

Robust Dual Control based on Gain Scheduling

Janani Venkatasubramanian, Johannes Köhler, Julian Berberich, Frank Allgöwer

Abstract—We present a novel strategy for robust dual control of linear time-invariant systems based on gain scheduling with performance guarantees. This work relies on prior results of determining uncertainty bounds of system parameters estimated through exploration. Existing approaches are unable to account for changes of the mean of system parameters in the exploration phase and thus to accurately capture the *dual* effect. We address this limitation by selecting the future (uncertain) mean as a scheduling variable in the control design. The result is a semi-definite program-based design that computes a suitable exploration strategy and a robust gain-scheduled controller with probabilistic quadratic performance bounds after the exploration phase.

I. INTRODUCTION

The *dual control* paradigm established research interest in simultaneous learning and control of uncertain dynamic systems [1]. This pioneering work recognized that control inputs to an uncertain system have a ‘probing’ effect to learn the uncertainty in the system, and a ‘directing’ effect to control the dynamical system. However, these two effects are naturally conflicting, drawing attention to the trade-off between ‘exploration’ (learning system uncertainty) and ‘exploitation’ (controlling the system to achieve optimal performance), which is also the subject of contemporary literature on Reinforcement Learning [2]. Dual control relies on stochastic Dynamic Programming (DP) which is, however, computationally intractable. Either approximations of stochastic DP, or heuristic probing methods are typically adopted to solve the problem of tractability [3]. A detailed survey of dual control methods is provided in [4].

Early works of implicit dual control methods that involve approximations of DP are based on the *wide-sense* property [5], but they require linearization of system dynamics and approximation of the conditional probability of the states by its mean and covariance [6], [7]. These methods were extended to nonlinear systems that could handle input constraints, nonetheless based on some approximations [8]. This laid the foundation of balancing exploration with caution [9].

Explicit dual control methods use heuristic probing techniques for active learning without the need to introduce approximations of DP [10]. Explicit dual control methods are closely related to *Optimal Experiment Design* in closed

loop [11], [12]. These methods utilize the control inputs to regulate system dynamics and to probe the closed-loop system dynamics by solving a combined problem. This led to application-oriented strategies for dual control that promoted reducing uncertainty that would be beneficial for optimizing cost [13]. Some recent application-oriented strategies are discussed in [14], [15], however, they consider a special class of systems and their control strategies are not robust to model uncertainties. The ‘coarse-ID’ family of methods study robustness guarantees in system identification based design methods [16]–[18], however, the control policies are not optimized to balance exploration and exploitation.

Recently, in [19] a high probability uncertainty bound has been derived that is applicable to both robust control synthesis and *targeted* exploration. This bound is then used in dual control by predicting the influence of the controller on the future uncertainty. In particular, the work in [20], building on [19], [21], proposes a dual control strategy that minimizes worst-case cost attained by a robust controller that is synthesized with reduced model uncertainty. This dual control strategy with *targeted* exploration seems to perform better than strategies with common greedy random exploration. The approach in [22], further extending [19], retains an application-oriented strategy, however, adopts a more realistic finite horizon problem setting that captures the trade-offs between exploration and exploitation better.

The results presented in the methods in [19]–[22] seem to have lower conservatism compared to previous works [11]–[14], however, only numerically without any performance guarantees. In particular, the approaches in [19], [20] do not account for changes in the mean of uncertain system parameters during exploration. Therefore, this paper seeks to address these drawbacks by designing a dual control scheme based on gain scheduling. The mean of future uncertain system parameters is selected as a scheduling variable. This leads to a linear matrix inequality (LMI) based design with guarantees under relaxed assumptions. In particular, the resulting controller is a state feedback, which depends on the parameter estimates after exploration and thus on the data, and it *guarantees* robust closed-loop performance after an initial exploration phase.

The remainder of the paper is structured as follows. In Section II we state the problem setting, and in Section III we provide important results from the literature that we employ for our approach. Section IV contains the proposed dual controller design procedure as well as a proof of robust closed-loop performance guarantees. Finally, we conclude the paper in Section V.

The authors are with the Institute for Systems Theory and Automatic Control, University of Stuttgart, 70550 Stuttgart, Germany. (email: {janani.venkatasubramanian, johannes.koehler, julian.berberich, frank.allgower}@ist.uni-stuttgart.de)

This work was funded by Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy - EXC 2075 - 390740016. The authors thank the International Max Planck Research School for Intelligent Systems (IMPRS-IS) for supporting Janani Venkatasubramanian and Julian Berberich, and the International Research Training Group Soft Tissue Robotics (GRK 2198/1).

II. PROBLEM STATEMENT

Notation: The transpose of a matrix A is denoted as A^\top . The value of the Chi-squared distribution with n degrees of freedom and probability p is denoted as $\chi_n^2(p)$. \mathcal{L}_2 denotes the space of square-summable functions.

Setting: Consider a discrete-time linear time-invariant dynamical system of the form

$$x_{k+1} = A_{tr}x_k + B_{tr}u_k + w_k, \quad w_k \sim \mathcal{N}(0, \sigma_w^2 I) \quad (1)$$

where $x_k \in \mathbb{R}^{n_x}$ is the state, $u_k \in \mathbb{R}^{n_u}$ is the control input, and $w_k \in \mathbb{R}^{n_x}$ is the normally distributed process noise. The true values of the system dynamics, A_{tr} and B_{tr} , are unknown.

Control goal - proposed approach: The main goal is to design a stabilizing state-feedback $u_k = K_{\text{new}}x_k$ which meets some desired performance specifications. Since the system is unknown, we first apply some exciting input to estimate the parameters and then use bounds on the estimation error to design a robust controller. Our goal is to simultaneously design a suitable exploration strategy and a controller for the system in (1), such that applying the feedback after the exploration phase provides desired quadratic performance guarantees with high probability. The main challenge is to accurately capture the *dual* effect of performance improvement through exploration. We solve this problem by interpreting the new parameter estimate as a scheduling variable, which influences the control law K_{new} , using tools from gain-scheduling. The corresponding necessary preliminaries regarding uncertainty bounds for parameter estimation, gain-scheduling and structured exploration are shown in Section III and the overall approach is presented in Section IV.

III. PRELIMINARIES

A. Uncertainty Bound

This subsection discusses preliminary results from [19] adopted in our work that quantify uncertainty in the system dynamics that are estimated, given some data. The unknown matrices A_{tr} and B_{tr} can be estimated through observed data $\mathcal{D} = \{x_k, u_k\}_{k=0}^N$ of length N . In particular, we consider the least squares estimates of A_{tr} and B_{tr} , similar to [19], which are given by,

$$(\hat{A}, \hat{B}) = \arg \min_{A, B} \sum_{k=0}^{N-1} \|x_{k+1} - Ax_k - Bu_k\|_2^2. \quad (2)$$

The following lemma provides a high probability credibility region for the uncertain system matrices.

Lemma 1. [19, Lemma 3.1] *Given data set \mathcal{D} and $0 < \delta < 1$, let $D = \frac{1}{\sigma_w^2 c_\delta} \sum_{k=1}^{N-1} \begin{bmatrix} x_k \\ u_k \end{bmatrix} \begin{bmatrix} x_k \\ u_k \end{bmatrix}^\top$ with $c_\delta = \chi_{n_x^2 + n_x n_u}^2(\delta)$. Suppose we have a uniform prior over the parameters (A, B) . Then, $[A_{tr}, B_{tr}] \in \Theta$ with probability $1 - \delta$, where*

$$\Theta := \left\{ A, B : \begin{bmatrix} (\hat{A} - A)^\top \\ (\hat{B} - B)^\top \end{bmatrix}^\top D \begin{bmatrix} (\hat{A} - A)^\top \\ (\hat{B} - B)^\top \end{bmatrix} \preceq I \right\}. \quad (3)$$

This lemma provides a data-dependent uncertainty bound. Given this uncertainty bound, the approaches in [19], [20], [22] synthesize a robust controller by minimizing a worst-case cost. This controller facilitates *targeted* exploration for dual control strategies by predicting the future uncertainty, depending on the exploring controller. However, their approach does not take into account that the estimate of the system parameters are subject to change through the process of exploration.

B. Gain Scheduling Approach

To account for the change in the estimates of the system parameters, we model the system in (1) as a linear parameter varying (LPV) system. The varying system parameters can be measured after exploration and are selected as the *scheduling block*. The goal is to design a gain-scheduling controller that ensures that the closed-loop system is stable while also satisfying a quadratic performance bound, e.g. \mathcal{L}_2 gain, from the disturbance input w to the performance output z with high probability, compare [23], [24]. The performance specification is imposed on the channel $w \rightarrow z$, where the performance output z_k at time k is the generalized error that depends on the state, control input and disturbance:

$$z_k = Cx_k + Du_k + D_w w_k, \quad (4)$$

where C , D and D_w are known. In this setup, since the dynamics are unknown, we have the following assumption from which an initial error bound of the form given in Lemma 1 can be derived.

Assumption 1. *An initial data set $\mathcal{D}_0 = \{x_t, u_t\}_{t=-N_0}^{-1}$ is available and a uniform prior over the parameters (A, B) is assumed. Moreover, it holds that*

$$D_0 := \frac{1}{\sigma_w^2 c_\delta} \sum_{k=-N_0}^{-1} \begin{bmatrix} x_k \\ u_k \end{bmatrix} \begin{bmatrix} x_k \\ u_k \end{bmatrix}^\top \succ 0. \quad (5)$$

From the data \mathcal{D}_0 , initial estimates of the system parameters \hat{A}_0 and \hat{B}_0 can be derived. The matrix D_0 quantifies the uncertainty associated with these initial estimates for a given probability $1 - \delta$ and can be determined from \mathcal{D}_0 as given in (5). This initial data can be acquired through some random persistently exciting input, while the later exploration will use the initially obtained model knowledge to provide a more *targeted* exploration strategy. Through the exploration process for T time steps, data $\mathcal{D}_T = \{x_t, u_t\}_{t=0}^T$ will be observed. The new estimates \hat{A}_T and \hat{B}_T will be computed from data $\mathcal{D}_0 \cup \mathcal{D}_T$ and made available at time T . The matrix $D_T := D_0 + \frac{1}{\sigma_w^2 c_\delta} \sum_{k=0}^{T-1} \begin{bmatrix} x_k \\ u_k \end{bmatrix} \begin{bmatrix} x_k \\ u_k \end{bmatrix}^\top$ will quantify the uncertainty associated with the estimates \hat{A}_T and \hat{B}_T . Existing approaches such as [20] rely on the assumption that $\hat{A}_0 \approx \hat{A}_T$ and $\hat{B}_0 \approx \hat{B}_T$. In the following, we propose a gain scheduling-based approach to provide closed-loop guarantees for the case $\hat{A}_0 \neq \hat{A}_T$ and $\hat{B}_0 \neq \hat{B}_T$. Since the system parameters will be updated through the process

of exploration, we proceed now by rewriting (1) as,

$$\begin{aligned} x_{k+1} &= A_{tr}x_k + B_{tr}u_k + w_k \\ &= \hat{A}_0x_k + \hat{B}_0u_k + (\hat{A}_T - \hat{A}_0)x_k + (\hat{B}_T - \hat{B}_0)u_k \\ &\quad + (A_{tr} - \hat{A}_T)x_k + (B_{tr} - \hat{B}_T)u_k + w_k. \end{aligned} \quad (6)$$

From (6), the scheduling and uncertainty blocks can be selected as,

$$\begin{aligned} \Delta_s &= [\hat{A}_T - \hat{A}_0 \quad \hat{B}_T - \hat{B}_0], \\ \Delta_u &= [A_{tr} - \hat{A}_T \quad B_{tr} - \hat{B}_T]. \end{aligned} \quad (7)$$

Since the estimates at time T affect both Δ_s and Δ_u , the latter blocks can be viewed as time-varying parameters, and the uncertain system combining (4) and (6) can be written as an LPV system:

$$\begin{aligned} \begin{bmatrix} x_{k+1} \\ z_k^s \\ z_k^u \\ z_k \end{bmatrix} &= \begin{bmatrix} \hat{A}_0 & I & I & I & \hat{B}_0 \\ I & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & I \\ 0 & 0 & 0 & 0 & I \\ C & 0 & 0 & D_w & D \end{bmatrix} \begin{bmatrix} x_k \\ w_k^s \\ w_k^u \\ w_k \\ u_k \end{bmatrix}, \\ w_k^s &= \Delta_s z_k^s, \\ w_k^u &= \Delta_u z_k^u, \end{aligned} \quad (8)$$

where $w^s \rightarrow z^s$ is the scheduling channel and $w^u \rightarrow z^u$ is the uncertainty channel. After the exploration phase, the control input can now be defined as

$$u_k = Kx_k + K_s w_k^s. \quad (9)$$

The goal is to design K and K_s such that the closed-loop system is stable and the specified performance criterion is met. The robust gain-scheduling configuration is illustrated in Figure 1. As can be seen in Figure 1, the open-loop system

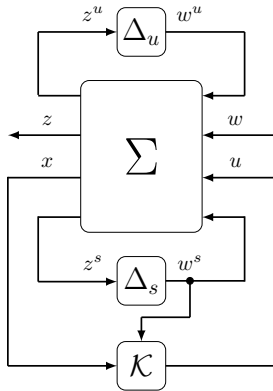


Fig. 1. Generalized plant view of the robust gain-scheduling problem.

has two uncertainty channels, affected by Δ_u and Δ_s . The latter uncertainty Δ_s is taken into account for the controller via w_k^s , compare (9), and hence plays the role of a scheduling variable. This accounts for changes in the mean of the system parameters through data gathered in the exploration phase that is available after exploration at time T , and thereby

learning the system dynamics better. The closed-loop system can be written as

$$\begin{aligned} \begin{bmatrix} x_{k+1} \\ z_k^s \\ z_k^u \\ z_k \end{bmatrix} &= \begin{bmatrix} \hat{A}_0 + \hat{B}_0 K & I + \hat{B}_0 K_s & I & I \\ I & 0 & 0 & 0 \\ K & K_s & 0 & 0 \\ C + DK & DK_s & 0 & D_w \end{bmatrix} \begin{bmatrix} x_k \\ w_k^s \\ w_k^u \\ w_k \end{bmatrix}, \\ w_k^s &= \Delta_s z_k^s, \\ w_k^u &= \Delta_u z_k^u. \end{aligned} \quad (10)$$

Given this formulation and suitable bounds on the blocks Δ_s and Δ_u , we can use established methods from robust control and gain-scheduling to guarantee a desired performance specification. In particular, we consider the case where a desired quadratic performance specification is given on the performance channel $w_k \rightarrow z_k$, i.e. for initial condition $x = 0$ for all signals $w \in \mathcal{L}_2$ with output z of the closed loop, the following inequality should hold with some $\epsilon > 0$:

$$\sum_{k=0}^{\infty} \begin{pmatrix} w_k \\ z_k \end{pmatrix}^\top \begin{pmatrix} Q_p & S_p \\ S_p^\top & R_p \end{pmatrix} \begin{pmatrix} w_k \\ z_k \end{pmatrix} \leq -\epsilon \sum_{k=0}^{\infty} w_k^\top w_k, \quad (11)$$

where $R_p \succ 0$ is assumed. We note that standard design goals, such as a desired \mathcal{L}_2 -gain of γ are contained as a special case with $S_p = 0$, $R_p = \frac{1}{\gamma}I$ and $Q_p = -\gamma I$ (c.f. [25, Prop. 3.12]). The following lemma provides a matrix inequality to design a robust gain scheduling controller satisfying such a performance specification, given suitable bounds on the blocks Δ_s, Δ_u .

Lemma 2. Suppose $\Delta_s \in \mathbf{\Delta}_s := \{\Delta : \Delta^\top Q_s \Delta + R_s \succ 0\}$, $\Delta_u \in \mathbf{\Delta}_u := \{\Delta : \Delta^\top Q_u \Delta + R_u \succ 0\}$, with $R_u, R_s \succ 0$. If there exists matrices K_s, M, N and scalars $\lambda_s, \lambda_u > 0$ satisfying the matrix inequality (13), the closed loop (10) satisfies the quadratic performance bound (11) with $K = MN^{-1}$, i.e., $u_k = MN^{-1}x_k + K_s w_k^s$.

Proof. The proof follows the arguments in [23] for LPV control. First, note that the set definitions $\mathbf{\Delta}_s, \mathbf{\Delta}_u$ are linear in Q_s, Q_u, R_s, R_u and thus remain valid if (Q_s, R_s) and (Q_u, R_u) are multiplied by some positive scalar $\lambda_s, \lambda_u > 0$, respectively. Define $X = N^{-1}$ and $K = MN^{-1}$. The Schur complement of the LMI (13) is multiplied from left and right by $\text{diag}(N^{-1}, I, I, I)$ to obtain

$$\begin{bmatrix} * \\ * \\ * \\ * \\ * \\ * \\ * \end{bmatrix}^\top \begin{bmatrix} -X & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & X & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & \lambda_s P_s & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & \lambda_u P_u & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & Q_p & S_p \\ 0 & 0 & 0 & 0 & 0 & 0 & S_p^\top & R_p \end{bmatrix} \times \begin{bmatrix} I & 0 & 0 & 0 \\ \hat{A}_0 + \hat{B}_0 K & I + \hat{B}_0 K_s & I & I \\ 0 & I & 0 & 0 \\ \begin{bmatrix} I \\ K \end{bmatrix} & \begin{bmatrix} 0 \\ K_s \end{bmatrix} & 0 & 0 \\ 0 & 0 & I & 0 \\ \begin{bmatrix} I \\ K \end{bmatrix} & \begin{bmatrix} 0 \\ K_s \end{bmatrix} & 0 & 0 \\ 0 & 0 & 0 & I \\ C + DK & DK_s & 0 & D_w \end{bmatrix} \prec 0, \quad (12)$$

where $P_s = \begin{bmatrix} Q_s & 0 \\ 0 & R_s \end{bmatrix}$ and $P_u = \begin{bmatrix} Q_u & 0 \\ 0 & R_u \end{bmatrix}$.

Using [23, Thm. 2], quadratic performance is guaranteed if there exists a positive definite matrix $X = X^\top \succ 0$ satisfying the matrix inequality (12). \square

We note that for λ_s, λ_u constant, inequality (13) is an LMI and thus can be efficiently solved using line-search like techniques for $(\lambda_s, \lambda_u) \in \mathbb{R}^2$.

C. Exploration and parameter estimation bounds

In this paper, we consider a dual control objective where, during an initial exploration phase, uncertainty is reduced in order to design a robust controller based on the data collected during exploration. The exploration controller is computed such that it excites the system sufficiently with a minimal *robust* LQR cost, based on initial parameter estimates. The exploration controller takes the form

$$u_k = K_e x_k + e_k, \quad k = 0, \dots, T \quad (14)$$

with a robustly stabilizing K_e and noise $e_k \sim \mathcal{N}(0, \Sigma)$. Based on initial estimates of the system dynamics and the associated uncertainty bound, K_e and Σ are computed such that they minimize a robust LQR cost $\sum_{k=0}^{\infty} x_k^\top Q x_k + u_k^\top R u_k$. Similar to [20], this robust LQR cost can be computed as the \mathcal{H}_2 -norm of the uncertain closed loop system $x_{k+1} = (\hat{A}_0 + \hat{B}_0 K_e)x_k + w_k$ and $y_k = \begin{bmatrix} Q^{\frac{1}{2}} \\ R^{\frac{1}{2}} K_e \end{bmatrix} x_k$. To be more precise, the robust LQR cost of the exploration controller is computed as

$$\begin{aligned} & \min_{t_e, Z_e, Y_e, W_e} \quad \text{tr } Y_e \\ & \text{s.t.} \quad S_1(W_e, Y_e, Z_e) \succeq 0, \quad t_e > 0 \\ & \quad S_e(t_e, Z_e, W_e, \Sigma, D_0, \hat{A}_0, \hat{B}_0) \succeq 0 \end{aligned} \quad (15)$$

where S_1, S_e are defined as

$$S_1(W_e, Y_e, Z_e) = \left[\begin{array}{c|c} Y_e & Q^{\frac{1}{2}} W_e \\ \hline W_e Q^{\frac{1}{2}} & R^{\frac{1}{2}} Z_e^\top \end{array} \right],$$

$$S_e(t_e, Z_e, W_e, \Sigma) = \begin{bmatrix} H_e & F_e & G_e \\ F_e^\top & C_e - t_e I & 0 \\ G_e^\top & 0 & t_e D_0 \end{bmatrix}$$

with

$$H_e = \begin{bmatrix} W_e & 0 \\ 0 & \Sigma \end{bmatrix}, \quad F_e = \begin{bmatrix} W_e \hat{A}_0^\top + Z_e \hat{B}_0^\top \\ \Sigma \hat{B}_0^\top \end{bmatrix},$$

$$G_e = \begin{bmatrix} -W_e & -Z_e \\ 0 & -\Sigma \end{bmatrix}, \quad Z_e = W_e K_e^\top, \quad C_e = W_e - \sigma_w^2 I.$$

A more detailed derivation and additional explanations are provided in [20]. In (15), W_e denotes the controllability Gramian, which plays an essential role to propagate the influence of the exploration phase on the parameters estimates based on the new data \mathcal{D}_T . Similar to [20], we make the following assumption.

Assumption 2. For the system (1) evolving under an exploration controller (14) with K_e, W_e satisfying the constraints in (15), the empirical covariance can be approximated via a solution W_e of (15) as

$$\sum_{t=0}^{T-1} \begin{bmatrix} x_t \\ u_t \end{bmatrix} \begin{bmatrix} x_t \\ u_t \end{bmatrix}^\top \approx T \begin{bmatrix} W_e & W_e K_e^\top \\ K_e W_e & K_e W_e K_e^\top + \Sigma \end{bmatrix}. \quad (16)$$

Assumption 2 implies that the empirical covariance can be approximated via its stationary distribution $\Sigma_{xx} = \mathbb{E}[xx^\top]$, which is in turn approximated by the worst-case state covariance, i.e., by W_e satisfying (15). Clearly, this is an assumption that is not guaranteed to hold, but it is usually a good approximation. Based on (16), the uncertainty bound D_T can be computed as

$$D_T = D_0 + \frac{T}{\sigma_w^2 c_\delta} \begin{bmatrix} W_e & Z_e \\ Z_e^\top & Z_e^\top W_e^{-1} Z_e + \Sigma \end{bmatrix}. \quad (17)$$

Since the uncertainty bound D_T is a nonlinear function of Z_e and W_e , we compute an affine lower bound of it as in [20, Lemma 1]

$$\begin{bmatrix} W_e & Z_e \\ Z_e^\top & Z_e^\top W_e^{-1} Z_e \end{bmatrix} \succeq \begin{bmatrix} W_e \\ Z_e^\top \end{bmatrix} V + V^\top \begin{bmatrix} W_e \\ Z_e^\top \end{bmatrix}^\top - V^\top W_e V. \quad (18)$$

For a fixed V , this leads to an affine lower bound on D_T , denoted as \bar{D}_T , via (17). The bound is tight when $\begin{bmatrix} W_e & Z_e \end{bmatrix} = W_e V$, so it is optimal to choose $V = W_e^{-1} \begin{bmatrix} W_e & Z_e \end{bmatrix} = \begin{bmatrix} I & K_e^\top \end{bmatrix}$. However, K_e is not known at this point, and hence a candidate K_0 is used instead to compute V , which can be computed, e.g., based on a robust LQR for the nominal model [20].

An essential ingredient of the proposed approach is the handling of the uncertainty bounds D_0, D_T , which influence the uncertain parameters Δ_u, Δ_s in (10), as derived in the following proposition.

Proposition 1. Let Assumption 1 hold, where \hat{A}_0 and \hat{B}_0 are the initial estimates. Let

$$\Delta_0 = [\hat{A}_0 - A_{tr} \quad \hat{B}_0 - B_{tr}].$$

$$\left(\begin{array}{c|c} \begin{bmatrix} -N & 0 & 0 & (CN + DM)^\top S_p^\top \\ 0 & \lambda_s Q_s & 0 & (DK_s)^\top S_p^\top \\ 0 & 0 & \lambda_u Q_u & 0 \\ S_p(CN + DM) & S_p DK_s & 0 & Q_p + D_w^\top S_p^\top + S_p D_w \end{bmatrix} & \star \\ \hline \begin{bmatrix} \hat{A}_0 N + \hat{B}_0 M & I + \hat{B}_0 K_s & I & I \\ \begin{bmatrix} N \\ M \end{bmatrix} & \begin{bmatrix} 0 \\ K_s \end{bmatrix} & 0 & 0 \\ \begin{bmatrix} N \\ M \end{bmatrix} & \begin{bmatrix} 0 \\ K_s \end{bmatrix} & 0 & 0 \\ CN + DM & DK_s & 0 & D_w \end{bmatrix} & \begin{bmatrix} -N & 0 & 0 & 0 \\ 0 & -\frac{1}{\lambda_s} R_s^{-1} & 0 & 0 \\ 0 & 0 & -\frac{1}{\lambda_u} R_u^{-1} & 0 \\ 0 & 0 & 0 & -R_p^{-1} \end{bmatrix} \end{array} \right) \prec 0. \quad (13)$$

Then, with probability $1 - \delta$, we have

$$\Delta_0 \in \mathbf{\Delta}_0 := \{\Delta_0 : \Delta_0^\top \Delta_0 \preceq D_0^{-1}\}, \quad (19)$$

and with probability $1 - \delta$, we have

$$\Delta_u \in \mathbf{\Delta}_u := \{\Delta_u : \Delta_u^\top \Delta_u \preceq D_T^{-1}\}. \quad (20)$$

If (19) and (20) hold, then for any $\epsilon > 0$ we have

$$\Delta_s \in \mathbf{\Delta}_s = \left\{ \Delta_s : \Delta_s^\top \Delta_s \preceq \left(1 + \frac{1}{\epsilon}\right) D_0^{-1} + (1 + \epsilon) D_T^{-1} \right\}. \quad (21)$$

Proof. Let

$$\Delta_0 = [\hat{A}_0 - A_{tr} \quad \hat{B}_0 - B_{tr}]. \quad (22)$$

By Lemma 1, with $0 < \delta < 1$, the following hold with probability $1 - \delta$.

$$\begin{aligned} \Delta_0^\top \Delta_0 &\preceq D_0^{-1}, \\ \Delta_u^\top \Delta_u &\preceq D_T^{-1}. \end{aligned} \quad (23)$$

The scheduling block can now be represented as

$$\Delta_s = -(\Delta_0 + \Delta_u). \quad (24)$$

To derive a probabilistic bound for Δ_s , we have

$$\begin{aligned} \Delta_s^\top \Delta_s &= (\Delta_0 + \Delta_u)^\top (\Delta_0 + \Delta_u) \\ &= \Delta_0^\top \Delta_0 + \Delta_0^\top \Delta_u + \Delta_u^\top \Delta_0 + \Delta_u^\top \Delta_u \\ &\preceq \left(1 + \frac{1}{\epsilon}\right) \Delta_0^\top \Delta_0 + (1 + \epsilon) \Delta_u^\top \Delta_u. \end{aligned} \quad (25)$$

The third inequality follows by Young's inequality which implies that for every $\epsilon > 0$, $\Delta_0^\top \Delta_u + \Delta_u^\top \Delta_0 \leq \left(\frac{1}{\epsilon}\right) \Delta_0^\top \Delta_0 + \epsilon \Delta_u^\top \Delta_u$. Therefore, the bound for Δ_s is

$$\Delta_s^\top \Delta_s \preceq \left(1 + \frac{1}{\epsilon}\right) D_0^{-1} + (1 + \epsilon) D_T^{-1}. \quad \square$$

The relation between the different sets is visualized in Figure 2. Using Lemma 1 with the initial data \mathcal{D}_0 , we know that the true system parameters θ_{tr} are in some ellipse $\mathbf{\Delta}_0$ around the initial parameter estimate $\hat{\theta}_0$. Using Lemma 1 after the exploration, we know that the true parameter θ_{tr} is contained in an ellipse $\mathbf{\Delta}_u$ around the new point estimate $\hat{\theta}_T$. Combining both of these bounds we know that the new point estimate $\hat{\theta}_T$, and thus the scheduling variable Δ_s , is

contained in an ellipse $\mathbf{\Delta}_s$ around the initial point estimate.

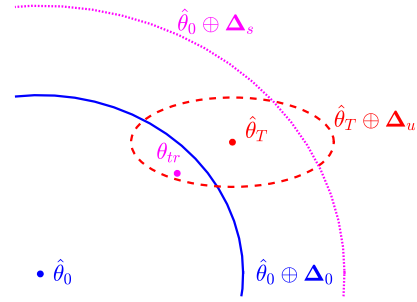


Fig. 2. Illustration of the sets $\mathbf{\Delta}_0, \mathbf{\Delta}_s, \mathbf{\Delta}_u$ (shifted w.r.t $\hat{\theta}_0, \hat{\theta}_T$) with the true parameters θ_{tr} , the initial parameter estimate $\hat{\theta}_0$ and the estimate resulting from the exploration $\hat{\theta}_T$.

IV. DUAL CONTROL

In this section, the proposed robust dual control strategy is presented. The overall proposed algorithm is summarized in Algorithm 1 in Section IV-A and its theoretical properties are analyzed in Theorem 1 in Section IV-B. Finally, Section IV-C discusses the approach relative to existing methods.

A. Proposed Algorithm

The overall goal of the proposed dual approach is to design a structured exploration strategy, such that the designed controller satisfies some desired quadratic performance (11) with high probability $1 - \delta \in (0, 1)$. Since we assume that we do not have a prior on the model, we start with some random (ideally persistently exciting) exploration over $N_0 \in \mathbb{N}$ steps to obtain initial data \mathcal{D}_0 (Assumption 1). Then we use the least mean squares estimate (2) to obtain initial estimates \hat{A}_0, \hat{B}_0 and a high probability uncertainty bound D_0^{-1} (c.f. Lemma 1 and (19) in Prop. 1). Before solving the dual control problem, we first seek some initial candidate feedback $K_0 \in \mathbb{R}^{m \times n}$. This step is in principle not necessary, however, in case $K_0 - K_e$ is small, this step can greatly reduce the conservatism in the convex relaxation (18). Thus, as also suggested in [20, Sec. IV.C], we compute K_0 as a robust LQR for the nominal model (15).

In order to emphasize its dependence on the variables and uncertainty parameters, we denote the matrix in (13) with

$Q_s = -I, Q_u = -I$ by

$$S_2(K_s, M, N, \lambda_s, \lambda_u, R_s^{-1}, R_u^{-1}).$$

Furthermore, the satisfaction of (13) is ensured if $\Delta_s^\top \Delta_s \prec D_s^{-1}$, $\Delta_u^\top \Delta_u \prec D_T^{-1}$, where D_s denotes the uncertainty bound associated with Δ_s . Since the set inclusions (19)–(21) from Proposition 1 hold, it suffices to show

$$D_s^{-1} \succ \frac{1+\epsilon}{\epsilon} D_0^{-1} + (1+\epsilon) D_T^{-1}, \quad (26)$$

By applying the Woodbury matrix identity to (26) and multiplying by $(1+\epsilon)$, we get

$$(1+\epsilon) D_s \succ \epsilon D_0 - \epsilon D_0 (D_T + \epsilon D_0)^{-1} \epsilon D_0. \quad (27)$$

Applying the Schur complement to (27) results in the following equivalent LMI

$$S_3(\epsilon, D_0, D_T, D_s) = \begin{bmatrix} \epsilon D_0 - (1+\epsilon) D_s & \epsilon D_0 \\ \epsilon D_0 & D_T + \epsilon D_0 \end{bmatrix}.$$

Given $K_0, D_0, \hat{A}_0, \hat{B}_0, \delta, Q_p, S_p, R_p, T$ and some fixed $\epsilon, t_e, \lambda_s, \lambda_u > 0$, we solve the following semi-definite program (SDP), which is a combination of the robust gain-scheduling problem (13) and the exploration inequalities (15), (17):

$$\inf_{\substack{W_e, Z_e, Y_e, \Sigma \\ K_s, M, N, \bar{D}_T, D_s}} \text{tr}(Y_e) \quad (28a)$$

$$\text{s.t. } S_1(W_e, Y_e, Z_e) \succeq 0 \quad (28b)$$

$$S_e(t_e, Z_e, W_e, \Sigma, D_0, \hat{A}_0, \hat{B}_0) \succeq 0 \quad (28c)$$

$$S_2(K_s, M, N, \lambda_s, \lambda_u, D_s, \bar{D}_T) \prec 0 \quad (28d)$$

$$S_3(\epsilon, D_0, \bar{D}_T, D_s) \succ 0 \quad (28e)$$

$$\frac{T}{\sigma_w^2 c_\delta} \begin{bmatrix} W_e & Z_e \\ Z_e^\top & Z_e^\top K_0^\top + K_0 Z_e - K_0 W_e K_0^\top + \Sigma \end{bmatrix} + D_0 - \bar{D}_T \succ 0. \quad (28f)$$

Solving this optimization problem directly leads to the controller parameters required for the implementation, i.e., the exploration controller $K_e = Z_e^\top W_e^{-1}$, the exploration variance Σ , and the robust gain scheduled controller parameters K_s and $K = MN^{-1}$. Essentially, (28b)–(28c) are needed to compute the cost $\text{tr}(Y_e)$ of the controller during the exploration phase, compare (15). Moreover, (28d) contains the main robust control LMI (compare Lemma 2) which returns a common Lyapunov function $N \succ 0$ as well as controller parameters M, K_s which guarantee robust performance of the closed loop (10) for all uncertainties Δ_u, Δ_s satisfying $\Delta_u^\top \Delta_u \prec D_T^{-1}$ and $\Delta_s^\top \Delta_s \prec D_s^{-1}$. In this context, (a bound on) the data obtained during exploration is approximated via $\bar{D}_T - D_0$, which implies that the uncertainties for robust controller design, i.e., the values \bar{D}_T and D_s , in turn depend on the controller during the exploration phase K_e through (28e) and (28f). This couples the exploration and robust control, thus resulting in a *dual effect* of the proposed controller.

Regarding the computational complexity of (28), we note that for $\epsilon, t_e, \lambda_s, \lambda_u > 0$ fixed, this is a standard (small-scale) semi-definite program (SDP), which can be efficiently

solved. Hence, the optimization problem can be solved by using a line-search like procedure (or gridding) for the variables $\epsilon, t_e, \lambda_s, \lambda_u > 0$ and solving the SDP in an inner loop.

After solving (28), we apply the targeted exploration sequence $u_t = K_e x_t + e_t$, $e_t \sim \mathcal{N}(0, \Sigma)$ for $t = 0, \dots, T$. Next, with the new data $\mathcal{D}_0 \cup \mathcal{D}_T$, we use the least mean square estimate (2) to obtain an improved/updated estimate \hat{A}_T, \hat{B}_T and a new bound D_T^{-1} on the uncertainty. Then, we can directly apply the designed gain-scheduling controller with the new scheduling variable $\Delta_s = (\hat{A}_T - \hat{A}_0 \quad \hat{B}_T - \hat{B}_0)$. Using (10), this controller can be explicitly written as a state feedback control law K_{new} using

$$\begin{aligned} u_k &= K x_k + K_s w_k^s \\ &= K x_k + K_s ((\hat{A}_T - \hat{A}_0) x_k + (\hat{B}_T - \hat{B}_0) u_k) \\ &= (I_m - K_s (\hat{B}_T - \hat{B}_0))^{-1} (K + K_s (\hat{A}_T - \hat{A}_0)) x_k \\ &=: K_{\text{new}} x_k. \end{aligned} \quad (29)$$

We note that $(I - K_s (\hat{B}_T - \hat{B}_0))$ is non-singular (with high probability) due to the equivalence in [23, Thm. 2]. The overall procedure is summarized in Algorithm 1.

Algorithm 1 Dual control using gain-scheduling

- 1: Specify confidence level $\delta \in (0, 1)$, quadratic performance (Q_p, S_p, R_p) (11), exploration cost $Q, R \succ 0$, initial and targeted exploration length N_0, T .
 - 2: Random exploration to obtain initial data \mathcal{D}_0 (Ass. 1).
 \Rightarrow Initial estimates \hat{A}_0, \hat{B}_0 and uncertainty bound D_0^{-1} , compute robust LQR controller K_0 (15).
 - 3: Solve the optimization problem (28) for different values $\epsilon, t_e, \lambda_s, \lambda_u > 0$ (e.g., via line-search in an outer loop).
 \Rightarrow Exploration sequence $K_e = Z^\top W^{-1}$, Σ and gain-scheduled controller $K_s, K = MN^{-1}$.
 - 4: Apply the exploration input $u_k = K_e x_k + e_k$, $e_k \sim \mathcal{N}(0, \Sigma)$ for $k = 0, \dots, T$.
 - 5: Update estimates \hat{A}_T, \hat{B}_T using new data.
 - 6: Compute the equivalent state-feedback K_{new} and apply the feedback $u_k = K_{\text{new}} x_k, k > T$.
-

B. Theoretical analysis

The following result proves that Algorithm 1 leads to a controller with closed-loop guarantees.

Theorem 1. *Let Assumptions 1–2 hold, suppose (28) is feasible and Algorithm 1 is applied. Assume further that the set inclusions (19)–(21) from Proposition 1 hold. Then the state-feedback K_{new} from (29) is well-defined and satisfies the quadratic performance bound (11).*

Proof. First, we recap that Lemma 2 guarantees the performance bound (11), assuming suitable bounds on Δ_s, Δ_u . Then, we show that exploration inequalities in combination with Assumption 2 ensure the bounds on Δ_s, Δ_u .

Part I. According to Lemma 2, satisfaction of the matrix inequality (28d) guarantees that the robust gain-scheduling controller $u = MN^{-1} x_k + K_s w_k$ ensures the quadratic

performance bound (11), if $\Delta_s^\top \Delta_s \prec D_s^{-1}$, $\Delta_u^\top \Delta_u \prec \bar{D}_T^{-1}$. Moreover, it is a direct consequence of the synthesis LMI that K_{new} is well-posed. Thus, it only remains to show that $\Delta_s^\top \Delta_s \prec D_s^{-1}$, $\Delta_u^\top \Delta_u \prec \bar{D}_T^{-1}$.

Part II. Since the set inclusions (19)–(21) from Proposition 1 hold, it suffices to show $\bar{D}_T^{-1} \succ D_T^{-1}$ and

$$D_s^{-1} \succ \frac{1+\epsilon}{\epsilon} D_0^{-1} + (1+\epsilon) \bar{D}_T^{-1}, \quad (30)$$

Assumption 2 ensures that the bound (17) holds. The convex relaxation (18) (c.f. [20, Lemma 1]) in combination with inequality (28f) ensures that $D_T \succ \bar{D}_T$ and thus $\Delta_u^\top \Delta_u \prec \bar{D}_T^{-1}$. Finally, as shown earlier, inequality (28e) is equivalent to (30), which implies $\Delta_s^\top \Delta_s \prec D_s^{-1}$. \square

We point out that, since the properties in Proposition 1 only hold with some probability $1 - \delta$, the quadratic performance (11) only holds with some probability, which is inherent in the considered stochastic/Gaussian setup.

C. Discussion

The proposed method detailed in Algorithm 1 combines structured exploration techniques as developed in [19], [20] and robust gain scheduling controller design. Given an initial data set (compare Assumption 1) and a quadratic performance specification Q_p, S_p, R_p on the channel $w \mapsto z$, Theorem 1 implies that Algorithm 1 guarantees robust performance for the closed loop with input $u_k = K_{\text{new}} x_k$, after an initial exploration phase whose worst-case cost is minimized simultaneously with the controller design. The influence of the exploring controller $u_k = K_e x_k + e_k$, $e \sim \mathcal{N}(0, \Sigma)$, on the performance after exploration is quantified by (approximately) predicting the future uncertainty depending on K_e and Σ via (17)–(18).

Compared to previous works [19]–[22], the key difference of the present approach is that the mean of the parameter estimates after exploration is taken into account by considering it as a scheduling variable via $w_k^s = \Delta_s z_k^s$. Initially, it is only known that $\Delta_s \in \mathbf{\Delta}_s$ (compare Lemma 2), but after exploration Δ_s is available and can hence be exploited for controller design. This is in contrast to existing works, which simply assumed $\Delta_s = 0$, i.e., the mean value of the parameter estimates does not change over time. An important observation is that, according to (29), the state-feedback K_{new} depends on \hat{A}_T, \hat{B}_T and hence, on the data \mathcal{D}_T obtained during time steps 0 through T . This means that the proposed controller explicitly exploits measurements during the exploration phase, which was not the case in [19]–[22]. Furthermore, [19]–[22] require a repeated LMI based design after the exploration phase, which is not the case in our formulation wherein we pre-compute a closed-form solution that guarantees quadratic performance based on a predicted bound of the exploration data.

Theorem 1 requires that Assumption 1 holds, which is a non-restrictive condition on the initial data and parameters. On the contrary, Assumption 2 is essentially an approximation on the empirical covariance, which is required to predict

the influence of the exploration phase on the parameter estimates. While Assumption 2 is generally not guaranteed to hold, it is approximately satisfied in practice and its validity can be verified a posteriori, i.e., after the exploration phase.

We briefly wish to elaborate on the impact of different values c_δ, σ_w corresponding to different noise and confidence levels. Assuming a fixed initial uncertainty D_0 is given, c_δ, σ_w have the same effect and only appear in (28f) to determine \bar{D}_T . In case we increase c_δ and/or σ_w^2 (assuming D_0 is fixed), the optimal controller parameters K, K_s resulting from (28) remain unchanged and only the cost of the exploration (Y_e, W_e, Z_e, Σ) increases proportionally. This is natural, as a higher noise level and/or a higher desired confidence level requires a stronger excitation to yield the required model quality. Thus, since the magnitude of noise and/or confidence level does not directly impact the resulting controller K, K_s (although K_{new} may change), the main structural property that would *qualitatively* change the shape of the resulting robust dual control strategy would be varying noise levels for the different states.

V. CONCLUSION

In this paper, we formulate a novel dual control approach for linear time-invariant systems with performance guarantees based on gain-scheduling. We propose an LMI-based controller design procedure which simultaneously computes a controller to apply during an exploration phase as well as a robust controller for closed-loop performance after exploration. Similar to [20], the influence of the exploration on the closed loop is quantified by predicting the future uncertainty of the system parameters. The key difference is that we account for the change in the mean estimate of system parameters after exploration by formulating an LPV closed-loop system and selecting the uncertain system parameters as a scheduling variable. In contrast to existing methods, the robust controller takes the estimates after exploration into account and therefore, it depends explicitly on the data obtained during exploration. Finally, we prove desirable theoretical properties of the proposed approach. An interesting issue for future research is a detailed comparison of the presented dual controller to existing alternatives.

REFERENCES

- [1] A. A. Feldbaum, “Dual control theory,” *Automation and Remote Control*, vol. 21, no. 9, pp. 874–1039, 1960.
- [2] B. Recht, “A tour of reinforcement learning: The view from continuous control,” *Annual Review of Control, Robotics, and Autonomous Systems*, vol. 2, pp. 253–279, 2019.
- [3] N. M. Filatov and H. Unbehauen, “Survey of adaptive dual control methods,” *IEEE Proceedings-Control Theory and Applications*, vol. 147, no. 1, pp. 118–128, 2000.
- [4] A. Mesbah, “Stochastic model predictive control with active uncertainty learning: a survey on dual control,” *Annual Reviews in Control*, vol. 45, pp. 107–117, 2018.
- [5] J. L. Doob, *Stochastic processes*. New York Wiley, 1953, vol. 101.
- [6] E. Tse, Y. Bar-Shalom, and L. Meier, “Wide-sense adaptive dual control for nonlinear stochastic systems,” *IEEE Transactions on Automatic Control*, vol. 18, no. 2, pp. 98–108, 1973.
- [7] E. Tse and Y. Bar-Shalom, “Actively adaptive control for nonlinear stochastic systems,” *Proceedings of the IEEE*, vol. 64, no. 8, pp. 1172–1181, 1976.

- [8] D. S. Bayard and M. Eslami, "Implicit dual control for general stochastic systems," *Optimal Control Applications and Methods*, vol. 6, no. 3, pp. 265–279, 1985.
- [9] Y. Bar-Shalom and E. Tse, "Caution, probing, and the value of information in the control of uncertain systems," in *Annals of Economic and Social Measurement, Volume 5, number 3*. NBER, 1976, pp. 323–337.
- [10] B. Wittenmark, "Adaptive dual control methods: An overview," in *Adaptive Systems in Control and Signal Processing 1995*. Elsevier, 1995, pp. 67–72.
- [11] M. Gevers, "Identification for control: From the early achievements to the revival of experiment design," *European journal of control*, vol. 11, pp. 1–18, 2005.
- [12] H. Hjalmarsson, "From experiment design to closed-loop control," *Automatica*, vol. 41, no. 3, pp. 393–438, 2005.
- [13] M. Annergren, C. A. Larsson, H. Hjalmarsson, X. Bombois, and B. Wahlberg, "Application-oriented input design in system identification: Optimal input design for controls," *IEEE Control Systems Magazine*, vol. 37, no. 2, pp. 31–56, 2017.
- [14] C. A. Larsson, A. Ebadat, C. R. Rojas, X. Bombois, and H. Hjalmarsson, "An application-oriented approach to dual control with excitation for closed-loop identification," *European Journal of Control*, vol. 29, pp. 1–16, 2016.
- [15] T. A. N. Heirung, B. E. Ydstie, and B. Foss, "Dual adaptive model predictive control," *Automatica*, vol. 80, pp. 340–348, 2017.
- [16] S. Dean, H. Mania, N. Matni, B. Recht, and S. Tu, "Regret bounds for robust adaptive control of the linear quadratic regulator," in *Advances in Neural Information Processing Systems*, 2018, pp. 4188–4197.
- [17] S. Dean, S. Tu, N. Matni, and B. Recht, "Safely learning to control the constrained linear quadratic regulator," in *Proceedings of the 2019 American Control Conference (ACC)*. IEEE, 2019, pp. 5582–5588.
- [18] S. Dean, H. Mania, N. Matni, B. Recht, and S. Tu, "On the sample complexity of the linear quadratic regulator," *Foundations of Computational Mathematics*, pp. 1–47, 2019.
- [19] J. Umenberger, M. Ferizbegovic, T. B. Schön, and H. Hjalmarsson, "Robust exploration in linear quadratic reinforcement learning," in *Advances in Neural Information Processing Systems*, 2019, pp. 15 310–15 320.
- [20] M. Ferizbegovic, J. Umenberger, H. Hjalmarsson, and T. B. Schön, "Learning robust LQ-controllers using application oriented exploration," *IEEE Control Systems Letters*, vol. 4, no. 1, pp. 19–24, 2019.
- [21] M. Barenthin and H. Hjalmarsson, "Identification and control: Joint input design and H_∞ state feedback with ellipsoidal parametric uncertainty via lmis," *Automatica*, vol. 44, no. 2, pp. 543–551, 2008.
- [22] A. Iannelli, M. Khosravi, and R. S. Smith, "Structured exploration in the finite horizon linear quadratic dual control problem," in *Proc. 21st IFAC World Congress*, 2020.
- [23] C. W. Scherer, "LPV control and full block multipliers," *Automatica*, vol. 37, no. 3, pp. 361–375, 2001.
- [24] J. Veenman and C. W. Scherer, "A synthesis framework for robust gain-scheduling controllers," *Automatica*, vol. 50, no. 11, pp. 2799–2812, 2014.
- [25] C. Scherer and S. Weiland, "Linear matrix inequalities in control," *Lecture Notes, Dutch Institute for Systems and Control, Delft, The Netherlands*, vol. 3, 2000.