter Deisenroth. Learning to control a low-cost manipulator using data-efficient reinforcement learning. *Robotics: Science and Systems (RSS)*, 2012.
- [6] Marc Peter Deisenroth, Dieter Fox, and Carl Edward Rasmussen. Gaussian processes for data-efficient learning in robotics and control. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 2013.
- [7] Andreas Doerr, Christian Daniel, Duy Nguyen-Tuong, Alonso Marco, Stefan Schaal, Toussaint Marc, and Sebastian Trimpe. Optimizing long-term predictions for model-based policy search. *Conference on Robot Learning*, 2017.
- [8] Chelsea Finn, Pieter Abbeel, and Sergey Levine. Model-agnostic meta-learning for fast adaptation of deep networks. *International Conference on Machine Learning (ICML)*, 2017.
- [9] Chelsea Finn and Sergey Levine. Deep visual foresight for planning robot motion. *IEEE International Conference on Robotics and Automation (ICRA)*, 2017.
- [10] Jeremy H Gillula, Haomiao Huang, Michael P Vitus, and Claire J Tomlin. Design of guaranteed safe maneuvers using reachable sets: Autonomous quadrotor aerobatics in theory and practice. *IEEE International Conference on Robotics and Automation (ICRA)*, 2010.
- [11] Agathe Girard, Carl Edward Rasmussen, J Quinonero-Candela, R Murray-Smith, O Winther, and J Larsen. Multiple-step ahead prediction for non linear dynamic systems—a Gaussian process treatment with propagation of the uncertainty. *Neural Information Processing Systems (NIPS)*, 15:529–536, 2002.
- [12] Shixiang Gu, Timothy Lillicrap, Ilya Sutskever, and Sergey Levine. Continuous deep q-learning with model-based acceleration. *International Conference on Machine Learning (ICML)*, 2016.
- [13] Abhishek Gupta, Russell Mendonca, YuXuan Liu, Pieter Abbeel, and Sergey Levine. Meta-reinforcement learning of structured exploration strategies. *Neural Information Processing Systems (NIPS)*, 2018.
- [14] James Harrison, Animesh Garg, Boris Ivanovic, Yuke Zhu, Silvio Savarese, Li Fei-Fei, and Marco Pavone. AdaPT: Zero-shot adaptive policy transfer for stochastic dynamical systems. *International Symposium on Robotics Research (ISRR)*, 2017.
- [15] James Harrison, Apoorva Sharma, and Marco Pavone. Meta-learning priors for efficient online bayesian regression. *Workshop on the Algorithmic Foundations of Robotics (WAFR)*, 2018.
- [16] Boris Ivanovic, James Harrison, Apoorva Sharma, Mo Chen, and Marco Pavone. BaRC: Backward reachability curriculum for robotic reinforcement learning. *arXiv:1806.06161*, 2018.
- [17] Edgar D Klenske and Philipp Hennig. Dual control for approximate bayesian reinforcement learning. *Journal of Machine Learning Research*, 2016.
- [18] Sergey Levine, Chelsea Finn, Trevor Darrell, and Pieter Abbeel. End-to-end training of deep visuomotor policies. *Journal of Machine Learning Research*, 2016.
- [19] Lennart Ljung and Torsten Söderström. *Theory and practice of recursive identification*. MIT press, 1983.

- [20] Nicolas Schweighofer and Kenji Doya. Meta-learning in reinforcement learning. *Neural Networks*, 2003.
- [21] Sumeet Singh, Anirudha Majumdar, Jean-Jacques Slotine, and Marco Pavone. Robust online motion planning via contraction theory and convex optimization. *IEEE International Conference on Robotics and Automation (ICRA)*, 2017.
- [22] Bradly C Stadie, Ge Yang, Rein Houthoofd, Xi Chen, Yan Duan, Yuhuai Wu, Pieter Abbeel, and Ilya Sutskever. Some considerations on learning to explore via meta-reinforcement learning. *arXiv:1803.01118*, 2018.
- [23] Yuval Tassa, Tom Erez, and Emanuel Todorov. Synthesis and stabilization of complex behaviors through online trajectory optimization. *IEEE International Conference on Intelligent Robots and Systems (IROS)*, 2012.
- [24] Emanuel Todorov and Weiwei Li. A generalized iterative lqg method for locally-optimal feedback control of constrained nonlinear stochastic systems. *American Control Conference (ACC)*, 2005.
- [25] Grady Williams, Nolan Wagener, Brian Goldfain, Paul Drews, James M Rehg, Byron Boots, and Evangelos A Theodorou. Information theoretic MPC for model-based reinforcement learning. *IEEE International Conference on Robotics and Automation (ICRA)*, 2017.