

---

# Safe online nonstochastic control

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 Online nonstochastic control has emerged as a promising strategy for online convex  
2 optimization of control policies for linear systems subject to adversarial distur-  
3 bances and time-varying cost functions. However, ensuring safety in these systems  
4 remains a significant open problem, especially when the system parameters are  
5 unknown. Practical nonstochastic control algorithms for real-world systems must  
6 adhere to safety constraints without becoming overly conservative or relying on  
7 exact models. We address this challenge by presenting a safe nonstochastic control  
8 algorithm for systems with unknown parameters subject to state and input  
9 constraints. Given data of a single disturbed input-state trajectory, we design non-  
10 conservative constraint sets for the policy parameters and develop a robust strongly  
11 stabilizing controller. By drawing a connection to model predictive control, we pro-  
12 pose a new analysis perspective and show how a slight change in the nonstochastic  
13 control algorithm can drastically improve performance if disturbances are constant  
14 or slowly time-varying.

## 15 1 Introduction

16 In reinforcement learning, gradient-based policy optimization has shown great success in practice  
17 Schulman et al. [2017]. For learning-based control, the paradigm of online convex optimization  
18 offers a powerful framework for iteratively updating control policies based on gradients and observed  
19 data. Nonstochastic control is such a gradient-based control method that has been proven effective  
20 for the control of linear dynamical systems in the face of deterministic, possibly adversarial, bounded  
21 disturbances and adversarially chosen cost functions [Agarwal et al., 2019, Hazan et al., 2020,  
22 Simchowitz, 2020]. At each time step, a convex cost function is revealed to the learner and the  
23 policy gradient is approximated by applying the cost function to the terminal state and action of  
24 a model-based rollout (simulation). Since optimizing over the function space of state or output  
25 feedback policies is computationally intractable, see for example [Goulart et al., 2006], disturbance  
26 feedback policies are employed. Nonstochastic control algorithms have been adapted or extended  
27 for different settings such as partial observability [Simchowitz et al., 2020], changing dynamics  
28 [Minasyan et al., 2021], bandit loss [Sun et al., 2023] or fully unknown linear systems [Chen and  
29 Hazan, 2021]. However, one critical challenge is the inclusion of a safety guarantee in the sense of  
30 adherence to state and input constraints, particularly in the presence of model uncertainty.

31 Little research on nonstochastic control so far has considered the addition of input and state constraints.  
32 In the related literature, this problem setting has only been considered with access to an exact model  
33 Li et al. [2021], Nonhoff et al. [2024], Liu et al. [2023], Zhou and Tzoumas [2023], Martin et al.  
34 [2023] or achieved results in high probability with i.i.d. disturbances and conservative fixed parameter  
35 constraints with careful transitions in between updates Li et al. [2024]. For systems with unknown  
36 parameters, most works propose a sequential approach of system identification via least squares  
37 estimation (LSE) and control. Recent advances in statistical learning theory and the related discovery  
38 of high-probability finite time guarantees for models obtained via LSE [Wagenmaker and Jamieson,

2020, Sarkar and Rakhlin, 2019, Foster and Simchowitz, 2020, Simchowitz et al., 2018] was leveraged by a multitude of works to obtain high-probability regret guarantees in the nonstochastic control setting [Hazan et al., 2020, Chen and Hazan, 2021, Simchowitz, 2020].

In this work, we change perspective from statistical learning LSE to data-driven robust control based on a set of unfalsified models Berberich et al. [2020], Van Waarde et al. [2023], Teutsch et al. [2024]. In the nonstochastic control setting of linear systems subject to bounded disturbances, such a set may be constructed by set membership identification (SMI). Informally, SMI begins by considering the whole space of model parameters and continually discards those that could not have reproduced the seen data. By leveraging the disturbance bounds, the resulting sets can be much smaller than LSE confidence regions Li et al. [2023], and always contain the system’s true parameters, which allows for the design of robust controllers with certainty instead of high probability.

**Contribution:** This work presents a safe online optimal control algorithm for unknown linear systems subject to nonstochastic disturbances. Given an input-state data trajectory, we bridge the gap between low-regret nonstochastic control and safe data-driven robust control by designing safety constraints for online policy updates that hold for all models that may have produced the data. By drawing from concepts in model predictive control Rawlings et al. [2017], Lorenzen et al. [2019], we establish recursive feasibility of the safety constraints for all models and propose a subtle but effective change to the initial state of the rollouts used for the policy gradient, which leads to the elimination of steady-state errors in the case of constant or slowly time-varying disturbances. We show the practical potential of the approach in a small simulation example.

## 2 Preliminaries and problem setting

In the nonstochastic control setting, the learner is presented with a linear time-invariant dynamical system

$$x_{t+1} = Ax_t + Bu_t + w_t \quad (1)$$

where  $x \in \mathbb{R}^{n_x}$  is the state of the system and  $u \in \mathbb{R}^{n_u}$  is the input or action taken by the learner. The disturbance  $w \in \mathbb{R}^{n_x}$  represents uncertainty and is not subject to any assumed stochastic properties, but may be chosen from a known compact set  $\mathcal{W}$  by an adversary at each time step and remains unknown to the learner. In this work, we assume  $\mathcal{W}$  is a convex polytope  $\mathcal{W} = \{w \in \mathbb{R}^{n_x} \mid G_w w \leq g_w\}$ . At each time step  $t$ , the learner measures the current state  $x_t$  and a cost function  $c_t : \mathbb{R}^{n_x} \times \mathbb{R}^{n_u} \rightarrow \mathbb{R}$  is revealed. The goal is to learn a policy that chooses inputs which minimize the cumulative costs  $\sum_t c_t(x_t, u_t)$ .

### 2.1 Disturbance-action policies and the Gradient perturbation controller

The considered policies are from the class of *disturbance-action policies* Agarwal et al. [2019], also called *affine disturbance feedback*, see for example [Goulart et al., 2006]. Instead of basing decisions on the current state directly, these policies compute the input based on estimates of past disturbances  $\hat{w}_t$ . These estimates are based on a system model  $(\hat{A}, \hat{B}) \approx (A, B)$ . At time step  $t$ , the disturbance estimate  $\hat{w}_{t-1}$  is computed as the prediction error

$$\hat{w}_{t-1} = x_t - (\hat{A}x_{t-1} + \hat{B}u_{t-1}). \quad (2)$$

**Definition 1.** A *disturbance-action policy (DAP)*  $\pi_{\text{DAP}}(\underline{M})$  chooses inputs based on parameter matrices  $M_i$  via

$$v_t = m_0 + \sum_{i=1}^L M_i \hat{w}_{t-i} = \underline{M} \hat{w}_t \quad (3)$$

where  $L$  is the memory length and  $\underline{M} = [m_0, M_1 \dots, M_L]$ ,  $\hat{w}_t = [1, \hat{w}_{t-1}^T, \dots, \hat{w}_{t-L}^T]^T$  allow for shorter notation.

In order to guarantee stability, DAPs are often used together with a fixed stabilizing state feedback controller  $u_t = Kx_t + v_t$ . We will do the same and abbreviate  $(A + BK) = A_K$  in the following.

A *gradient perturbation controller (GPC)* iteratively updates the policy (3) based on gradients computed via a model rollout. Online, at each time step  $t$ , the state  $x_t$  is measured, the last disturbance  $\hat{w}_{t-1}$  is estimated and the control parameters  $\underline{M}_t$  are updated by taking a gradient-step as

$$M_{t+1,i} = M_{t,i} - \eta_t \nabla l_t(M_{t,i}) \quad (4)$$

84 where  $\eta_t > 0$  is the learning rate. The loss  $l_t$  approximates the that would have been obtained  
 85 under the fixed policy  $\pi_{\text{DAP}}(\underline{M}_t)$  and is defined based on the terminal state and input of an  $H$ -step  
 86 simulation of the current with a given model  $(\hat{A}, \hat{B})$  and the most recent disturbance estimates  
 87  $\hat{w}_{t-L-H:t}$ . That is,  $l_t(\underline{M}_t) = c_t(x_{H|t}(\underline{M}_t), u_{H|t}(\underline{M}_t))$  where  $(x_{k|t}, u_{k|t})$  denotes the simulated  
 88 states and inputs running from  $k = 0, \dots, L$  computed at time step  $t$  via

$$x_{0|t} = 0, \quad w_{k|t} = \hat{w}_{t-H+k}, \quad \underline{w}_{k|t} = [1, w_{k-1|t}^T, \dots, w_{k-L|t}^T]^T, \quad u_{k|t} = Kx_{k|t} + \underline{M}_t \underline{w}_{k|t}, \quad (5)$$

$$x_{k+1|t} = \hat{A}x_{k|t} + \hat{B}u_{k|t} + w_{k|t}, \quad k = 0, \dots, H-1.$$

89 The justification is here that the actual state may be well approximated by such a simulation, since it  
 90 evolves as  $x_t = A_K^t x_0 + \sum_{k=0}^{t-1} A_K^k B v_{t-1-k} + A_K^k w_{t-1-k}$  whereas a simulation with horizon  $H$ ,  
 91 initial state zero, and the latest  $H$  disturbances reads  $\tilde{x}_t = \sum_{k=0}^{H-1} A_K^k B v_{t-1-k} + A_K^k w_{t-1-k} \approx x_t$ .  
 92 The resulting approximation error reads

$$x_t - \tilde{x}_t = \sum_{k=H}^{t-1} A_K^k (B v_{t-1-k} + w_{t-1-k}) = A_K^H \sum_{k=0}^{t-1-H} A_K^k (B v_{t-1-k} + w_{t-1-k}) = A_K^H x_{t-H} \quad (6)$$

93 and is small for stable  $A_K$  and large memory  $H$ . How small is captured by the following quantitative  
 94 notion of stability introduced in Cohen et al. [2018].

95 **Definition 2.**  $K$  is a  $(\kappa, \gamma)$ -strongly stable controller for  $(A, B)$  if  $\|A_K^t\| \leq \kappa \gamma^t$  for all  $t \geq 0$ .

96 Equipped with convergence bounds for the dynamics  $A_K$ , the presented gradient perturbation controller  
 97 (3)-(5) enjoys sublinear regret against the best fixed policy  $\underline{M}^*$  in hindsight, and thereby sublinear  
 98 regret against an expressive class of controllers, see Hazan and Singh [2023] for an overview of  
 99 results.

## 100 2.2 Problem setting: Uncertain system and safety constraints

101 While the presented control scheme has been extended in many directions, for example to bandit loss  
 102 functions Sun et al. [2023] and partial observations Simchowit et al. [2020], one challenge for the  
 103 application of GPC to safety-critical systems is the adherence to input and state constraints in the face  
 104 of model uncertainty. In this work, we consider a setting where the true system parameters  $(A, B)$   
 105 are unknown, and only an input-state trajectory  $\{u_t, x_t\}_{t=0}^{T_D}$  is available. Furthermore, we restrict  
 106 actions to a set

$$u_t \in \mathcal{U} = \{u \in \mathbb{R}^{n_u} \mid G_{\mathcal{U}} u \leq g_{\mathcal{U}}\} \quad \forall t \geq 0. \quad (7)$$

107 and subject the state to polytopic safety constraints

$$x_t \in \mathcal{S} = \{x \in \mathbb{R}^{n_x} \mid G_{\mathcal{S}} x \leq g_{\mathcal{S}}\} \quad \forall t \geq 1 \quad (8)$$

108 where both  $\mathcal{U}$  and  $\mathcal{S}$  are known user-specified convex compact sets that contain the origin. In order to  
 109 render the problem of safety tractable, we assume that there exists a state feedback controller that  
 110 can keep the system safe from initial state  $x_0 = 0$  no matter which disturbances are chosen by the  
 111 adversary.

112 **Definition 3** (Safe control policies). A control policy is called safe if it generates inputs  $u_t \in \mathcal{U}$  for  
 113 which the state of (1) satisfies  $x_t \in \mathcal{S}$  for all time  $t \geq 0$ .

114 **Assumption 1.** There exists  $K_{\text{safe}}$  such that given  $x_0 = 0$ , the state feedback  $u_t = K_{\text{safe}} x_t$  is safe  
 115 for all disturbance realizations  $w_t \in \mathcal{W}$  and all time  $t \geq 0$ .

116 **Remark 1.** On first glance, Assumption 1 may seem restrictive, but note that 1) we do not have  
 117 access to  $K_{\text{safe}}$ , and 2) the application of  $u_t = K_{\text{safe}} x_t$  may incur high costs without the disturbance  
 118 feedback, whose addition can in turn cause a loss of safety. Informally, Assumption 1 guarantees that  
 119 the disturbances in  $\mathcal{W}$  are not too large compared to the set of safe states  $\mathcal{S}$ .

## 120 3 Safe nonstochastic control

121 A guarantee of safety constraints (8) during operation requires that the control input  $u_t$ , or more  
 122 specifically the control parameters  $\underline{M}_t$  are always chosen such that  $x_{t+1} \in \mathcal{S}$ . Since most non-  
 123 stochastic control algorithms in the literature already include projections into a set of parameters  $\mathcal{M}$   
 124 for regret guarantees, it is a natural adaptation to enforce input and constraint satisfaction by similarly  
 125 projecting into the set of safe control parameters  $\mathcal{M}_{\text{safe}, t}$  and applying  $u_t(\Pi_{\mathcal{M}_{\text{safe}, t}} M_{\text{opt}, t})$  instead of  
 126  $u_t(\underline{M}_{\text{opt}})$  for control. The resulting algorithm presented in this paper is shown in Algorithm 1.

---

**Algorithm 1** Safe online optimal control

---

**Identification phase:**Collect data  $(u_t, x_t)_{t=0}^{T_{\text{ini}}}$ .Construct set of unfalsified models  $\Omega$  as in (9).Compute  $(\kappa, \gamma)$ -strongly stable state feedback gain  $K$  for all  $\Omega$  as in Lemma 1.Choose a nominal model  $(\hat{A}, \hat{B}) \in \Omega$ .**Control phase:****for** each time step  $t = T_{\text{ini}}, \dots, T_{\text{ini}} + T$  **do**    Record state  $x_t$  and construct latest disturbance estimate  $\hat{w}_{t-1} := x_t - (\hat{A}x_t + \hat{B}u_t)$ .    Receive cost function and update policy  $\underline{M}_{t,\text{opt}} = \underline{M}_{t-1} - \eta_t \nabla l_t(\underline{M}_{t-1})$ .    Project to closest safe policy  $\underline{M}_t = \Pi_{\mathcal{M}_{\text{safe},t}} \underline{M}_{t,\text{opt}}$ .    Apply control  $u_t(\underline{M}_t)$  (3)127 **3.1 From data to a set of models**

128 Instead of identifying one best-fit system, we consider the set of models  $(A, B)$  that agree with (may  
129 have produced) the given or recorded input-state data. Let  $\mathbb{Z}_{i,j} = \{i, i+1, \dots, j\}$ . Given an input-state  
130 data trajectory  $\{x_t, u_t\}_{t=0}^T$  resulting from the application of  $T$  arbitrary inputs to system (1), and  
131 assuming that the unknown disturbances  $\{w_t\}_{t=0}^{T-1}$  were always in the known set  $\mathcal{W}$ , the resulting set  
132 of *consistent* or *unfalsified* models is given by

$$\Omega_{[0,T]} = \{[A \ B] \in \mathbb{R}^{n_x \times (n_x + n_u)} \mid x_{t+1} - [A \ B] \begin{bmatrix} x_t \\ u_t \end{bmatrix} \in \mathcal{W}, t = 0, \dots, T-1\}. \quad (9)$$

133 The set  $\Omega_{[0,T]}$  inherits convexity and closedness from  $\mathcal{W}$ , and can be directly constructed in half-space  
134 representation by reorganizing the inequality constraints that represent  $\mathcal{W}$ . If the data trajectory  
135 is sufficiently informative (see Lemma 2 in the Appendix), then  $\Omega_{[0,T]}$  is also bounded and may  
136 be described as convex hull of its vertices  $\Omega = \text{conv}(\{[A_i \ B_i]\}_{i=1}^{N_v})$ . As a representation of model  
137 uncertainty,  $\Omega$  behaves nicely: First, as new data streams in, new constraints are added to the set  
138 and therefore updates never increase the uncertainty set in size. Second, crucially, as long as the  
139 assumed disturbance bound  $\mathcal{W}$  holds,  $\Omega$  always contains the true data-generating system matrices by  
140 construction, and every statement that holds for all models  $[A \ B]$  inside  $\Omega$  necessarily holds for the  
141 actual unknown system. Since  $\Omega$  is defined by input-state data (and the disturbance bound  $\mathcal{W}$ ) alone,  
142 these statements can be inferred directly from data. In this work, we will use the set of models  $\Omega$  to  
143 construct constraints on the control parameters with which safety can be guaranteed.

144 **3.2 From a set of models to safety constraints**

145 In this work, safety is defined as constraints in input and state space. In order to derive a set of  
146 safe control parameters, we need to map the state space constraints  $x_t \in \mathcal{S}$  into constraints on the  
147 policy parameters  $\underline{M}_t$ . As intermediate mapping, we may consider the space of inputs since  $v_t$  spans  
148 all of  $\mathcal{R}^m$  in the sense that any desired safe input  $u$  can be reproduced by some choice of control  
149 parameters  $\underline{M}$  such that  $u = Kx + v(\underline{M}_t)$ . The challenge is that the constraints on  $\underline{M}$  1) need to be  
150 recursively feasible, i.e., the state is only steered to where constraint satisfaction remains possible, 2)  
151 need to consider all possible models in  $\Omega$  need to be considered, and 3) should not be conservative  
152 but restrict the space of parameters as little as possible.

153 Consider the set of models  $\Omega$  containing the true system and the disturbance bound  $\mathcal{W}$ . The state  
154 evolution of the unknown system (1) satisfies the inclusion

$$x_{t+1} \in \Omega \begin{bmatrix} x_t \\ u_t \end{bmatrix} \oplus \mathcal{W}, \quad \Omega \begin{bmatrix} x_t \\ u_t \end{bmatrix} = \{[A \ B] \begin{bmatrix} x_t \\ u_t \end{bmatrix} \mid [A \ B] \in \Omega\} \quad (10)$$

155 where  $\oplus$  denotes the Minkowski set addition. With the Minkowski (Pontryagin) set difference  $\ominus$ , we  
156 can reformulate the above into a sufficient condition on the state and input at the current time step for  
157 satisfaction of safety constraints at the next time step,

$$\Omega \begin{bmatrix} x_t \\ u_t \end{bmatrix} \in \mathcal{S} \ominus \mathcal{W} \implies x_{t+1} \in \mathcal{S}. \quad (11)$$

158 **Remark 2.** *Non-emptiness of  $\mathcal{S} \ominus \mathcal{W}$  is covered by Assumption 1 since the existence of  $K_{\text{safe}}$  implies*  
 159 *that the safe set  $\mathcal{S}$  is specified large enough to contain the disturbance set, i.e.,  $\mathcal{S} \supseteq \mathcal{W}$  and*  
 160 *consequently  $\mathcal{S} \ominus \mathcal{W} \neq \emptyset$ .*

161 In order to guarantee that the left-hand-side of (11) remains feasible during operation, we construct a  
 162 (maximal) robust control invariant subset of  $\mathcal{S}$  [Blanchini, 1999, Rawlings et al., 2017].

163 **Definition 4.** *A set  $\mathcal{X}$  is robust control invariant (RCI) for dynamics  $x_{t+1} = Ax_t + Bu_t + w_t$ ,*  
 164  *$w_t \in \mathcal{W}$ , if for all  $x \in \mathcal{X}$  there exists  $u \in \mathcal{U}$  such that  $Ax + Bu + w \in \mathcal{X}$  for all  $w \in \mathcal{W}$ .*

165 A maximal RCI subset of the safe state set  $\mathcal{S}$  is a set that contains all other RCI subsets of  $\mathcal{S}$ . The  
 166 maximal RCI subset is well defined since the set property of robust control invariance is closed  
 167 under the union. For the present discrete-time linear dynamics, maximal RCI sets are computed via  
 168 recursive erosion and expansion [Blanchini and Miani, 2015], see Appendix for details.

169 In the following, let  $\mathcal{X}$  be the maximal subset of  $\mathcal{S}$  that is RCI for all models in  $\Omega$ . That is, let  $\mathcal{X}$  be  
 170 such that

$$(\forall x \in \mathcal{X})(\exists u \in \mathcal{U}) \Omega \begin{bmatrix} x \\ u \end{bmatrix} \in \mathcal{X} \ominus \mathcal{W}. \quad (12)$$

171 Since  $\Omega$ ,  $\mathcal{S}$  and  $\mathcal{W}$  are compact and convex polytopes, so is  $\mathcal{X}$  and we can write  $\mathcal{X} = \{x \in \mathbb{R}^n \mid$   
 172  $G_{\mathcal{X}}x \leq g_{\mathcal{X}}\}$ . Similarly, define  $\mathcal{X} \ominus \mathcal{W} = \{x \in \mathbb{R}^n \mid G_{\mathcal{X}}x \leq g_{\mathcal{X} \ominus \mathcal{W}}\}$ .

173 **Remark 3.** *For the true system,  $\mathcal{X}$  is nonempty by Assumption 1 and contains the origin.*

174 Substituting the DAC policy (3) into (12) and reformulating based on vertices of  $\Omega$  leads to linear  
 175 constraints on the control parameters

$$G_{\mathcal{X}}B_iM\hat{w}_t \leq g_{\mathcal{X} \ominus \mathcal{W}} - G_{\mathcal{X}}(A_i + B_iK)x_t \quad i = 1, \dots, N_v, \quad (13)$$

$$G_{\mathcal{U}}M\hat{w}_t \leq g_{\mathcal{U}} - G_{\mathcal{U}}Kx_t, \quad (14)$$

176 which define a convex constraint set  $\mathcal{M}(\hat{w}_t, x_t) = \{M \in \mathbb{R}^{n_u \times L n_x} \mid (13), (14) \text{ are satisfied}\}$  that is  
 177 parameterized by the past estimated disturbances in  $\hat{w}_t$  and the current state  $x_t$ .

178 **Remark 4.** *Note that  $\mathcal{M}(\hat{w}_t, x_t)$  also depends on the chosen state feedback gain  $K$ , which is however*  
 179 *constant throughout. If  $K$  is safe as per Definition 3 (such that  $K = K_{\text{safe}}$  from Assumption 1), then*  
 180  *$M_t = 0$  is a safe parameter choice for all time and  $\{0\}$  is a common subset of all sets  $\mathcal{M}(\hat{w}_t, x_t)$*   
 181 *with  $x_t \in \mathcal{X}$ ,  $\hat{w}_{t-i} \in \mathcal{W}$ . In general, the set of control parameters which are safe for all possible*  
 182 *states and disturbances is the intersection of all such  $\mathcal{M}(\hat{w}, x)$ . In other words, more restrictive but*  
 183 *fixed safety constraints as in [Li et al., 2021, 2024] are recovered by minimizing the RHS (13), (14)*  
 184 *of over all  $x_t \in \mathcal{X}$  and requiring the inequalities to hold for all  $\hat{w}_t = [1, \hat{w}_{t-1}^T, \dots, \hat{w}_{t-L}^T]^T$  with*  
 185  *$\hat{w}_{t-i} \in \mathcal{W}$ .*

### 186 3.3 Theoretical guarantees

187 If at each time step  $t$ , the control parameters  $M_t$  are projected into  $\mathcal{M}_{\text{safe},t} = \mathcal{M}(\hat{w}_t, x_t)$ , we may  
 188 guarantee safety as shown in the following result.

189 **Lemma 1** (Recursive feasibility). *Let  $x_t \in \mathcal{X}$ . Then  $\mathcal{M}_{\text{safe},t} \neq \emptyset$  and any choice of control*  
 190 *parameters  $M_t \in \mathcal{M}_{\text{safe},t}$  leads to a nonempty constraint set in the next time step,  $\mathcal{M}_{t+1,\text{safe}} \neq \emptyset$ .*

191 *Proof.* Since  $\mathcal{X}$  is RCI,  $\mathcal{M}_{\text{safe},t} \neq \emptyset$  for all  $x \in \mathcal{X}$  by construction. Moreover, any choice of control  
 192 parameters  $M_t \in \mathcal{M}_{\text{safe},t}$  leads to  $x_{t+1} \in \mathcal{X}$ . Since  $\mathcal{X} \subseteq \mathcal{F}$ , the safety condition (12) is feasible for  
 193 all states in  $\mathcal{X}$ .  $\square$

194 **Theorem 1** (Constraint satisfaction). *Assume that the model uncertainty in  $\Omega$  is small enough such*  
 195 *that  $\mathcal{X} \neq \emptyset$ . Then, for any  $T \geq 0$  and all possible disturbance sequences  $w_{0:T} \in \mathcal{W}$ , the proposed*  
 196 *control strategy in Algorithm 1 is safe in the sense of Definition 3.*

197 *Proof.* Since  $x_0 = 0$ , we have  $x_0 \in \mathcal{X}$  as long as  $\mathcal{X} \neq \emptyset$ . By Lemma 1 the set of control parameters  
 198  $\mathcal{M}_{\text{safe},0}$  is not empty and any choice  $M_0 \in \mathcal{M}_{\text{safe},0}$  satisfies input constraints by construction of  
 199  $\mathcal{M}_{\text{safe},0}$  and leads to a next state  $x_1 \in \mathcal{X}$ . Since  $\mathcal{X} \subseteq \mathcal{S}$ , the next state  $x_1$  is safe. Safety for all time  
 200 follows by induction.  $\square$

201 **Remark 5.** Besides safety with certainty, the proposed approach based on SMI offers another distinct  
 202 advantage: Online adaptation of the above safety constraints is trivial. With every new data triple  
 203  $(x_{t-1}, u_{t-1}, x_t)$ ,  $\Omega$  may be updated by adding the constraints representing  $x_t - Ax_{t-1} - Bu_{t-1} \in \mathcal{W}$ .  
 204 Since  $\Omega_t \subseteq \Omega_{t-1}$ , a newly computed maximal RCI set  $\mathcal{X}_t$  will always contain the prior version,  
 205  $\mathcal{X}_t \supseteq \mathcal{X}_{t-1}$  and the constraints on the policy parameters are relaxed without loss of recursive  
 206 feasibility. If computational time is an issue, update computations may happen asynchronously by  
 207 computing an update of  $\mathcal{X}$  on a batch of new data and injecting it into the control algorithm once  
 208 the computation is finished. By contrast, a similarly easy adaptation is not possible if error bounds  
 209 around a least square estimate replace the set of models: the error bound of LSE decreases with more  
 210 data, but the change of the estimate itself may cause the new set of models to not be contained in the  
 211 prior one. Consequently, a careful transition between updates is necessary Li et al. [2024], which is  
 212 not the case in the proposed approach.

213 If the gradient perturbation controller presented in Section 2 runs with an approximate model  
 214  $(\hat{A}, \hat{B})$ , the only difference between the loss simulation (5) and the approximation in (6) is that an  
 215 additional error is introduced due to the model error. By also bounding this additional error, setting  
 216 an appropriate learning rate, and restricting control parameters  $\underline{M}_t$  to a special set, the gradient  
 217 perturbation controller for uncertain systems achieves sublinear regret with respect to the class of  
 218 state feedback controllers  $u_t = Kx_t$  Hazan et al. [2020], linear dynamical controllers Simchowitz  
 219 [2020], or disturbance action policies with fixed control parameters Chen and Hazan [2021]. In  
 220 order to recover similar regret guarantees with the additional projection to safety, we too require a  
 221  $(\kappa, \gamma)$ -strongly stabilizing controller as in Definition 2. In the foll we show how such a controller  
 222 may be constructed for all models in  $\Omega$ .

223 **Synthesis of a strongly stabilizing controller** A sufficient condition for stability  $\rho(A + BK) < 1$   
 224 of all models  $(A, B) \in \Omega$  is given by existence of a common quadratic Lyapunov function  $V(x) =$   
 225  $x^T Px$  for all hypothetical closed loop systems  $x^+ = (A + BK)x$ ,  $(A, B) \in \Omega$ . Computationally,  
 226 this check requires solving a finite system of linear matrix inequalities (LMI) in a semi-definite  
 227 program. Since regret bounds in the literature depend on the notion of  $(\kappa, \gamma)$ -strong stability, we  
 228 provide a semi-definite program for the direct synthesis of a  $(\sqrt{c}, \gamma)$ -strongly stable controller with  
 229 specified rate  $\gamma < 1$  and minimal constant  $\sqrt{c}$  in the following. The idea is to combine a bound on  
 230 the norm powers of  $A_K$  based on the positive invariance of Lyapunov sublevel sets [Ahiyevich et al.,  
 231 2018] with the fact since  $\rho(rA) = r\rho(A)$  for any matrix  $A$ , stability of  $\frac{1}{\gamma}A_K$  (i.e.,  $\rho(\frac{1}{\gamma}A_K) \leq 1$ )  
 232 implies  $\rho(A_K) \leq \gamma$ . Recall that  $\Omega = \text{conv} \{[A_i \ B_i]\}_{i=1}^{N_v}$ .

233 **Proposition 1.** Choose a desired spectral radius  $0 \leq r < 1$  and let  $(c, Z, Y)$  be the solution of

$$\underset{c, Z, Y}{\text{minimize}} \quad c \tag{15}$$

$$\text{subject to} \quad I_n \preceq Z \preceq cI_n, \tag{16}$$

$$\begin{bmatrix} rZ & A_i Z + B_i Y \\ * & rZ \end{bmatrix} \succ 0 \quad \forall i = 1, \dots, N_v. \tag{17}$$

234 Then the controller  $K = YZ^{-1}$  is  $(\sqrt{c}, r)$ -strongly stable for all  $(A, B) \in \Omega$ .

235 Please see Appendix for the proof.

236 **On regret bounds with safety constraints** The presented algorithm allows to run a safe variant  
 237 of GPC with any nominal model  $(\hat{A}, \hat{B}) \in \Omega$ , for example chosen via LSE and projection or as  
 238 Chebyshev center of  $\Omega$ . The computation of a strongly stable controller in Proposition 1 allows for  
 239 a recovery of GPC regret bounds in literature, as long as the safety constraints are not active. The  
 240 presented design of safety constraints restricts control parameters as little as possible. In fact, it was  
 241 motivated by the following Proposition.

242 **Proposition 2.** Every causally safe control policy (without foreknowledge of  $w_t$ ) needs to keep the  
 243 state in the maximal RCI subset  $\mathcal{X}_{\max} \subseteq \mathcal{S}$ .

244 *Proof.* If starting from  $x_t$  there exists an input sequence that keeps the state inside  $\mathcal{S}$  for all possible  
 245 disturbance sequences  $w_t$ : and all time, then the resulting state trajectory would be part of the maximal  
 246 RCI subset of  $\mathcal{S}$ . Since  $x_t$  is not, the proof follows by contradiction.  $\square$

247 In other words, enforcing the state to stay within the maximal RCI subset of  $\mathcal{S}$  does not lead to a  
 248 meaningful change of regret bounds if the comparator class is restricted to causally safe policies. In  
 249 the continuation of this work, we are interested in formalizing regret bounds in case of active safety  
 250 constraints.

#### 251 4 Better policy gradients by adaptive initial state - An MPC perspective

252 In essence, the gradient perturbation controller presented above takes decisions that minimize the  
 253 loss of model-based predictions. Recall the definition of the loss  $l_t$  based on a model rollout. If  
 254 instead of updating parameter towards the minimzer of the loss function, the policy parameters were  
 255 chosen directly as the minimzer in each time step, the scheme could be interpreted as parameterized  
 256 model predictive control (MPC): At each time step, choose the policy parameters parameterized by  
 257 solving the finite-time optimal control problem (OCP)  $\underline{M}_t^* = \operatorname{argmin}_{\underline{M}} l_t(\underline{M})$ , where compared to  
 258 classical MPC formulations the costs act only on the terminal state. In other words, GPC tries to  
 259 emulate a parameterized MPC by always updating the parameters towards the MPC solution. As  
 260 such, MPC lends itself as analysis tool for GPC and existing results in MPC may carry over. One  
 261 difference between classical MPC formulations and the present nonstochastic control version defined  
 262 by  $l_t$  comes from the fact that in MPC, the simulation (or rollout) is interpreted as *prediction*, instead  
 263 of *loss approximation in hindsight*. As such, the initial state in (5) would be updated to the current  
 264 state at each time step, i.e., set to  $x_{0|t} = x_t$ .

265 Note that with  $x_{0|t} = 0$ , the optimal solution  $\underline{M}_t^*$  depends only on the current cost function  $c_t$  and the  
 266 past disturbances  $\hat{w}_{-L:H-1}$ . Imagine the case where the cost function is fixed and the disturbances  
 267 are constant or very slowly time-varying (compared to the update rate of GPC). Then,  $\underline{M}_t^*$  is constant  
 268 and GPC converges quickly to fixed parameters, representing a very simple constant policy. If  
 269 instead, the initial state of (5) was set to  $x_{0|t} = x_t$ , the OCP would implicitly represent a linear affine  
 270 map from  $x_t$  to  $\underline{M}_t^*$  [Goulart et al., 2006], with the map being parameterized by the disturbances.  
 271 As a consequence, GPC with varying initial state (for the loss simulation) could still influence the  
 272 dynamics.

273 **A pathological example for the gradient perturbation controller** Consider a simple integrator  
 274 system with constant disturbance where the first and second component of the state may denote  
 275 the position and velocity of a point mass, control inputs change the velocity, and the disturbance  
 276 represents unknown changes in acceleration and velocity in between time steps,

$$x_{t+1} = \begin{bmatrix} 1 & 0.1 \\ 0 & 1 \end{bmatrix} x_t + \begin{bmatrix} 0 \\ 1 \end{bmatrix} u_t + \begin{bmatrix} 1 \\ 0 \end{bmatrix}. \quad (18)$$

277 Let  $K = [-1 \quad -1]$  stabilize the system, imagine the objective is to keep the point mass at the origin,  
 278 and let the learner's system model be exact so that the resulting predictions (loss simulation) used to  
 279 compute the gradient are exact. Since the estimated disturbances in  $\hat{w}_t$  are constant, so is any DAP  
 280  $\underline{M}\hat{w}_t$  and we choose a minimal disturbance memory of  $L = 1$  without loss of generality. For ease of  
 281 exposition, set the horizon to  $H = 2$ . In this simple setting, we would expect GPC to perform quite  
 282 well. However, it does not, as seen in Figure 1 (a), where the position  $x_1$  tends to  $-10$  instead of zero.  
 283 As shown by the behavior of the associated MPC algorithm, this is not an issue of convergence, but of  
 284 a loss function disconnected to the problem at hand. Figure 1 (c),(d) shows the disconnect between  
 285 loss, which tends to zero, and cumulative costs, which grow unbounded. GPC takes gradient steps  
 286 that minimize  $\bar{x}_{2|t,1}^2 = (0.1v + 2)^2$  and converges to a constant input  $\underline{M}\hat{w} = -20$ . The resulting  
 287 steady state  $x_\infty = (A + BK)x_\infty + [1 \quad -20]^T$  is  $x_\infty = [-10 \quad -10]^T$ . With larger horizons  $H$ ,  
 288 the steady state error of GPC shrinks, but only tends to zero for the maximal choice  $H = t$ , i.e.,  
 289 if the full horizon is taken into account. For example a horizon of  $H = 50$  leads to a steady state  
 290  $[-0.0127, -10]^T$ .

291 If the loss simulation instead starts at the current state  $x_{0|t} = x_t$ , the steady state error vanishes  
 292 and MPC even beats the best fixed DAP  $\underline{M}^*$  computed in hindsight (and denoted by Opt). If  $x_t$  is  
 293 accounted for in the loss, GPC minimizes  $x_{2|t,1} = ([1 \quad 0] A^2 z_0 + 0.1v + 2)^2 = ([0.9 \quad 0.1] x_{0|t} +$   
 294  $0.1v + 2)^2$  and no longer tends to a constant input, but towards an affine linear state feedback  
 295  $v_t = -[9 \quad 1] x_t - 20$  under which the steady state  $x_\infty = [0 \quad -10]^T$  incurs zero cost. Regret against  
 296 the best fixed DAP in hindsight is not only sublinear, but bounded.

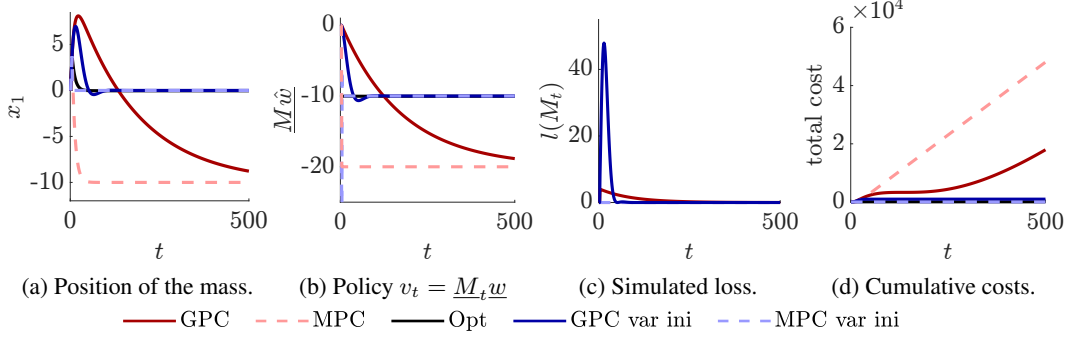


Figure 1: A simple pathological example of the basic nonstochastic control algorithm (OGD) as proposed in the literature. GPC’s loss tends to zero while the costs do not. With varying initial condition (var ini), the costs tend to zero.

297 **A generalization** Considering the MPC variants lets us generalize this example. In the follow-  
 298 ing, consider constant disturbances  $w_t = w$  and fixed costs  $c(x, u)$  with minimizing steady state  
 299  $(x^*, u^*) = \operatorname{argmin} c(x, u)$  such that  $x^* = A_K x^* + B u^* + w$ . Assume that  $x^*$  is reachable in  $H$   
 300 time steps and that  $x_{H|t} = x^*$  is the terminal state of the solution trajectory to the OCP such that  
 301  $x_{H|t} = A_K^H x_{0|t} + S_{H-1} B v + S_{H-1} w$ , where  $S_{H-1} = I + A_K + \dots + A_K^{H-1}$ . At every time step  $t$ ,  
 302 solving the OCP with  $x_{0|t} = 0$  leads to a constant input  $v_t$  where

$$B v_t = S_{H-1}^{-1} x^* - w, \quad x_{t+1} = A_K x_t + S_{H-1}^{-1} x^*. \quad (19)$$

303 The state thus converges, since  $A_K$  is stable, but setting  $x_t = x_{t+1} = x_\infty$  leads to

$$x_\infty = (I - A_K)^{-1} S_{H-1}^{-1} x^* = (I - A_K)^{-1} (I - A_K) (I - A_K^H) x^* = (I - A_K^H) x^* \quad (20)$$

304 so that  $x_t$  only converges (close) to  $x^*$  for very large horizons  $H$  where  $A_K^H \approx 0$ .

305 This is different in the case where the initial state is updated to the current state,  $x_{0|t} = x_t$ .

306 **Proposition 3.** Consider constant disturbances  $w_t = w$  and assume the predicted terminal state  
 307 satisfies  $x_{H|t} = x^*$  for all  $t \geq 0$ . Then the closed-loop dynamics induced by MPC with  $x_{0|t} = x_t$  are  
 308 stable and  $x_t$  converges to  $x^*$ .

309 The technical proof of Proposition 3 is in the Appendix. We note here that with the change of initial  
 310 state in the OCP, the first (optimal) predicted state  $x_{1|t}$  is the actual next state  $x_{t+1}$ . So that if the  
 311 state ever converges, i.e.,  $x_t = x_{t+1}$ , we had  $x_{1|t} = x_t$  which implies  $x_{k+1|t} = x_{k|t}$  (since the inputs  
 312  $v_{k|t}$  are constant) so that  $x_{H|t} = \dots = x_t$  which implies  $x_t = x^*$  by assumption. In short, the  
 313 state can *only* converge to the optimal state. As a consequence of Proposition 3, GPC with varying  
 314 initial state chases an optimal policy that achieves bounded  $O(1)$  regret, instead of one that induces a  
 315 steady-state error.

## 316 5 Simulation Example

317 Consider the numerical example of a linearized DC-DC converter from Section V.B in [Lorenzen  
 318 et al., 2016], where  $A = \begin{bmatrix} 1 & 0.0075 \\ -0.143 & 0.996 \end{bmatrix}$ ,  $B = \begin{bmatrix} 4.798 \\ 0.115 \end{bmatrix}$ , the state is subject to constraints  
 319  $|x_1| \leq 2$ ,  $|x_2| \leq 3$ , and the disturbance is bounded as  $\|w\|_\infty \leq 0.2$ . We let  $\mathcal{U} = \{u \in \mathbb{R} \mid |u| \leq 4\}$   
 320 and generate an input-state data trajectory of length  $T_{\text{Data}} = 15$  starting from zero initial state with  
 321 inputs and disturbances sampled uniformly from  $\mathcal{U}$  and  $\mathcal{W}$ , respectively. After building the set of  
 322 models  $\Omega$  from the data, we solve (15) with  $r = 0.6$  and receive a controller  $K = [-0.33 \ 0.78]$   
 323 that is  $(8.6, 0.6)$ -strongly stable for all models in  $\Omega$ . We choose the Chebyshev center of  $\Omega$  as  
 324 nominal model  $(\hat{A}, \hat{B})$ , set  $H = 10$ ,  $L = 1$ , pick a learning rate  $\eta = 0.1$  and transition to the control  
 325 phase of Algorithm 1. During the control phase, we let the disturbance vary at constant rate from  
 326 zero to  $[0.2 \ 0.2]^T$  and back to zero over  $T = 500$  time steps. The cost functions are defined as  
 327  $c_t(x_t, u_t) = (x_{t,2} - x_{t,2}^*)^2$  where  $x_{t,2}^* = 1.5$  for the first 250 time steps and  $x_{t,2}^* = -1.5$  for the



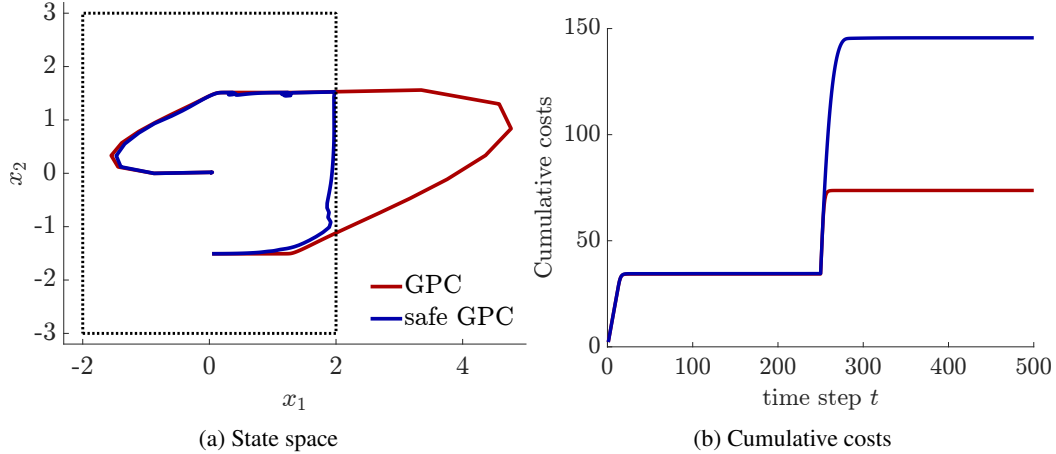


Figure 2: Behavior of safe GPC (blue) and GPC without state and input constraints (red). Both methods first steer to and stabilize the state at the optimal  $x_2 = 1.5$  in the first 250 time steps, and the optimal  $-1.5$  in the second 250 time steps. Since safe GPC needs to adhere to the state constraints on  $x_1$ , it takes more time steps to transition and suffers higher cost along the way.

328 last 250 time steps. Recall that both the disturbances and future cost functions are unknown to the  
 329 control algorithm. Figure 2 shows the resulting trajectories and cumulative costs for the proposed  
 330 safe nonstochastic control algorithm running with varying initial state as proposed in Section 4. For  
 331 comparison, the equivalent nonstochastic control algorithm without safety constraints is also shown.  
 332 In the transition from  $x_2 = 1.5$  to  $x_2 = -1.5$ , high values of  $x_1$  are necessary. As seen in Figure 2(a),  
 333 the proposed algorithm satisfies the safety constraints with virtually no conservatism.

## 334 6 Conclusion

335 This work addressed the challenge of ensuring safety in online nonstochastic control for linear  
 336 systems with unknown parameters. By leveraging a data-driven robust control approach based on  
 337 set membership identification, we derived non-conservative constraint sets for policy parameters  
 338 and constructed a strongly stabilizing controller. In contrast to existing works, both safety and  
 339 strong stability are guaranteed for all unfalsified models and hold with certainty. In simulation, we  
 340 demonstrated that our approach can effectively maintain system safety and performance from data  
 341 alone. By integrating principles from model predictive control, we ensured recursive feasibility of the  
 342 safety constraints and showed how updating the initial state of policy gradient rollouts effectively  
 343 eliminates steady-state errors under constant or slowly varying disturbances. Beyond the above, this  
 344 work left certain questions unanswered. First and foremost, we left a formal regret bound against an  
 345 expressive class of causally safe policies open for future work. We hypothesize that sublinear regret  
 346 against an expressive class of noncausally safe policies is unattainable in general, since a policy with  
 347 foreknowledge of future disturbances may lead the state outside of the maximal RCI set and rely on  
 348 the disturbances to stay safe.

349 The MPC perspective also poses new questions. What role would intermediate costs play if applied  
 350 to policy gradient rollouts? Moreover, if rollouts are interpreted as predictions, could not a learned  
 351 disturbance model, instead of simply the last few disturbance estimates, be included in policy gradient  
 352 rollouts without losing convexity? The lessons also go in the other direction, as most works in robust  
 353 MPC either consider nominal predictions without disturbances, implicitly hoping that disturbances  
 354 average out over time, or defend against the worst case, as in min-max MPC. As a consequence,  
 355 these algorithms perform poorly if disturbances are constant or slowly-time-varying, a setting which  
 356 nonstochastic control (with varying initial states) handles gracefully. Another exciting connection  
 357 to explore is that of nonstochastic control and real-time iterative MPC [Gros et al., 2020], where  
 358 at each time step, the (sub-)optimal input sequence is computed by updating the prior solution,  
 359 instead of recomputing anew. Overall, this work highlights the potential of combining online convex  
 360 optimization-based policy search with robust and predictive control techniques to achieve both safety  
 361 and performance in real-world control systems.

362 **References**

- 363 N. Agarwal, B. Bullins, E. Hazan, S. Kakade, and K. Singh. Online control with adversarial  
364 disturbances. In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th Inter-*  
365 *national Conference on Machine Learning*, volume 97 of *Proceedings of Machine Learning*  
366 *Research*, pages 111–119. PMLR, 09–15 Jun 2019. URL [https://proceedings.mlr.press/](https://proceedings.mlr.press/v97/agarwal19c.html)  
367 [v97/agarwal19c.html](https://proceedings.mlr.press/v97/agarwal19c.html).
- 368 U. M. Ahiyevich, S. E. Parsegov, and P. S. Shcherbakov. Upper bounds on peaks in discrete-time  
369 linear systems. *Automation and Remote Control*, 79(11):1976–1988, Nov 2018. ISSN 1608-3032.  
370 doi: 10.1134/S0005117918110036. URL <https://doi.org/10.1134/S0005117918110036>.
- 371 J. Berberich, J. Köhler, M. A. Müller, and F. Allgöwer. Robust constraint satisfaction in data-driven  
372 mpc. In *2020 59th IEEE Conference on Decision and Control (CDC)*, pages 1260–1267. IEEE,  
373 2020.
- 374 A. Bisoffi, C. De Persis, and P. Tesi. Controller design for robust invariance from noisy data. *IEEE*  
375 *Transactions on Automatic Control*, 68(1):636–643, 2023. doi: 10.1109/TAC.2022.3170373.
- 376 F. Blanchini. Set invariance in control. *Automatica*, 35(11):1747–1767, 1999. ISSN 0005-1098. doi:  
377 [https://doi.org/10.1016/S0005-1098\(99\)00113-2](https://doi.org/10.1016/S0005-1098(99)00113-2). URL [https://www.sciencedirect.com/](https://www.sciencedirect.com/science/article/pii/S0005109899001132)  
378 [science/article/pii/S0005109899001132](https://www.sciencedirect.com/science/article/pii/S0005109899001132).
- 379 F. Blanchini and S. Miani. *Invariant sets*, pages 121–191. Springer International Publishing,  
380 Cham, 2015. ISBN 978-3-319-17933-9. doi: 10.1007/978-3-319-17933-9\_4. URL [https://](https://doi.org/10.1007/978-3-319-17933-9_4)  
381 [doi.org/10.1007/978-3-319-17933-9\\_4](https://doi.org/10.1007/978-3-319-17933-9_4).
- 382 X. Chen and E. Hazan. Black-box control for linear dynamical systems. In M. Belkin and S. Kpotufe,  
383 editors, *Proceedings of Thirty Fourth Conference on Learning Theory*, volume 134 of *Proceedings*  
384 *of Machine Learning Research*, pages 1114–1143. PMLR, 15–19 Aug 2021. URL [https://](https://proceedings.mlr.press/v134/chen21c.html)  
385 [proceedings.mlr.press/v134/chen21c.html](https://proceedings.mlr.press/v134/chen21c.html).
- 386 A. Cohen, A. Hasidim, T. Koren, N. Lazic, Y. Mansour, and K. Talwar. Online linear quadratic  
387 control. In J. Dy and A. Krause, editors, *Proceedings of the 35th International Conference on*  
388 *Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 1029–1038.  
389 PMLR, 10–15 Jul 2018. URL <https://proceedings.mlr.press/v80/cohen18b.html>.
- 390 D. Foster and M. Simchowitz. Logarithmic regret for adversarial online control. In H. D. III and  
391 A. Singh, editors, *Proceedings of the 37th International Conference on Machine Learning*, volume  
392 119 of *Proceedings of Machine Learning Research*, pages 3211–3221. PMLR, 13–18 Jul 2020.  
393 URL <https://proceedings.mlr.press/v119/foster20b.html>.
- 394 P. J. Goulart, E. C. Kerrigan, and J. M. Maciejowski. Optimization over state feedback policies  
395 for robust control with constraints. *Automatica*, 42(4):523–533, 2006. ISSN 0005-1098. doi:  
396 <https://doi.org/10.1016/j.automatica.2005.08.023>. URL [https://www.sciencedirect.com/](https://www.sciencedirect.com/science/article/pii/S0005109806000021)  
397 [science/article/pii/S0005109806000021](https://www.sciencedirect.com/science/article/pii/S0005109806000021).
- 398 S. Gros, M. Zanon, R. Quirynen, A. Bemporad, and M. Diehl. From linear to nonlinear mpc: bridging  
399 the gap via the real-time iteration. *International Journal of Control*, 93(1):62–80, 2020. doi:  
400 [10.1080/00207179.2016.1222553](https://doi.org/10.1080/00207179.2016.1222553).
- 401 E. Hazan and K. Singh. Introduction to online nonstochastic control, 2023.
- 402 E. Hazan, S. Kakade, and K. Singh. The nonstochastic control problem. In A. Kontorovich and  
403 G. Neu, editors, *Proceedings of the 31st International Conference on Algorithmic Learning Theory*,  
404 volume 117 of *Proceedings of Machine Learning Research*, pages 408–421. PMLR, 08 Feb–11  
405 Feb 2020. URL <https://proceedings.mlr.press/v117/hazan20a.html>.
- 406 Y. Li, S. Das, and N. Li. Online optimal control with affine constraints. *Proceedings of the AAAI*  
407 *Conference on Artificial Intelligence*, 35(10):8527–8537, May 2021. doi: 10.1609/aaai.v35i10.  
408 17035. URL <https://ojs.aaai.org/index.php/AAAI/article/view/17035>.
- 409 Y. Li, J. Yu, L. Conger, and A. Wierman. Learning the uncertainty sets for control dynamics via set  
410 membership: A non-asymptotic analysis, 2023.

- 411 Y. Li, S. Das, J. Shamma, and N. Li. Safe adaptive learning-based control for linear quadratic  
412 regulators constraints, 2024.
- 413 X. Liu, Z. Yang, and L. Ying. Online nonstochastic control with adversarial and static constraints. In  
414 A. Krause, E. Brunskill, K. Cho, B. Engelhardt, S. Sabato, and J. Scarlett, editors, *Proceedings of*  
415 *the 40th International Conference on Machine Learning*, volume 202 of *Proceedings of Machine*  
416 *Learning Research*, pages 22277–22288. PMLR, 23–29 Jul 2023. URL <https://proceedings.mlr.press/v202/liu23at.html>.
- 418 M. Lorenzen, F. Dabbene, R. Tempo, and F. Allgöwer. Constraint-tightening and stability in stochastic  
419 model predictive control. *IEEE Transactions on Automatic Control*, 62(7):3165–3177, 2016.
- 420 M. Lorenzen, M. Cannon, and F. Allgöwer. Robust MPC with recursive model update. *Automatica*,  
421 103:461–471, 2019.
- 422 A. Martin, L. Furieri, F. Dörfler, J. Lygeros, and G. Ferrari-Trecate. Regret optimal control for  
423 uncertain stochastic systems, 2023.
- 424 E. Minasyan, P. Gradu, M. Simchowitz, and E. Hazan. Online control of unknown time-varying  
425 dynamical systems. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P. Liang, and J. W. Vaughan,  
426 editors, *Advances in Neural Information Processing Systems*, volume 34, pages 15934–15945. Cur-  
427 ran Associates, Inc., 2021. URL [https://proceedings.neurips.cc/paper\\_files/paper/](https://proceedings.neurips.cc/paper_files/paper/2021/file/856b503e276cc491e7e6e0ac1b9f4b17-Paper.pdf)  
428 [2021/file/856b503e276cc491e7e6e0ac1b9f4b17-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/856b503e276cc491e7e6e0ac1b9f4b17-Paper.pdf).
- 429 M. Nonhoff, E. Dall’Anese, and M. A. Müller. Online convex optimization for robust control of  
430 constrained dynamical systems, 2024.
- 431 J. B. Rawlings, D. Q. Mayne, and M. Diehl. *Model Predictive Control: Theory, Computation, and*  
432 *Design*. Nob Hill Publishing, Madison, Wisconsin, 2 edition, 2017. ISBN 9780975937730.
- 433 T. Sarkar and A. Rakhlin. Near optimal finite time identification of arbitrary linear dynamical systems.  
434 In K. Chaudhuri and R. Salakhutdinov, editors, *Proceedings of the 36th International Conference*  
435 *on Machine Learning*, volume 97 of *Proceedings of Machine Learning Research*, pages 5610–5618.  
436 PMLR, 09–15 Jun 2019. URL <https://proceedings.mlr.press/v97/sarkar19a.html>.
- 437 J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization  
438 algorithms, 2017.
- 439 M. Simchowitz. Making non-stochastic control (almost) as easy as stochastic. In *Proceedings of the*  
440 *34th International Conference on Neural Information Processing Systems*, NIPS’20, Red Hook,  
441 NY, USA, 2020. Curran Associates Inc. ISBN 9781713829546.
- 442 M. Simchowitz, H. Mania, S. Tu, M. I. Jordan, and B. Recht. Learning without mixing: Towards a  
443 sharp analysis of linear system identification. In S. Bubeck, V. Perchet, and P. Rigollet, editors,  
444 *Proceedings of the 31st Conference On Learning Theory*, volume 75 of *Proceedings of Machine*  
445 *Learning Research*, pages 439–473. PMLR, 06–09 Jul 2018. URL [https://proceedings.mlr.](https://proceedings.mlr.press/v75/simchowitz18a.html)  
446 [press/v75/simchowitz18a.html](https://proceedings.mlr.press/v75/simchowitz18a.html).
- 447 M. Simchowitz, K. Singh, and E. Hazan. Improper learning for non-stochastic control. In J. Abernethy  
448 and S. Agarwal, editors, *Proceedings of Thirty Third Conference on Learning Theory*, volume 125  
449 of *Proceedings of Machine Learning Research*, pages 3320–3436. PMLR, 09–12 Jul 2020. URL  
450 <https://proceedings.mlr.press/v125/simchowitz20a.html>.
- 451 Y. J. Sun, S. Newman, and E. Hazan. Optimal rates for bandit nonstochastic control. In  
452 A. Oh, T. Naumann, A. Globerson, K. Saenko, M. Hardt, and S. Levine, editors, *Advances in*  
453 *Neural Information Processing Systems*, volume 36, pages 21908–21919. Curran Associates,  
454 Inc., 2023. URL [https://proceedings.neurips.cc/paper\\_files/paper/2023/file/](https://proceedings.neurips.cc/paper_files/paper/2023/file/45591d6727f0e127295f8d16adba6b23-Paper-Conference.pdf)  
455 [45591d6727f0e127295f8d16adba6b23-Paper-Conference.pdf](https://proceedings.neurips.cc/paper_files/paper/2023/file/45591d6727f0e127295f8d16adba6b23-Paper-Conference.pdf).
- 456 J. Teutsch, S. Kerz, D. Wollherr, and M. Leibold. Sampling-based stochastic data-driven predictive  
457 control under data uncertainty. *arXiv preprint arXiv:2402.00681*, 2024.

458 H. J. Van Waarde, J. Eising, M. K. Camlibel, and H. L. Trentelman. The informativity approach:  
 459 To data-driven analysis and control. *IEEE Control Systems Magazine*, 43(6):32–66, 2023. doi:  
 460 10.1109/MCS.2023.3310305.

461 A. Wagenmaker and K. Jamieson. Active learning for identification of linear dynamical systems. In  
 462 J. Abernethy and S. Agarwal, editors, *Proceedings of Thirty Third Conference on Learning Theory*,  
 463 volume 125 of *Proceedings of Machine Learning Research*, pages 3487–3582. PMLR, 09–12 Jul  
 464 2020. URL <https://proceedings.mlr.press/v125/wagenmaker20a.html>.

465 H. Zhou and V. Tzoumas. Safe non-stochastic control of linear dynamical systems. In *IEEE*  
 466 *Conference on Decision and Control (CDC)*, 2023. URL [https://arxiv.org/pdf/2308.](https://arxiv.org/pdf/2308.12395.pdf)  
 467 12395.pdf.

## 468 A Appendix

469 Variants of the following lemma are well known in the system identification literature. The following  
 470 version is adapted from Bisoffi et al. [2023].

471 **Lemma 2.** *The set of consistent models  $\Omega$  is convex and closed. It is bounded if and only if its*  
 472 *generating data satisfies  $\text{rank} \begin{bmatrix} x_0 & \dots & x_{N-1} \\ u_0 & \dots & u_{N-1} \end{bmatrix} = n_x + n_u$ .*

473 In practice, the rank condition of Lemma 2 is easily satisfied by long enough trajectories with random  
 474 inputs.

475 The following is a classical result in control due to Lyapunov.

476 **Proposition 4.** *A system  $x_{t+1} = Ax_t$  is stable in the sense that  $\lim_{t \rightarrow \infty} x_t = 0$  if and only if there*  
 477 *exists  $P \succ 0$  such that*

$$P - A^T P A \succ 0. \quad (21)$$

478 Proposition 4 implies the existence of a scalar *Lyapunov* function  $V(x) = x^T P x$  which attains  
 479 its minimum at the origin ( $V(x) > 0$  for all  $x \neq 0$  and  $V(0) = 0$ ) and descents with time  
 480 ( $V(x_{t+1}) < V(x_t)$ ) until  $x_t = 0$  since for all  $x_t \neq 0$  the condition 21 guarantees

$$V(x_{t+1}) - V(x_t) = x_t^T A^T P A x_t - x_t^T P x_t = x_t^T (A^T P A - P) x_t < 0. \quad (22)$$

481 Informally, this implies  $\lim_{t \rightarrow \infty} V(x_t) = V(x_{t \rightarrow \infty}) = \min_x V(x) = V(0)$  and the state tends to  
 482 the origin.

483 **Construction of the maximal RCI subset** A maximal RCI subset  $\mathcal{X}$  of  $\mathcal{F}_x$  can be constructed by  
 484 recursion [Blanchini and Miani, 2015], where the idea is to first set  $\mathcal{X}_0 = \mathcal{F}_x$  and iteratively compute  
 485  $\mathcal{X}_{k+1}$  as the set of all states from which  $\mathcal{X}_k$  can be surely reached (for all disturbances in  $\mathcal{W}$ ). That is,  
 486  $\mathcal{X}_{k+1}$  contains all states for which there exists an admissible input which drives the nominal state  
 487 (without disturbance)  $A_* x + B_* u$  into  $\mathcal{X}_k \ominus \mathcal{W}$ ,

$$\mathcal{X}_{k+1} = \text{proj}_{1:n_x} \{z \in \text{col}(\mathcal{X}_k, \mathcal{U}) \mid \Omega z \in \mathcal{X}_k \ominus \mathcal{W}\}. \quad (23)$$

488 Crucially,  $x \in \mathcal{X}_{k+1}$  guarantees the existence of *one* input that drives *all* models of  $\Omega$  into  $\mathcal{X}_k \ominus \mathcal{W}$   
 489 and may be computed similar to  $\mathcal{F}_x$  above based on vertices of  $\Omega$ , yielding again a convex polytope.  
 490 Note that  $\mathcal{X}_{k+1} \subseteq \mathcal{X}_k$  by construction. As soon as  $\mathcal{X}_{k+1} = \mathcal{X}_k$  the computation is stopped and  
 491  $\mathcal{X} := \mathcal{X}_k$  is RCI for the true system following (12).

492 **Proof of Proposition 1** The proof makes use the well-known fact that sublevel sets of Lyapunov  
 493 functions are positive invariant, which we formally define next before proving the result.

494 **Definition 5.** *A set  $\mathcal{X}$  is positive invariant for dynamics  $x_{t+1} = f(x_t)$  if  $f(x) \in \mathcal{X}$  for all  $x \in \mathcal{X}$ .*

495 **Lemma 3.** *Consider a system  $x_{t+1} = f(x_t)$  and let  $V(x)$  be a Lyapunov function such that*  
 496  *$V(0) = 0$ ,  $V(x) > 0 \forall x \neq 0$ , and  $V(f(x)) \leq V(x) \forall x \in \mathbb{R}^{n_x}$ . Then any sublevel set*  
 497  *$\mathcal{E}_c = \{x \in \mathbb{R}^{n_x} \mid V(x) \leq c, c \geq 0\}$  of  $V(x)$  is positive invariant for dynamics  $x_{t+1} = f(x_t)$ .*

498 *Proof.* We first show that condition 17 implies stability of of

$$x_{t+1} = \frac{1}{r}(A + BK)x_t. \quad (24)$$

499 By the Schur complement, it is equivalent to  $rZ \succ 0$  and  $rZ - (\hat{A}^{(i)}Z + \hat{B}^{(i)}Y)^T(rZ)^{-1}(\hat{A}^{(i)}Z +$   
500  $\hat{B}^{(i)}Y) \succ 0$ . Multiplying from left and right by  $Z^{-1}$  yields

$$rZ^{-1} - (\hat{A}^{(i)} + \hat{B}^{(i)}YZ^{-1})^T \frac{1}{r} Z^{-1} (\hat{A}^{(i)} + \hat{B}^{(i)}YZ^{-1}) \succ 0. \quad (25)$$

501 Substituting  $K = YZ^{-1}$  and  $P = Z^{-1}$  and dividing by  $r$  leads to

$$P - \frac{1}{r} (\hat{A}^{(i)} + \hat{B}^{(i)}K)^T P \frac{1}{r} (\hat{A}^{(i)} + \hat{B}^{(i)}K) \succ 0 \quad (26)$$

502 so that  $V(x) = x^T P x$  is a Lyapunov function certifying stability for each closed loop system  
503  $x_{t+1} = \frac{1}{r} (\hat{A}^{(i)} + \hat{B}^{(i)}K)x_t$  by Proposition 4. By convexity of  $\Omega$ , stability of the models at the  
504 vertices of  $\Omega$  implies stability for all models in  $\Omega$ , which in turn implies  $\rho(\frac{1}{r}(A + BK)) \leq 1$  or  
505 equivalently  $\rho(A + BK) \leq r$  for all  $(A, B) \in \Omega$ .

506 Since  $V(x_{t+1}) \leq V(x_t)$  the ellipsoidal sublevel set  $\mathcal{V} = \{x \in \mathbb{R}^{n_x} \mid x^T P x \leq 1\}$  of  $V(x)$  is positive  
507 invariant for system 24, i.e.,  $x_0 \in \mathcal{V} \implies r^{-t}(A + BK)^t x_0 \in \mathcal{V}$  for all time  $t \geq 0$ . Multiplying all  
508 sides in condition 16 by  $P = Z^{-1}$  yields  $P \preceq I$  and  $I \preceq cP$  which in turn implies  $x^T x \leq x^T P x$   
509 and  $x^T P x \leq \frac{1}{c} x^T x$ . As a consequence, whenever  $x^T x \leq 1$  then  $x^T P x \leq 1$  and thus  $\mathcal{V}$  contains  
510 the unit norm ball  $\mathcal{B}_1 = \{x \in \mathbb{R}^{n_x} \mid \|x\| = \sqrt{x^T x} \leq 1\}$ . Additionally, whenever  $x^T P x \leq 1$  then  
511  $\frac{1}{c} x^T x \leq 1$  (and equivalently  $x^T x \leq c$ ) so that  $\mathcal{V}$  is contained in a ball around the origin with radius  
512  $\sqrt{c}$  denoted by  $\mathcal{B}_{\sqrt{c}} = \{x \in \mathbb{R}^{n_x} \mid \|x\| \leq \sqrt{c}\}$ . By positive invariance of  $\mathcal{V}$  then,  $x \in \mathcal{B}_1 \subseteq \mathcal{V}$   
513 implies  $r^{-t}(A + BK)^t x \in \mathcal{V} \subseteq \mathcal{B}_{\sqrt{c}}$  for all time  $t \geq 0$ . In other words

$$\|x\| \leq 1 \implies \frac{1}{r^t} \|(A + BK)^t x\| \leq \sqrt{c} \iff \sup_{\|x\| \leq 1} \frac{\|(A + BK)^t x\|}{\|x\|} \leq \sqrt{c} r^t \quad (27)$$

514 which proves the last part of the result by definition of the induced matrix norm.  $\square$

### 515 Proof of Proposition 3

516 *Proof.* Rearranging  $x^* = x_{H|t} = A_K^H x_{0|t} + S_{H-1} B v + S_{H-1} w$  with  $x_{0|t} = x_t$  leads to

$$B v_t = S_{H-1}^{-1} (x^* - A_K^H x_t) - w, \quad x_{t+1} = (A_K - S_{H-1}^{-1} A_K^H) x_t + S_{H-1}^{-1} x^*. \quad (28)$$

517 Note  $(A_K - S_{H-1}^{-1} A_K^H) = S_{H-1}^{-1} (S_{H-1} A_K - A_K^H)$  and

$$S_{H-1} A_K - A_K^H = (I + \dots + A_K^{H-1}) A_K - A_K^H = A_K + \dots + A_K^{H-1} = S_{H-1} - I \quad (29)$$

518 so that the dynamics induced by MPC can be rewritten as

$$x_{t+1} = (I - S_{H-1}^{-1}) x_t + S_{H-1}^{-1} x^*. \quad (30)$$

519 Letting  $x_t = x_{t+1} = x_\infty$  immediately leads to  $S_{H-1}^{-1} x_\infty = S_{H-1}^{-1} x^*$  which implies  $x_\infty = x^*$  since  
520  $S_{H-1}^{-1}$  has full rank. It remains to show that  $x_t$  actually converges, i.e., the closed-loop dynamics  
521 (30) are stable. Let  $\lambda$  be an eigenvalue of  $A_K$  such that  $A_K v = \lambda v$  for some  $v \in \mathbb{R}^{n_x}$ . Then  
522  $S_{H-1} v = (1 + \lambda + \dots + \lambda^{H-1}) v$  so that  $(1 + \lambda + \dots + \lambda^{H-1})^{-1}$  is an eigenvalue of  $S_{H-1}^{-1}$  and  
523  $1 - (1 + \lambda + \dots + \lambda^{H-1})^{-1} = \frac{\lambda - \lambda^H}{1 - \lambda^H}$  is an eigenvalue of the closed-loop dynamics (30). Since  
524  $\frac{\lambda - \lambda^H}{1 - \lambda^H} \in [0, \lambda)$  for all  $H \in \mathbb{N}$  we have  $\rho(I - S_{H-1}^{-1}) < \rho(A_K)$  and the closed-loop dynamics are  
525 stable by (strong) stability of  $A_K$ .  $\square$