

# A Proximal-Point Lagrangian Based Parallelizable Nonconvex Solver for Bilinear Model Predictive Control

Yingzhao Lian, *Graduate Student Member, IEEE*, Yuning Jiang, *Member, IEEE*, Daniel F. Oplia, *Senior Member, IEEE*, and Colin N. Jones, *Senior Member, IEEE*

**Abstract**—Nonlinear model predictive control has been widely adopted to manipulate bilinear systems with dynamics that include products of the inputs and the states. These systems are ubiquitous in chemical processes, mechanical systems, and quantum physics, to name a few. Running a bilinear MPC controller in real time requires solving a non-convex optimization problem within a limited sampling time. This paper proposes a novel parallel proximal-point Lagrangian based bilinear MPC solver via an interlacing horizon-splitting scheme. The resulting algorithm converts the non-convex MPC control problem into a set of parallelizable small-scale multi-parametric quadratic programs (mpQPs) and an equality-constrained linear-quadratic regulator problem. As a result, the solutions of mpQPs can be pre-computed offline to enable efficient online computation. The proposed algorithm is validated on a simulation of an HVAC system control. It is deployed on a TI LaunchPad XL F28379D microcontroller to execute speed control on a field-controlled DC motor, where the MPC updates at 10 ms and solves the problem in 1.764 ms on average and at most 2.088 ms.

## I. INTRODUCTION

Bilinear systems were originally introduced in [1], [2] to model systems where the dynamics involve products of the inputs and the states. These dynamics may result from linearizing a nonlinear input affine system and are most commonly used to model convection and spinning in chemical processes and mechanical systems [3], [4]. Additionally, using the concept of Carleman linearization [5], it has been shown that bilinear systems can model general nonlinear systems [6]. Meanwhile, with the help of various sophisticated tools such as Lie algebra [7, Chapter 2] and Volterra series [8], bilinear control theory has been explored in-depth and has found various successful applications [9]–[13].

Nonlinear model predictive control (NMPC) is one of the most successful approaches to control bilinear systems [14]–[17]. The main idea of NMPC is to achieve the desired performance by optimizing the input in a receding horizon scheme while enforcing state and input constraints [18]. This requires the solution of a nonlinear optimal control problem

This work has received support from the Swiss National Science Foundation under the RISK project (Risk Aware Data-Driven Demand Response, grant number 200021\_175627) and the NCCR Automation project (grant agreement 51NF40\_180545). (Corresponding author: Yuning Jiang)

Yingzhao Lian, Yuning Jiang, and Colin N. Jones are with the Automatic Control Lab, EPFL, Switzerland. (e-mail: yingzhao.lian, colin.jones@epfl.ch, yuning.jiang@ieee.org)

Daniel F. Oplia is with the United States Naval Academy, Electrical and Computer Engineering Department. (email: oplia@usna.edu)

(OCP) online within a limited update time. Therefore, an efficient solver is critical to running NMPC in real time<sup>1</sup>.

Among various real-time NMPC methods, designing and executing an online solver that can run in parallel via distributed algorithms has been a trend over the past decade [19]. Compared to centralized solution approaches, parallelizable methods split the problem into multiple smaller problems such that the computational resources can be utilized more efficiently by exploiting the structure of the OCP being solved. A classical approach used in distributed optimization is based on dual decomposition, where, for example, a gradient-based method [20], [21] or a semi-smooth Newton method [22] have been used to solve the concave dual problem. Another famous approach is the Alternating Direction Method of Multipliers, which parallelizes the computation by introducing auxiliary variables [19], [23]. These two methods lack convergence guarantees for nonconvex problems and hence are only formally applicable to linear systems. In [24], an augmented Lagrangian based distributed optimization algorithm is proposed, which has been applied to parallelize the computation of MPC problems in [25], [26]. However, despite being parallelized, these algorithms require a solution to multiple non-convex optimization problems in each iteration, which are still numerically intense.

Decomposing an NMPC problem into a set of small-scale problems mainly leverages the linear equality constraints that appear in NMPC problems, which can reflect the topology of a network system or that naturally emerge in the temporal direction via the introduction of auxiliary variables. The latter approach is the horizon splitting method [27], [28], or sometimes termed Schwarz decomposition [29]. It splits the predictive trajectories into short sequential sequences, where linear couplings naturally enforce the equality between the initial and terminal states of two adjacent short sequences, hence the name. Within the scope of horizon splitting, tools beyond distributed algorithms have been leveraged to improve efficiency further. The banded structure of the KKT system is the most investigated object in this setup. A binary-tree-structured algorithm summarizes In [30], a general parallel solver, and in [28], an approximation scheme is introduced to develop a parallel Ricatti solver. However, these algorithms still handle the nonconvex problem directly and, as such, are still numerically challenging.

<sup>1</sup>In this work, real-time means that the MPC solver should return the solution fast enough to enable a desirable operation of the targeted system. Based on our experience, for a mechatronic/mechanical system, the MPC solver should be at least five times faster than the sampling frequency.

Another category of methods widely used in real-time NMPC leverages the super-linear local convergence property of Newton-type methods to accelerate online convergence, given a good initialization of the decision variables. This category of methods roughly defines the “warm-start” strategy, whose initialization usually derives from the solution information gathered from the preceding time step. A basic approach directly shifts the solution from the last iteration [31], and then a Newton iteration ensures efficient local convergence. Under the umbrella of sequential quadratic programming (SQP), the sensitivity information of the local solution is further used to initialize the KKT system, where an initial guess of active constraints is the most challenging object. In [32], the piecewise affine property of linear model predictive control (MPC) is used to estimate the change of the active constraint. This idea is generalized in [33] under the name of real-time iteration (RTI), where a sensitivity analysis of the local solution is used to give a piece-wise affine update of the control law.

Instead of solving the NMPC directly online, explicit MPC shifts the online computational burden offline. It treats the MPC control law as a nonlinear mapping from the initial state to control input, and this control law is precomputed offline to enable efficient online calls. In a linear MPC setup, the optimal control law is locally affine [34], [35], and this piecewise affine parametric solution is first used to pre-compute the MPC control law offline in [35]. However, this algebraic property only holds for linear systems, and its application to nonlinear MPC is limited without approximation [36].

This work proposes a new proximal-point Lagrangian based algorithm, which combines the ideas of horizon splitting, explicit MPC, and real-time SQP. In contrast to a standard horizon-splitting approach, a novel interlacing horizon-splitting scheme is introduced. The advantages of the proposed controller are summarized as follows:

- 1) The proposed algorithm runs computationally efficient iterations, which only require an evaluation of a multi-parametric QP (mpQP) solution and to solve a sparse linear equation system.
- 2) The detection of the active set is shifted to the mpQP solution, whose problem size is independent of the prediction horizon.
- 3) A novel interlacing horizon splitting scheme is introduced. The resulting problem has the same number of decision variables as the original NMPC problem without introducing auxiliary variables.
- 4) The proposed algorithm will not abort even when an infeasible initial state is given. It will output a solution that at least satisfies the input constraint.

After introducing notation and background knowledge in the rest of this Section I, the bilinear MPC control problem is presented in Section II, after which the parallelizable non-convex solver is proposed in Section III. In particular, Section III-A introduces a novel interlacing horizon splitting scheme, based on which the solver is detailed in Section III-B. Convergence properties of the proposed solver are studied in Section III-C. After introducing the proposed algorithm, a dual space interpretation of the proposed algorithm is given in

Section III-D, after which a comparison with related results follows in Section III-E. The numerical details of the proposed algorithm are investigated in Section IV. The efficacy of the proposed algorithm is studied in Section V, where the efficient real-time MPC solver *acados* and the nonconvex parallel primal-dual solver augmented Lagrangian based alternating direction inexact Newton (ALADIN) method are used as a benchmark. Meanwhile, the proposed algorithm is deployed on a Texas Instruments C2000 Delfino LaunchPad XL F28379 microcontroller to control a field-controlled DC motor.

## A. Notation

We use the symbols  $\mathbb{S}_+^n$  and  $\mathbb{S}_{++}^n$  to denote the set of symmetric, positive semi-definite, and symmetric, positive definite matrices in  $\mathbb{R}^{n \times n}$ . For a given matrix  $\Sigma \in \mathbb{S}_{++}^n$  the notation  $\|x\|_\Sigma = \sqrt{x^\top \Sigma x}$  is used, and  $\|x\|$  denotes the Euclidean norm. Moreover, a function  $c : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$  is called strongly convex with matrix parameter  $\Sigma \in \mathbb{S}_+^n$ , if the inequality

$$c(tx + (1-t)y) \leq tc(x) + (1-t)c(y) - \frac{1}{2}t(1-t)\|x-y\|_\Sigma^2$$

is satisfied for all  $x, y \in \mathbb{R}^n$  and all  $t \in [0, 1]$ . Notice that all convex functions in this paper are assumed to be closed and proper [37]. For a vector  $x \in \mathbb{R}^n$ , we denote by  $[x]_i$  its  $i$ -th element. Set  $\mathbb{Z}_i^j$  denotes the range of integers from  $i$  to  $j$  with  $i \leq j$ . The Kronecker product of two matrices  $A \in \mathbb{R}^{k \times l}$  and  $B \in \mathbb{R}^{m \times n}$  is:

$$\mathbb{R}^{km \times ln} \ni A \otimes B := \begin{bmatrix} a_{1,1}B & a_{1,2}B & \cdots & a_{1,l}B \\ \vdots & \vdots & \ddots & \vdots \\ a_{k,1}B & a_{k,2}B & \cdots & a_{k,l}B \end{bmatrix}.$$

$\text{vec}(A)$  denotes the vector that is obtained by stacking all columns of  $A$  into one long vector. The reverse operation is denoted by  $\text{mat}(\cdot)$ , such that  $\text{mat}(\text{vec}(A)) = A$ . The identity matrix in  $\mathbb{R}^{n \times n}$  is denoted by  $\mathbf{I}_n$  and the zero matrix in  $\mathbb{R}^{m \times n}$  is denoted by  $\mathbf{0}_{m \times n}$ . Notation  $\text{diag}(H_1, \dots, H_n)$  constructs a block-diagonal matrix whose  $i$ -th diagonal block is  $H_i$ . For a given function  $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , we use the Landau notation

$$f(x) = \mathbf{O}(\|x\|), \quad \text{if } \exists c \in \mathbb{R}, \lim_{x \rightarrow 0} \frac{f(x)}{\|x\|} = c.$$

## B. Preliminaries

We first recap some existing results from the field of multi-parametric quadratic programming (mpQP) used later in this paper. A generic convex mpQP can be written in the form of

$$\begin{aligned} \min_x \quad & \frac{1}{2}x^\top Qx + \theta^\top Sx & (1a) \\ \text{s.t.} \quad & Ax \leq b + C\theta, & (1b) \end{aligned}$$

with decision variables  $x \in \mathbb{R}^{n_x}$  and parameters  $\theta \in \mathbb{R}^{n_p}$ . Here, matrices  $Q \in \mathbb{S}_+^{n_x}$ ,  $S \in \mathbb{R}^{n_p \times n_x}$ ,  $A \in \mathbb{R}^{m \times n_x}$ ,  $C \in \mathbb{R}^{m \times n_p}$  and vector  $b \in \mathbb{R}^m$  are given data. Moreover, we denote by  $\Omega$  the set of all parameters  $\theta$  for which (1) is feasible. For a mpQP (1) with a strongly convex value function, it has been shown (see, e.g., [38]) that  $\Omega$  is a

polyhedron while the solution map  $x^*(\theta) : \mathbb{R}^{n_p} \rightarrow \mathbb{R}^{n_x}$  is a continuous piecewise affine (PWA) function of the parameters. Each affine piece is called a critical region [39, Chapter 7.1.2]. Meanwhile, the Lipschitz-continuity holds at  $x^*(\cdot)$ , i.e., there exists a positive constant  $\eta > 0$  such that for any  $\theta_1, \theta_2 \in \Omega$ , we have

$$\|x^*(\theta_1) - x^*(\theta_2)\| \leq \eta \|\theta_1 - \theta_2\|. \quad (2)$$

We now recall some definitions from the field of nonlinear programming (NLP). Let us consider NLPs in a generic form

$$\min_x f(x) \quad \text{subject to} \quad \begin{cases} g(x) = 0 & | \lambda, \\ h(x) \leq 0 & | \kappa. \end{cases} \quad (3)$$

Throughout the rest of this paper, we write Lagrangian multipliers right after the constraints such that  $\lambda \in \mathbb{R}^{n_g}$  and  $\mathbb{R}^{n_h} \ni \kappa \geq 0$  denote, respectively, the Lagrangian multipliers of the equality constraints and inequality constraints. Functions  $f : \mathbb{R}^{n_x} \rightarrow \mathbb{R}$ ,  $g : \mathbb{R}^{n_x} \rightarrow \mathbb{R}^{n_g}$  and  $h : \mathbb{R}^{n_x} \rightarrow \mathbb{R}^{n_h}$  are assumed twice continuously differentiable. A primal-dual solution  $(x^*, \lambda^*, \kappa^*)$  is called a Karush–Kuhn–Tucker (KKT) point of (3) if the following conditions are satisfied [40, Chapter 12.3]

$$\nabla f(x^*) + \nabla g(x^*)\lambda^* + \nabla h(x^*)\kappa^* = 0, \quad (4a)$$

$$g(x^*) = 0, \quad h(x^*) \leq 0, \quad (4b)$$

$$\forall i \in 1, \dots, n_h, \quad [\kappa^*]_i \cdot [h(x^*)]_i = 0, \quad [\kappa^*]_i \geq 0, \quad (4c)$$

where  $[\kappa^*]_i$  and  $[h(x^*)]_i$  define the  $i$ -th element of  $\kappa^*$  and  $h(x^*)$ , respectively. For a given feasible  $x$ , we denote by  $\mathcal{A}(x)$  the active set at  $x$ , i.e., the index set that includes the equality constraints and the inequality constraints that holds equality at  $x$ . When the set of active constraint gradients (i.e.,  $[\nabla g(x), \nabla h_{i \in \mathcal{A}(x)}(x)]$ ) is linearly independent at point  $x$ , the linear constraint qualification (LICQ) holds [40, Chapter 12.2]. Furthermore, we say the second-order sufficient condition (SOSC) holds at point  $x$  if its hessian  $\nabla^2 h(x)$  is positive definite semidefinite on the null space spanned by active constraint gradients [41]. Finally, we say the strict complementary condition (SCC) holds if a dual variable equals zero only when the corresponding constraint is inactive [40, Definition 12.5]. Then, we state the definition of regular KKT point for NLP (3).

**Definition 1** [41] *A given KKT point  $(x^*, \lambda^*, \kappa^*)$  is called a regular KKT point if the LICQ, SOSC, and the SCC hold.*

For a given KKT point  $(x^*, \lambda^*, \kappa^*)$ , if it is regular, then there exists an open neighborhood  $\mathcal{B}(x^*)$  around  $x^*$  such that the active set is fixed for any  $x \in \mathcal{B}(x^*)$ , (i.e.,  $\mathcal{A}(x) = \mathcal{A}(x^*)$ ) [41]). Regularity at KKT points guarantees the local convergence property when a Newton-type method is applied to solve (3) [40].

When the inequality constraint  $h(x) \leq 0$  defines a convex set, the first-order optimality condition (4) can be further simplified for the sake of compactness:

$$0 \in \nabla f(x^*) + \nabla g(x^*)\lambda^* + \mathcal{N}_{\mathcal{X}}(x^*),$$

with the convex set  $\mathcal{X} := \{x \in \mathbb{R}^{n_x} | h(x) \leq 0\}$  and  $\mathcal{N}_{\mathcal{X}}(x^*) := \{y | (y - x^*)^\top (x - x^*) \leq 0, \forall x \in \mathcal{X}\}$  denotes the normal cone of convex set  $\mathcal{X}$  at  $x^*$ .

In contrast to the standard Hestenes-Powell augmented Lagrangian method [42], [43], a variant of an augmented Lagrangian method, termed proximal-point Lagrangian [44], is used in this work. Given an equality constraint optimization problem

$$\min_x f(x) \quad \text{s.t.} \quad g(x) = 0,$$

its linearized proximal-point Lagrangian around  $\bar{x}$  is defined by

$$\mathcal{L}^\rho(x, \lambda, \bar{x}) := f(x) + (\nabla g(\bar{x})\lambda)^\top x + \frac{\rho}{2}\|x - \bar{x}\|^2. \quad (5)$$

Note that when  $\rho = 0$ , it recovers the linearized Lagrangian. For the sake of completeness, the Hestenes-Powell augmented Lagrangian is defined by  $f(x) + \lambda^\top g(x) + \frac{\rho}{2}\|g(x)\|^2$ .

## II. PROBLEM FORMULATION

This paper considers discrete-time, time-invariant bilinear dynamic systems:

$$x_{k+1} = Ax_k + Bu_k + \sum_{i=1}^{n_u} C_i x_k [u_k]_i + B_w w_k \quad (6)$$

with state  $x_k \in \mathbb{R}^{n_x}$ , control inputs  $u_k \in \mathbb{R}^{n_u}$  and disturbance  $w_k$  at time instant  $k$ . For the sake of simplicity, we group the bilinear coefficient matrices  $C = [C_1^\top, \dots, C_{n_u}^\top]^\top \in \mathbb{R}^{n_x \times n_u \times n_x}$  and assume that the states and control inputs are subject to the polyhedral constraints

$$x_k \in \mathcal{X} := \{x \in \mathbb{R}^{n_x} | P_x x \leq p_x\},$$

$$\text{and } u_k \in \mathcal{U} := \{u \in \mathbb{R}^{n_u} | P_u u \leq p_u\}.$$

The dynamics (6) also includes the case without process noise. For the case with process noise, to make the problem tractable, we consider solving a certainty equivalent problem, where the prediction of the nominal process noise  $w_k$  is available. This assumption holds in many energy-related applications: solar radiation in photovoltaic power systems, the outdoor climate in building control, and power generation in airborne wind energy systems, to name a few. An MPC controller can be designed by recursively solving the following optimal control problem in a receding horizon fashion,

$$\min_{x_0, x, u} \sum_{k=0}^{N-1} \ell_k(x_k, u_k) + \ell_N(x_N) \quad (7a)$$

subject to:

$$x_0 = x(t), \quad (7b)$$

$$\forall k \in \{0, 1, \dots, N-1\},$$

$$x_{k+1} = Ax_k + Bu_k + \sum_{i=1}^{n_u} C_i x_k [u_k]_i + B_w w_k \quad | \lambda_k, \quad (7c)$$

$$x_{k+1} \in \mathcal{X}, \quad u_k \in \mathcal{U} \quad | \kappa_k, \quad (7d)$$

with  $x = [x_1^\top, \dots, x_N^\top]^\top$ ,  $u = [u_0^\top, \dots, u_{N-1}^\top]^\top$  and prediction horizon  $N \in \mathbb{Z}_{>0}$ . For the sake of consistency, variables indexed by bracketed time, such as  $x(t)$ , denote the actual measurements read out from sensors. Meanwhile, variables indexed by a subscript, such as  $x_k$ , denote the predictive “virtual” variables used in the MPC problem. In problem (7),

the stage cost  $\ell_k(\cdot, \cdot)$ ,  $k \in \mathbb{Z}_0^{N-1}$  and terminal cost  $\ell_N(\cdot)$  are quadratic and strongly convex, i.e.,

$$\begin{aligned}\ell_k(x, u) &= \frac{1}{2}x^\top Q_k x + q_k^\top x + \frac{1}{2}u^\top R_k u + r_k^\top u, \\ \ell_N(x) &= \frac{1}{2}x^\top Q_N x + q_N^\top x\end{aligned}$$

with user-defined parameters  $Q, Q_N \in \mathbb{S}_{++}^{n_x}$ ,  $R \in \mathbb{S}_{++}^{n_u}$ ,  $q_k \in \mathbb{R}^{n_x}$ ,  $r_k \in \mathbb{R}^{n_u}$ . Notice that although its objective is strongly convex, solving the nonconvex Problem (7) is challenging due to the bilinear dynamics (7c).

### III. ALGORITHM DEVELOPMENT

In this work, we propose an algorithm to solve the bilinear MPC problem (7) efficiently. Before delving into the algorithmic details, we would first state the logic behind the design of the proposed algorithm. As reviewed in Section I, real-time MPC mainly applies an SQP solver with warm-start strategies or uses explicit MPC. When a good initialization is not available, detecting active inequality constraints becomes the major performance bottleneck for the SQP algorithm. This requires the design of sophisticated, active set detection strategies or the use of merit functions (see e.g. [40, Chapter 18.2] [45, Chapter 2.3] [46, Chapter 2.3]). In the worst case, a poor estimate of the active set will lead to an infeasible QP subproblem, which ultimately aborts the progress of the SQP algorithm. On the contrary, the information of the active set is implicitly saved as critical regions in the explicit solution of explicit MPC. However, its application is limited to linear systems (Section I).

This work aims at bringing the benefits of explicit MPC to the SQP method in the application of bilinear MPC. In particular, instead of using an explicit MPC, the explicit solution will play the role of an active set detector in this work. In the rest of this section, we will first introduce a novel interlacing horizon splitting scheme, after which the parallelizable parametric nonconvex solver is elaborated. The convergence properties of the proposed solver are studied in Section III-C. An interpretation of the proposed scheme in the dual space is given in Section III-D, followed by a comparison with related works in Section III-E.

#### A. Interlacing horizon splitting reformulation

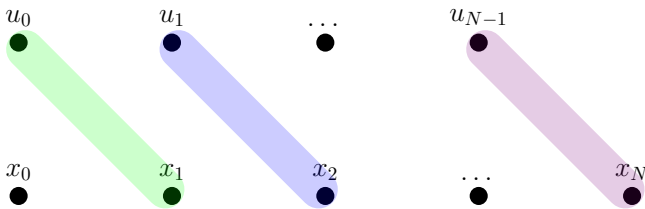


Fig. 1: Visualization of the interlacing horizon splitting.

This section presents the interlacing horizon splitting scheme used later to develop a parallelizable parametric solver to deal with (7). As depicted in Fig. 1, its main idea is to bind the  $k$ -th input  $u_k$  with state  $x_{k+1}$ . To this end, we introduce the

shorthand  $\xi_0 = x_0$  and  $\xi_k = [u_{k-1}^\top, x_k^\top]^\top$  for all  $k \in \mathbb{Z}_1^N$  with associated constraint sets  $\Xi_0 = \{\xi \in \mathbb{R}^{n_x} : \xi_0 = x_0 = x(t)\}$  and

$$\begin{aligned}\Xi_k &= \{\xi \in \mathbb{R}^{n_u+n_x} : \xi \in \mathcal{U} \times \mathcal{X}\}, \quad k \in \mathbb{Z}_1^N \\ &= \{\xi \in \mathbb{R}^{n_u+n_x} : P_\xi \xi \leq p_\xi\}\end{aligned}\quad (8)$$

with  $P_\xi = \text{diag}(P_u, P_x)$  and  $p_\xi = [p_u^\top, p_x^\top]^\top$ . Moreover, we denote the decoupled objective by

$$\begin{aligned}F_0(\xi_0) &= \frac{1}{2}\xi_0^\top Q_0 \xi_0 + q_0^\top \xi_0, \\ F_k(\xi_k) &= \frac{1}{2}\|\xi_k\|_{\text{diag}(R_{k-1}, Q_k)}^2 + [r_{k-1}^\top, q_k^\top]^\top \xi_k, \quad k \in \mathbb{Z}_1^N\end{aligned}$$

and summarize the bilinear dynamics (6) by

$$D_k \xi_k + E_k \xi_{k+1} + (S_{k+1} \xi_{k+1} \otimes \mathbf{I}_{n_x})^\top G_k \xi_k = d_k$$

with coefficients  $d_k = -B_w w_k$ ,  $k \in \mathbb{Z}_0^{N-1}$ ,

$$\begin{aligned}D_0 &= A, \quad D_k = [\mathbf{0}_{n_x \times n_u}, A], \quad k \in \mathbb{Z}_1^{N-1} \\ E_k &= [B, -\mathbf{I}_{n_x}], \quad S_k = [\mathbf{I}_{n_u}, \mathbf{0}_{n_u \times n_x}], \quad k \in \mathbb{Z}_0^{N-1}\end{aligned}$$

and  $G_0 = C$ ,  $G_k = [[\mathbf{0}_{n_x \times n_u}, C_1]^\top, \dots, [\mathbf{0}_{n_x \times n_u}, C_{n_u}]^\top]^\top$  for all  $k \in \mathbb{Z}_1^{N-1}$ . Accordingly, Problem (7) can be rewritten as

$$\min_{\xi} \sum_{k=0}^N F_k(\xi_k) \quad (9a)$$

$$\text{s.t. } (S_{k+1} \xi_{k+1} \otimes \mathbf{I}_{n_x})^\top G_k \xi_k + D_k \xi_k + E_k \xi_{k+1} = d_k \quad | \lambda_k, \quad k \in \mathbb{Z}_0^{N-1} \quad (9b)$$

$$\xi_k \in \Xi_k, \quad | \kappa_k, \quad k \in \mathbb{Z}_0^N \quad (9c)$$

#### B. Proximal-point Lagrangian Based Parallelizable Solver

Based on the interlacing splitting, the linearized proximal-point Lagrangian (5) is used to design the algorithm. On the one hand, it gives zero local duality gap even under the nonlinear/non-convex dynamics [47, Chapter 11.K]. On the other hand, it enables a **convex** QP, accordingly an mpQP, formulation of parallelizable local problems. In particular, for a given primal trajectory  $\bar{\xi}$  (i.e inputs and states) and a dual trajectory  $\{\lambda_k\}_{k=0}^{N-1}$ , the linearized proximal-point Lagrangian of (9) w.r.t the equality constraint (9b) is given by

$$\begin{aligned}\mathcal{L}^\rho(\xi, \lambda, \bar{\xi}) &= \mathcal{L}_0^\rho(\xi_0, \lambda_0, \bar{\xi}_0, \bar{\xi}_1) + \mathcal{L}_N^\rho(\xi_N, \lambda_{N-1}, \bar{\xi}_{N-1}, \bar{\xi}_N) \\ &\quad + \sum_{k=1}^{N-1} \mathcal{L}_k^\rho(\xi_k, \lambda_{k-1}, \lambda_k, \bar{\xi}_{k-1}, \bar{\xi}_k, \bar{\xi}_{k+1})\end{aligned}\quad (10)$$

with

$$\mathcal{L}_0^\rho(\xi_0, \lambda_0, \bar{\xi}_0, \bar{\xi}_1) := F_0(\xi_0) \quad (11a)$$

$$+ \lambda_0^\top [D_0 + (S_1 \bar{\xi}_1 \otimes \mathbf{I}_{n_x})^\top G_0] \xi_0 + \frac{\rho}{2} \|\xi_0 - \bar{\xi}_0\|^2,$$

$$\mathcal{L}_k^\rho(\xi_k, \lambda_{k-1}, \lambda_k, \bar{\xi}_{k-1}, \bar{\xi}_k, \bar{\xi}_{k+1}) := \quad (11b)$$

$$\begin{aligned}F_k(\xi_k) + \lambda_{k-1}^\top (E_{k-1} + \text{mat}(G_{k-1} \bar{\xi}_{k-1}) \cdot S_k) \xi_k \\ + \lambda_k^\top [D_k + (S_{k+1} \bar{\xi}_{k+1} \otimes \mathbf{I}_{n_x})^\top G_k] \xi_k + \frac{\rho}{2} \|\xi_k - \bar{\xi}_k\|^2,\end{aligned}$$

$$\mathcal{L}_N^\rho(\xi_k, \lambda_{N-1}, \bar{\xi}_{N-1}, \bar{\xi}_N) := F_N(\xi_k) + \frac{\rho}{2} \|\xi_N - \bar{\xi}_N\|^2 + \lambda_{N-1}^\top (E_{N-1} + \text{mat}(G_{N-1} \bar{\xi}_{N-1}) \cdot S_N) \xi_N \quad (11c)$$

with  $\rho > 0$  and  $\lambda := [\lambda_0^\top, \lambda_1^\top, \dots, \lambda_{N-1}^\top]^\top$ . From the proximal-point Lagrangian (10), the benefits of the interlacing horizon splitting become clear. Firstly, the problem is decomposed into  $N + 1$  independent subproblems in  $\xi$ . Secondly, each subproblem is a convex mpQP. Furthermore, the use of the proximal-point Lagrangian allows a simplification of  $\mathcal{L}_k^\rho(\cdot)$  to a modified proximal form (Section. IV-A) [48].

If the primal-dual solution  $(\xi^*, \lambda^*)$  of (9) is a regular KKT point, we have solving (9) equivalent to solving

$$\max_{\lambda} \left( - \sum_{k=0}^{N-1} \lambda_k^\top d_k + \min_{\xi} \mathcal{L}^\rho(\xi, \lambda, \bar{\xi} = \xi^*) \right) \quad (12)$$

subject to  $\xi \in \Xi = \Xi_0 \times \dots \times \Xi_N$ .

As  $\mathcal{L}^\rho$  is decoupled in  $\xi$ , our main idea to develop a parallelizable algorithm solving (9) is to design a primal-dual algorithm to solve the dual problem (12) with an iterative update in  $\bar{\xi}$ .

Algorithm 1 outlines the main steps of the proposed algorithm for solving (9). Step 1) deals with decoupled problem (13) in parallel, which has explicit solutions as **convex** mpQPs. Particularly, their solution maps are piece-wise affine functions and can be pre-computed offline (See Section IV-A). Based on the local solutions  $\xi$ , Step 2) evaluates the sensitivities, including the Hessian approximation of the Lagrangian  $\mathcal{L}^0$ , the gradients of the decoupled objective  $F_k$  and the bilinear dynamics residual  $c_k$ . The active Jacobian  $\hat{P}_\xi^k$  are constructed based on the active set at local solutions  $\xi_k$ . The terminal condition is given in Step 3). It is clear that if these termination conditions hold, we have the iterate  $(\xi, \lambda)$  satisfying the first-order optimality condition

$$\mathbf{O}(\epsilon) \in \nabla_{\xi} \mathcal{L}^0(\xi, \lambda, \bar{\xi}) + \mathcal{N}_{\Xi}(\xi)$$

with  $\Xi = \Xi_0 \times \Xi_1 \times \dots \times \Xi_N$  and the primal feasibility condition

$$\|D_k \xi_k + E_k \xi_{k+1} + (S_{k+1} \xi_{k+1} \otimes \mathbf{I}_{n_x})^\top G_k \xi_k\| = \mathbf{O}(\epsilon)$$

for all  $k \in \mathbb{Z}_0^{N-1}$  up to a small error of order  $\mathbf{O}(\epsilon)$ . Step 4) deals with a structured equality-constrained QP (15). To overcome the potential infeasibility caused by the linearization of nonlinear dynamic (6) in constraint (15c), we introduce a decoupled slackness  $s_k$  for each active constraint (15d). This makes QP (15) always feasible regardless of the feasibility of the original problem (9). Note that step 4) is similar to an SQP step, while the mpQP solutions directly generate its active sets. On top of this, the mpQPs in step 1) are also always feasible. Therefore, if one applies Algorithm 1 as an online solver for MPC, the resulting MPC controller is always feasible, i.e., the iteration of Algorithm 1 will never fail before termination, and it is independent of the initial condition  $x(t)$ . This property is desirable in real-world applications because handling infeasibility requires careful design/tuning of a relaxed problem. Furthermore, even for a feasible problem, the standard SQP algorithm may abort due to an incorrect estimate of the active sets. More specifically, if the estimated active set leads to an infeasible QP, the iterations of the SQP

---

### Algorithm 1 Proximal-point Lagrangian Based Online Solver for Bilinear MPC

---

**Input:** an initial guess of  $(\bar{\xi}, \lambda)$ , a stop tolerance  $\epsilon > 0$ , a proximal weight  $\rho$  and a slack penalty  $\mu$

**Repeat:**

- 1) Solve decoupled mpQP problems **sequentially or in parallel**,

$$\min_{\xi_0 \in \Xi_0} \mathcal{L}_0^\rho(\xi_0, \lambda_0, x(t), \bar{\xi}_1), \quad (13a)$$

$$\min_{\xi_k \in \Xi_k} \mathcal{L}_k^\rho(\xi_k, \lambda_{k-1}, \lambda_k, \bar{\xi}_{k-1}, \bar{\xi}_k, \bar{\xi}_{k+1}), \quad k \in \mathbb{Z}_1^{N-1}, \quad (13b)$$

$$\min_{\xi_N \in \Xi_N} \mathcal{L}_N^\rho(\xi_N, \lambda_{N-1}, \bar{\xi}_{N-1}, \bar{\xi}_N). \quad (13c)$$

In all the following steps,  $\xi_k, k \in \mathbb{Z}_0^N$  denote optimal solutions of the above QPs.

- 2) Evaluate sensitivities

$$H \approx \nabla_{\xi\xi} \mathcal{L}^0(\xi, \lambda, \xi), \quad (14a)$$

$$g_k = \nabla F_k(\xi_k) - \nabla \mathcal{L}_k^\rho(\xi, \lambda, \xi), \quad k \in \mathbb{Z}_0^N, \quad (14b)$$

$$c_k = D_k \xi_k + E_k \xi_{k+1} + (S_{k+1} \xi_{k+1} \otimes \mathbf{I}_{n_x})^\top G_k \xi_k - d_k, \quad (14c)$$

and the active Jacobian  $\hat{P}_\xi^k$  at local solution  $\xi_k$ . Here, we use simplified notation  $\mathcal{L}_k^\rho(\xi, \lambda, \xi), k \in \mathbb{Z}_0^N$ . for (11).

- 3) Terminate if  $\max_k \|c_k\| \leq \epsilon$  and  $\max_k \rho \|\xi_k - \bar{\xi}_k\| \leq \epsilon$  hold.
- 4) Solve equality-constrained QP

$$\min_{\Delta\xi, s} \frac{1}{2} \Delta\xi^\top H \Delta\xi + \sum_{k=1}^N \{g_k^\top \Delta\xi_k + \mu \|s_k\|^2\} \quad (15a)$$

$$\text{s.t. } \Delta\xi_0 = 0 \quad (15b)$$

$$E_k \Delta\xi_{k+1} + (S_{k+1} \xi_{k+1} \otimes \mathbf{I}_{n_x})^\top G_k \Delta\xi_k + \text{mat}(G_k \xi_k) \cdot S_{k+1} \Delta\xi_{k+1} + c_k + D_k \Delta\xi_k = 0 \mid \lambda_k^{\text{QP}}, \quad k \in \mathbb{Z}_0^{N-1} \quad (15c)$$

$$\hat{P}_\xi^k \Delta\xi_k = s_k, \quad k \in \mathbb{Z}_1^N. \quad (15d)$$

- 5) Update  $\bar{\xi}^+ = \xi + \Delta\xi$  and  $\lambda^+ = \lambda^{\text{QP}}$ .
- 

algorithm will stop. In summary, the interlacing horizon splitting scheme enables the mpQPs formulation. The proposed algorithm thereby iteratively calls the mpQP solutions, and the inputs to the mpQPs are iteratively updated in the SQP-style step 4).

**Remark 1** *As discussed above, the proposed solver is always feasible even when the initial state makes the NMPC problem infeasible. This property is also observed in the compositions of the projection operator [49], [50], whose convergence to a point pair that are closest to all the sets is proved. However, in a nonconvex setup, the property of the convergent results is unclear and remains open for future research.*

**Remark 2** *In Section II, we discussed that the proposed algorithm only considers process noise in a certainty equivalent form, assuming that the nominal process noise is available.*

However, the proposed algorithm is not able to provide an efficient solution to the robust NMPC problem, which is a challenging nonconvex problem that requires further investigation in the future. However, it is worth noting that the proposed algorithm still has practical benefits even in the current setup. For instance, in building control, even if weather forecasts are available, the actual weather may not align with the nominal forecast, causing the states of the building to evolve into an initial state that renders the NMPC infeasible. The property that the proposed algorithm remains feasible at all times ensures that the building's operations continue uninterrupted.

### C. Local Convergence Property

This section shows that the proposed Algorithm 1 asymptotically converges to the local optimal solution of (7) at a quadratic rate. The logic behind the constructive proof follows two facts: local mpQPs (13) have a Lipschitz-continuous solution map, and the coupled QP (15) is equivalent to a Newton-type method. To this end, we introduce the following lemma first to establish the quadratic contraction of the solution of (15).

**Lemma 1** *Let the KKT point  $(\xi^*, \lambda^*)$  of Problem (9) be regular such that the solution  $\xi^*$  is a local minimizer. Moreover, let Algorithm 1 be initialized locally<sup>2</sup>, whose Hessian evaluation  $H$  and parameter  $\mu$  in (15) satisfy*

$$H = \nabla_{\xi\xi} \mathcal{L}^0(\xi, \lambda, \xi) + \mathbf{O}(\|\xi - \bar{\xi}\|) \quad \text{and} \quad \frac{1}{\mu} \leq \mathbf{O}(\|\xi - \bar{\xi}\|), \quad (16)$$

respectively for every iterate  $(\xi, \bar{\xi})$ . Then, there exists  $\alpha > 0$ , the solution to (15) locally satisfies,

$$\left\| \bar{\xi}^+ - \xi^* \right\| \leq \alpha \|\xi - \xi^*\|^2, \quad \|\lambda^+ - \lambda^*\| \leq \alpha \|\xi - \xi^*\|^2. \quad (17)$$

*Proof:* Based on the definition 1 of regular KKT point, we have that the active sets are not changed locally [51]. Then, the standard analysis of Newton's method gives

$$\begin{aligned} \left\| \begin{bmatrix} \bar{\xi}^+ - \xi^* \\ \lambda^+ - \lambda^* \end{bmatrix} \right\| &\leq \|H - \nabla_{\xi\xi} \mathcal{L}^0(\xi, \lambda, \xi)\| \cdot \mathbf{O}(\|\xi - \xi^*\|) \\ &\quad + \mathbf{O}(\|\xi - \xi^*\|^2) \stackrel{(16)}{\leq} \mathbf{O}(\|\xi - \xi^*\|^2) \end{aligned}$$

as discussed in [40, Chapter 3.3], which concludes the proof. ■

Intuitively speaking, this lemma states that if the iterates given by (13) are close to the optimal solution to (7), then the distance between iterate given by (15) and the optimal solution contract quadratically. The following theorem intends to show that this quadratic contraction holds even when we consider the iterates given by (13). Based on condition (16), we have that there exists a constant  $\alpha > 0$  such that the local quadratic contraction (17) holds. Then, we define by  $\xi^+$  the solution

of (13) based on the updated primal-dual iterate  $(\bar{\xi}^+, \lambda^+)$  such that we can summarize the local convergence result as follows.

**Theorem 1** *Let all assumptions in Lemma 1 be satisfied. The iterates  $\xi$  of Algorithm 1 locally converges to the local minimizer  $\xi^*$  of Problem (9) with quadratic rate.*

*Proof:* As discussed in Section I-B, we have the local solution map  $\xi^*(\bar{\xi}, \lambda)$  of convex mpQPs (13) are Lipschitz continuous such that we have

$$\left\| \xi^*(\bar{\xi}^+, \lambda^+) - \xi^*(\xi^*, \lambda^*) \right\| = \|\xi^+ - \xi^*\| \leq \eta \left\| \begin{bmatrix} \bar{\xi}^+ - \xi^* \\ \lambda^+ - \lambda^* \end{bmatrix} \right\|$$

with a constant  $\eta > 0$ . Here, we use the fact  $\xi^* = \xi^*(\xi^*, \lambda^*)$ , i.e., if we initialize Algorithm 1 with the optimal solution  $(\xi^*, \lambda^*)$ , the solution of convex mpQPs (13) is equal to the local minimizer  $\xi^*$ . Moreover, iterate  $\xi^+$  is the solution of (13) if one starts the Algorithm 1 with  $(\bar{\xi}^+, \lambda^+)$  as the initial guess. Substituting (17) into the inequality above yields

$$\|\xi^+ - \xi^*\| \leq \alpha \cdot \eta \|\xi - \xi^*\|^2,$$

which is sufficient to establish the local quadratic convergence of iterates  $\xi$  to the local minimizer  $\xi^*$  [40, Chapter 3.3]. ■

It is worth mentioning that the same order of local convergence speed holds in a wide range of second-order algorithms, such as the SQP algorithm [40, Chapter 18.7] and the augmented Lagrangian based alternating direction inexact Newton (ALADIN) method [24]. The theoretical importance of Theorem 1 shows that such convergence rate is still preserved even when another layer of mpQPs is added (i.e., step 1) in Algorithm 1). Therefore, regarding the motivation discussed at the beginning of this Section III, the proposed algorithm not only achieves efficient convergence as the Newton-type algorithm but also achieves an efficient active set detection mechanism via the concept of explicit MPC (i.e., mpQPs). Hence, in comparison with the standard SQP, the detection of active sets via mpQPs makes the proposed algorithm advantageous. Additionally, such a benefit does not significantly increase the iteration cost, which retains a low absolute computational time in practice. This is not the case in other SQP-style extensions, such as the ALADIN method, and we postpone the detailed comparison with ALADIN to Section III-E.

### D. Dual Interpretation

In this part, we would show a different but more intuitive view of the proposed algorithm. The expansion of the first-order information  $g_k$  used in (15) gives

$$g_k = (\text{diag}(R_{k-1}, Q_k) + \rho \mathbf{I}_{n_x + n_u}) \xi_k + [r_{k-1}^\top; q_k^\top]^\top + P_\xi^k \kappa_k.$$

Moreover, recall that  $P_\xi$  is the parameter of the inequality constraints (8), and  $\kappa_k$  are the corresponding dual variables (9c), which are generated by the mpQPs solutions directly. By inspecting the objective function in (15), the quadratic penalty term  $\mu \|s_k\|^2$  and  $g_k$  together recover an augmented Lagrangian defined at  $\xi$ , where the active inequality constraints are dualized. This observation leads to a dual interpretation of the proposed algorithm. In the local problems (13) (i.e., step

<sup>2</sup>The term "local" in the statement throughout this paper means that the initial guess of primal-dual iterates is located within a small neighborhood around the local solution  $(\xi^*, \lambda^*)$ . Hence, the condition (16) is required to be satisfied locally. It is worth mentioning that the assumption of locality is standard and widely used in the local convergence analysis of Newton's type methods [40].

1) in Algorithm 1), the system dynamics are dualized with fixed dual variables (i.e.,  $\{\lambda_k\}$ ) based on the proximal-point Lagrangian. The dual variables to the inequality constraints (i.e.,  $\{\kappa_k\}$ ) are thereby updated. Accordingly, the coupled QP problem (15) (i.e., step 4) in Algorithm 1) updates the dynamics dual with the dual variables to the inequality constraints fixed. Hence, the proposed algorithm can be viewed as an alternating direction method in the dual space. Via the scope of this dual interpretation, the coupled QP (15) is not a relaxed problem, as the augmented Lagrangian is an exact penalty function [40, Chapter 17]. More specifically, due to the local convergence of the dual variables by Theorem 1, the augmented Lagrangian converges to an exact penalty.

With this dual interpretation at hand, we are ready to elaborate on the reasoning behind the use of proximal-point Lagrangian. Firstly, the dual variables model first-order local properties [52], and an aggressive update should therefore be avoided due to such locality. The proximal term (i.e.,  $\frac{\rho}{2}\|\xi_k - \bar{\xi}_k\|^2$  in the local steps (13)) realize this goal. This is important, especially when a good estimate of dual variables is not yet available. Secondly, as dual variables encompass first-order information, linearization is therefore needed. More specifically, the dual to the dynamics is updated to  $\bar{\xi}$  in the coupled QP step. Hence the proximal-point Lagrangian is linearized around  $\bar{\xi}$  in (13). It is worth noting that the design logic similar to the aforementioned one also appears in other nonconvex primal alternating direction methods, such as [53]. To the best of our knowledge, the proposed algorithm is the first algorithm to bring this idea to the dual space.

**Remark 3** *The penalty parameters  $\rho$  in (13) and  $\mu$  in (15) play a crucial role in determining the performance of the algorithm. A larger penalty typically leads to better convergence performance [52], but also results in more iterations and longer computational times. Although the convergence aspect of penalty weight has been extensively studied in the literature [52], [54], it remains unclear how to select a penalty that balances convergence performance and the absolute solution time. We leave this for future investigation.*

### E. Comparison with Related Work

In this part, we will compare the proposed scheme with other related results, particularly the augmented Lagrangian based alternating direction inexact Newton (ALADIN) method. The ALADIN method is also an extension of the SQP algorithm, but it can only handle linear coupling. Hence, auxiliary variables that duplicate the states are introduced to handle the bilinear dynamics. More specifically, ALADIN reformulates the bilinear MPC problem (7) to the following equivalent problem:

$$\min_{x_0, x, u} \sum_{k=0}^{N-1} \ell_k(x_k, u_k) + \ell_N(x_N)$$

subject to:

$$\begin{aligned} x_0 &= x(t), \\ \forall k &\in \{0, 1, \dots, N-1\}, \end{aligned}$$

$$\begin{aligned} z_{k+1} &= Ax_k + Bu_k + \sum_{i=1}^{n_u} C_i x_k [u_k]_i + B_w w_k \\ x_{k+1} &= z_{k+1} \quad | \quad \tilde{\lambda}_k, \\ z_{k+1} &\in \mathcal{X}, \quad x_k \in \mathcal{X}, \quad u_k \in \mathcal{U} \end{aligned} \quad (18a)$$

where auxiliary variables  $z_k$  duplicate the states  $x_k$ , and a linear coupling constraint is thereby introduced in (18a). This is the standard horizon splitting scheme used in other nonlinear MPC algorithms [26]–[29], where inputs and states of the same time step are grouped. This leads to the first advantage of this paper's proposed splitting scheme: no auxiliary variables are introduced, so the problem size remains unchanged. This benefit also leads to a limitation of the proposed algorithm: the proposed splitting scheme requires that the stage cost and the constraints are decoupled between states and inputs. How to overcome this limitation is left for future research.

In the ALADIN algorithm, following  $N$  nonconvex subproblems  $\{\mathcal{P}_k\}_{k=0}^{N-1}$  can be solved in parallel:

$$\forall k \in \{0, 1, \dots, N-2\}$$

$$\begin{aligned} \mathcal{P}_k := \min_{\substack{x_k, u_k \\ z_{k+1}}} & \ell_k(x_k, u_k) + \begin{bmatrix} \tilde{\lambda}_{k-1} \\ -\tilde{\lambda}_k \end{bmatrix}^\top \begin{bmatrix} x_k \\ z_{k+1} \end{bmatrix} \\ & + \frac{\rho}{2} \left\| \begin{bmatrix} x_k \\ z_{k+1} \end{bmatrix} - \begin{bmatrix} \bar{x}_k \\ \bar{z}_{k+1} \end{bmatrix} \right\|^2 \end{aligned}$$

subject to:

$$\begin{aligned} z_{k+1} &= Ax_k + Bu_k + \sum_{i=1}^{n_u} C_i x_k [u_k]_i + B_w w_k \\ z_{k+1} &\in \mathcal{X}, \quad x_k \in \mathcal{X}, \quad u_k \in \mathcal{U} \end{aligned}$$

$$\begin{aligned} \mathcal{P}_{N-1} := \min_{\substack{x_{N-1}, u_{N-1} \\ z_N}} & \ell_k(x_{N-1}, u_{N-1}) + \ell(z_N) + \tilde{\lambda}_{N-1}^\top x_{N-1} \\ & + \frac{\rho}{2} \|x_k - \bar{x}_k\|^2 \end{aligned}$$

subject to:

$$\begin{aligned} z_N &= Ax_{N-1} + Bu_{N-1} + \sum_{i=1}^{n_u} C_i x_k [u_k]_{N-1} + B_w w_{N-1} \\ z_N &\in \mathcal{X}, \quad x_{N-1} \in \mathcal{X}, \quad u_{N-1} \in \mathcal{U} \end{aligned}$$

After the parallel iteration, the ALADIN algorithm applies a relaxed SQP-style step to the reformulated problem (18) in order to update the coupling dual variables  $\{\tilde{\lambda}_k\}$ . The differences between the ALADIN and the proposed scheme now become clear:

- The proposed scheme only needs to solve a convex QP, whose solutions can be precomputed offline via mpQPs. Instead, the ALADIN algorithm has to solve a nonconvex problem online. Even though these nonconvex subproblems can be computed in parallel, the resulting computational cost per iteration is still significantly higher than the proposed scheme.
- The proposed scheme directly handles the nonlinear coupling (i.e., the bilinear dynamics) and therefore does

not need to introduce auxiliary variables to duplicate the states. As a result, the SQP step used in the proposed scheme solves a smaller problem than the one solved in the ALADIN.

Due to the use of an SQP-style update and the use of augmented Lagrangian methods [24, Section 4], the proposed algorithm and the ALADIN algorithm have a certain similarity. However, their focus during the algorithm design is different. ALADIN focuses more on allocating the computational complexity, while the proposed algorithm aims at efficient iteration with a good active detection scheme. That is why the ALADIN tends to handle the non-convexity directly, as this can be handled by different computational nodes. On the contrary, the proposed algorithm is customized to bilinear MPC to have a lower computational cost per iteration and fewer decision variables. In summary, even though both the ALADIN and the proposed scheme can be viewed as extensions to the SQP algorithm, the ALADIN is more tailored for distributed computation. The proposed scheme is instead tailored for efficient online computation. Finally, we wrap up this section by summarizing the benefits of the proposed scheme as follows:

- It brings the efficiency of explicit MPC into an NMPC setup. Integrating the explicit mpQP solution provides two benefits: it returns an accurate primal solution when good estimates of the dual variables  $\{\lambda_k\}$  are given, and it significantly improves the real-time efficiency by providing the active set estimation.
- It retains the SQP structure. This not only preserves the convergence rate of the SQP algorithm but also makes the proposed algorithm compatible with any existing acceleration strategy developed for real-time SQP, such as warm-start.
- It enjoys high computational efficiency even without parallelization. With proper implementation, this efficiency may be further improved by parallelization for some processors.

**Remark 4** *It is worth mentioning that the proposed scheme reduces to an ALADIN algorithm when the dynamics are linear (i.e.,  $C_k = 0$ ). In this case, the resulting algorithm is similar to the one studied in [26], where global convergence is also guaranteed [55].*

**Remark 5** *The alternating direction method of multipliers (ADMM) can be used to solve a bi-convex optimization problem [19, Chapter 9.2]. With the standard formulation given in [19, Chapter 9.2], the ADMM algorithm needs to solve a nonconvex QP problem in each iteration, which is proved to be NP-hard even for the calculation of a local minimizer [56], [57]. The resulting computational cost per iteration is significantly higher, and such ADMM formulation is, therefore not suitable for our comparison. Suppose the proposed interlacing horizon splitting scheme is applied instead. In that case, the resulting ADMM algorithm gets rid of the solution of a nonconvex QP, which is also one contribution of this work. However, the ADMM algorithm is still not suitable for comparison. On the one hand, a*

*convergence guarantee exists only when there is no state constraint, which is undesirable in MPC applications. Based on our test on the numerical example given in the following Section V-A, we did not observe the convergence of the ADMM after 3000 iterations (equivalently 1 minute in absolute time). On the other hand, the bilinear dynamics are squared in the augmented Lagrangian. The resulting problem is no longer an mpQP, and cannot be precomputed offline. As a result, multiple inequality constraint QPs are required to be solved in each iteration, leading to a much higher computational cost. Finally, even though we did not observe convergence in our numerical study, if it happens to converge for some specific cases, the convergence rate of a nonconvex ADMM algorithm is at most sublinear [58]. Therefore, it requires more iterations and accordingly, more computational time to converge.*

#### IV. IMPLEMENTATION DETAILS

This section elaborates on the implementation details of Algorithm 1 with a particular emphasis on run-time aspects and a limited memory requirement. Here, the implementation of Steps 3) and 5) turns out to be straightforward, such that we focus on the implementation of Steps 1), 2), and 4).

##### A. mpQP Subproblems

We summarize the local mpQPs (13) into a uniform form

$$\mathbb{P}(\theta_k) : \min_{\xi_k \in \Xi_k} \frac{1}{2} \xi_k^\top \mathcal{Q}_k \xi_k + \theta_k^\top \xi_k \quad (20)$$

with parametric inputs  $\theta_k \in \mathbb{R}^{n_x + n_u}$  and coefficient matrices  $\mathcal{Q}_k = \text{diag}(R_{k-1}, Q_k) + \rho \mathbf{I}_{n_x + n_u}$  for all  $k \in \mathbb{Z}_1^N$ . Here, the first problem is omitted as its solution is fixed by  $\xi_0 = x(t)$  due to the initial state constraint enforced by  $\Xi_0$ . Based on the formulation of  $\mathcal{L}_k^\rho$ , we can work out the explicit form of  $\theta_k$  as follows,

$$\theta_k = [r_{k-1}^\top; q_k^\top]^\top + (E_{k-1} + \text{mat}(G_{k-1} \bar{\xi}_{k-1}) \cdot S_k)^\top \lambda_{k-1} + [D_k + (S_{k+1} \bar{\xi}_{k+1} \otimes \mathbf{I}_{n_x})^\top G_k]^\top \lambda_k - \rho \bar{\xi}_k \quad (21a)$$

$$\theta_N = [r_{N-1}^\top; q_N^\top]^\top - \rho \bar{\xi}_N + (E_{N-1} + \text{mat}(G_{N-1} \bar{\xi}_{N-1}) \cdot S_N)^\top \lambda_{N-1}. \quad (21b)$$

Evaluating these parameters only requires matrix-vector multiplications such that the complexity is  $\mathcal{O}(N \cdot (n_x + n_u)^2)$ . In this paper, we use the enumeration-based multi-parametric QP algorithm from [59] for generating solution maps  $\xi_k^* : \mathbb{R}^{n_x + n_u} \rightarrow \mathbb{R}^{n_x + n_u}$  of (20). The complexity of pre-processing the small-scale QPs (20) depends on the number of critical regions  $N_{R,k}$  over which the PWA optimizers  $\xi_k^*(\cdot)$  are defined [38]. Here, we assume that each parametric QP is post-processed, off-line, to obtain binary search trees [60] in  $\mathcal{O}(N_{R,k}^2)$  time. Once the trees are constructed, they provide for a fast evaluation of the solution maps in (20) in time that is logarithmic in the number of regions, thus establishing the  $\mathcal{O}(\sum_k \log_2(N_{R,k}))$  on-line computational bound. The memory requirements are directly proportional to the number of critical regions, with each region represented by a finite number of affine half-spaces. Finally, it is worth mentioning that if  $Q_i = Q_j$ ,  $q_i = q_j$ ,  $R_i = R_j$ ,  $r_i = r_j$ ,  $\forall i \neq j$ ,



then the  $i$ -th mpQP subproblems (20) is identical to the  $j$ -th one, and one mpQP solution can therefore serve for two subproblems. Identical subproblems happen in many MPC applications as the stage cost are usually fixed throughout the prediction horizon.

## B. Sensitivities Evaluation

Step 2) of Algorithm 1 evaluates the sensitivities  $g_k$ ,  $c_k$ ,  $\hat{P}_\xi^k$  and  $H$ . As we consider the quadratic cost, the gradients  $g_k$  can be easily evaluated with analytical form. Moreover, the primal feasibility residual  $c_k$  and active Jacobian  $\hat{P}_\xi^k$  are also straightforward. Therefore, we focus on the computation of the Hessian matrix  $H$  in this subsection.

As we used an interlacing horizon splitting scheme, the exact Hessian  $\nabla_{\xi\xi} L^0(\xi, \lambda, \xi)$  is not block diagonal with respect to each  $y_k$  but the banded block diagonal. However, as the off-diagonal blocks only involve the bilinear dynamics, we can work out each block analytically as follows:

$$\nabla_{\xi\xi} L^0(y, \lambda, y) = \begin{bmatrix} \mathcal{Q}_0 & S_{0,1} & & & \\ S_{1,0} & \mathcal{Q}_1 & S_{1,2} & & \\ & \ddots & \ddots & \ddots & \\ & & & S_{N-1,N} & \mathcal{Q}_N \end{bmatrix}$$

with blocks

$$\begin{aligned} S_{0,1} &= S_{1,0}^\top = [C_1^\top \lambda_0, \dots, C_{n_u}^\top \lambda_0, \mathbf{0}_{n_x \times n_x}] \in \mathbb{R}^{n_x \times (n_u + n_x)} \\ S_{k,k+1} &= S_{k,k+1}^\top \\ &= \begin{bmatrix} C_1^\top \lambda_k & \dots & C_{n_u}^\top \lambda_k & \mathbf{0}_{n_x \times n_x} \\ \mathbf{0}_{n_u} & \dots & \mathbf{0}_{n_u} & \mathbf{0}_{n_u \times n_x} \end{bmatrix} \in \mathbb{R}^{(n_x + n_u) \times (n_x + n_u)} \end{aligned}$$

for all  $k \in \mathbb{Z}_1^{N-1}$ . It is clear that evaluating the exact Hessian is equivalent to evaluating  $C_i^\top \lambda_k$  for all  $i \in \mathbb{Z}_1^{n_u}$  and  $k \in \mathbb{Z}_0^{N-1}$ . Therefore, its computational complexity is only  $\mathcal{O}(Nn_un_x^2)$ . In practice, some heuristics can be adopted to achieve better numerical robustness on the convergence performance of Algorithm 1 such as enforcing  $H \approx \nabla_{\xi\xi} L^0 \succ 0$  by adding a regularization term, i.e.,  $H = \nabla_{\xi\xi} L^0 + \sigma \mathbf{I}$  with  $\sigma \geq 0$  [61].

## C. Coupled QP

The coupled QP (15) has no inequality constraints such that solving (15) is equivalent to solving linear equations defined by the KKT system:

$$\underbrace{\begin{bmatrix} H + \rho \hat{P}_\xi^\top \hat{P}_\xi & J^\top \\ J & \end{bmatrix}}_{\mathcal{H}} \underbrace{\begin{bmatrix} \Delta \xi \\ \lambda^{\text{QP}} \end{bmatrix}}_w = \begin{bmatrix} -g \\ -c \end{bmatrix} \quad (22)$$

with

$$J = \begin{bmatrix} \tilde{D}_0 & \tilde{E}_0 & & & \\ & \tilde{D}_1 & \tilde{E}_1 & & \\ & & \ddots & \ddots & \\ & & & \tilde{D}_N & \tilde{E}_N \end{bmatrix}, \quad \begin{aligned} \hat{P}_\xi &= \text{diag}(\hat{P}_\xi^1, \dots, \hat{P}_\xi^N) \\ g &= [g_0^\top, g_1^\top, \dots, g_N^\top]^\top \\ c &= [c_0^\top, c_1^\top, \dots, c_N^\top]^\top \end{aligned}$$

and for all  $k \in \mathbb{Z}_0^N$ ,

$$\begin{aligned} \tilde{D}_k &= D_k + (S_{k+1} \bar{\xi}_{k+1} \otimes \mathbf{I}_{n_x})^\top G_k, \\ \tilde{E}_k &= E_k + \text{mat}(G_k z_k) \cdot S_{k+1}. \end{aligned}$$

If we rearrange the KKT matrix  $\mathcal{H}$  by resorting  $w$  as

$$(\Delta \xi_0, \lambda_0^{\text{QP}}, \Delta \xi_1, \lambda_1^{\text{QP}}, \dots, \Delta \xi_{N-1}, \lambda_{N-1}^{\text{QP}}, \Delta \xi_N),$$

a tri-blocked-diagonal sparsity pattern appears in the KKT matrix  $\mathcal{H}$ , such that the Schur complement based back-forward sweeps can be used to solve the linear equation efficiently. To better illustrate this idea, we consider  $N = 2$  such that the resulting rearranged KKT system is

$$\left[ \begin{array}{ccc|cc} \mathcal{Q}_0 & \tilde{D}_0^\top & S_{0,1} & & \\ \tilde{D}_0 & & \tilde{E}_0 & & \\ S_{1,0} & \tilde{E}_0^\top & \tilde{\mathcal{Q}}_1 & \tilde{D}_1^\top & S_{1,2} \\ & & \tilde{D}_1 & & \tilde{E}_1 \\ \hline & & S_{2,1} & \tilde{E}_1^\top & \tilde{\mathcal{Q}}_2 \end{array} \right] \begin{bmatrix} \Delta \xi_0 \\ \lambda_0^{\text{QP}} \\ \Delta \xi_1 \\ \lambda_1^{\text{QP}} \\ \Delta \xi_2 \end{bmatrix} = \begin{bmatrix} -g_0 \\ -c_0 \\ -g_1 \\ -c_1 \\ -g_2 \end{bmatrix}$$

with  $\tilde{\mathcal{Q}}_k = \text{diag}(R_{k-1}, Q_k) + \mu(\hat{P}_\xi^k)^\top \hat{P}_\xi^k$ . We start the backward sweep by considering the whole KKT matrix as a 2x2 block matrix. Then, applying the Schur complement with respect to the lower left block  $\tilde{\mathcal{Q}}_2$  yields a reduced KKT matrix

$$\left[ \begin{array}{cc|c} \mathcal{Q}_0 & \tilde{D}_0^\top & S_{0,1} \\ \tilde{D}_0 & & \tilde{E}_0 \\ \hline S_{1,0} & \tilde{E}_0^\top & \tilde{\mathcal{Q}}_1 + S_{1,2} \tilde{\mathcal{Q}}_2^{-1} S_{2,1} \quad \tilde{D}_1^\top + S_{1,2} \tilde{\mathcal{Q}}_2^{-1} \tilde{E}_1^\top \\ & & \tilde{D}_1 + \tilde{E}_1 \tilde{\mathcal{Q}}_2^{-1} S_{2,1} \quad \tilde{E}_1 \tilde{\mathcal{Q}}_2^{-1} \tilde{E}_1^\top \end{array} \right]$$

Applying the Schur complement once more results in a reduced KKT system with respect to only  $(\Delta \xi_0, \lambda_0)$  such that the substitution of the initial condition  $\Delta \xi_0 = 0$  can enable a forward sweep to recover the primal-dual solution  $(\Delta \xi, \lambda)$ . This method has been shown that it is equivalent to the Riccati recursion in dealing with LQR problems [62]. As the update of the right-hand side of the KKT system only requires matrix-vector multiplication, we observe that the computational complexity of this linear solver is dominated by the matrix update (i.e., computation of the Schur complement), which is  $\mathcal{O}(N(n_x + n_u)^3)$ .

## V. NUMERICAL RESULTS

This section studies the proposed algorithm on two bilinear system examples. The proposed algorithm is first compared against other state-of-art solvers on a building control problem running on a laptop computer. The algorithm is then implemented in an embedded microcontroller for speed control of a DC motor. The binary search tree of the mpQP solutions used in the proposed algorithm is generated by the multi-parametric toolbox (MPT 3.0) [63].

### A. Bilinear Building Control

In this part, the proposed algorithm is compared with an efficient optimal control solver *acados* [45] and the ALADIN algorithm, which is implemented by *ALADIN- $\alpha$*  toolbox [64]. The code generation in *acados* is based on the SQP method with exact Hessians and without/with condensing. All the algorithms use the mirror method to regularize the indefinite QP problem [61]. It is worth mentioning that *acados* is highly optimized for MPC, whose linear algebra

subroutine `BLASFEO` [65] and QP solver [66] exploit the structure in MPC. On top of that, a sophisticated, active set detection scheme by exact penalty function is implemented in `acados` [45, Chapter 2.3]. Hence, this comparison can demonstrate the performance of the proposed algorithm.

We considered a multi-zone building model reported in [67] with room indices shown in Figure 2. Due to the space limit, the parameters of the model (i.e.,  $A$ ,  $B$ ,  $B_w$ ,  $C$ , matrices) are included in the supplementary material on GitHub. In this multi-zone building (Figure 2), room 2 is the corridor linking a large warehouse (room 1) and two offices (rooms 3 and 4). An independent HVAC system controls the indoor temperature of room 1, while another HVAC controls the temperature of all other rooms. The corresponding control inputs ( $u \in \mathbb{R}^2$ ) are the valve positions in the air handling unit, where the heat transfer between the air and the hot water flowing in the heating coil results in the bilinear term in the system dynamics. As a result, the control inputs can manipulate the supply air temperature in a nonlinear way, which accordingly controls the indoor temperatures. In summary, this is a 15-dimensional model (i.e.,  $x \in \mathbb{R}^{15}$ ) with two-dimensional control inputs, the states include the indoor temperature, wall temperature between two different rooms, wall temperature that stands between a specific room, and outdoor, and supply air temperature control. Process noises are outdoor temperature and solar radiation (i.e.,  $w \in \mathbb{R}^2$ ). In building control, a common practice is to apply certainty equivalence control [68], which uses weather forecast as the nominal disturbance in the MPC formulation. Meanwhile, the building evolves under the actual weather condition that is similar but not identical to the weather forecast. Real-world weather data is used for this numerical study.

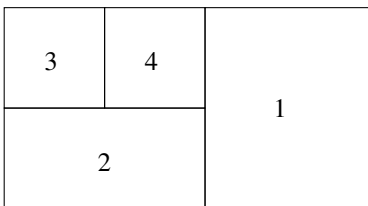


Fig. 2: Schematic diagram of the multi-zone building

Using this approach, 100 Monte-Carlo tests were conducted with recorded weather from tomorrow.io [69] in winter (Fig.3(a)-(b) plot one sampled weather condition). The weather forecast used in the MPC problem is the recorded weather perturbed by zero-mean random noise, while the simulation uses the recorded weather (i.e., the weather forecast curves in Figure 3 (a)-(b)). The prediction horizon is set to 8 with an objective of minimizing energy consumption, whose loss function is

$$\ell_k(x_k, u_k) = u_k^2.$$

To ensure occupant comfort, the indoor temperature is bounded within  $[22, 24]^\circ\text{C}$ . The control input (i.e., fractional valve position) is bounded within  $[0, 1]$ .

**Remark 6** It is possible to define an objective as  $\ell_k(x_k, u_k) = |u_k|$ . The resulting local problem can be

reformulated as a linear program and, thus, also an mpQP. We use a quadratic loss function here to avoid unnecessary confusion.

For this numerical test, all solvers use the same initialization in the first iteration and apply the same warm-start strategy to generate the initialization for the following iterations. In particular, the shifted solution from the last iteration is used to warm-start. The computational time is the sum of the CPU time returned by the solver. The results in this subsection are generated by a laptop with an Intel i7-11800H 16-core processor and 32 GB memory. Meanwhile, as Step 1) in the proposed algorithm 1 is a convex QP, the solution time without using mpQPs is also investigated. In particular, the mpQP solution in this example has 729 critical regions, and the resulting binary search tree is of depth 13 (Section IV-A). The parallelization of the proposed algorithm is done by using `OpenMP`, while the parallelization in `ALADIN` relies on the parallel computing toolbox from `Matlab`. The statistics of the solution time are summarized in Table I, where the maximal solution time indicates the solution time when a good initialization is not available (i.e., cold-start). In contrast, the mean solution time reflects the averaged performance when a good initialization is available.

Above all, Table I shows that the `ALADIN` method is not desirable for fast MPC applications, the need to solve multiple nonconvex problems significantly slows down its speed (Section III-E). We only report the parallelized solution time for `ALADIN`, and the non-parallelized solution time is at least three times slower. Regarding the proposed algorithm, the overhead caused by parallelization pays off only when mpQP solutions are not used. In this case, step 1) in the proposed algorithm requires solving a QP whose computational cost is significantly higher than calling the mpQP solution. Thus, it is easier to improve performance by parallelization in this case. On the contrary, as calling mpQP solution is already computationally highly efficient, improving performance by parallelization may require more involved code design, such as caching. We believe that is the reason why the use of `OpenMP` does not accelerate the mpQP-based implementation in this case. Note that this observation does not negate the benefit of parallelization. On the one hand, the efficiency of parallelization depends on the computing unit and the compiler. The use of `OpenMP` and a general purpose Intel CPU in this numerical example may not be the most efficient implementation. On the other hand, constructing the mpQP solutions may not be computationally affordable for large-scale systems. If solving convex QPs online is needed, then performance improvement is easy to achieve by parallelization, which is also justified by this numerical test.

The comparison with `acados` is shown in Table I. When the proposed scheme uses mpQPs solution without parallelization, its maximal solution time is on average 71% faster than that of `acados` without condensing. Regarding the mean solution time, even though the proposed scheme is faster than `acados` without condensing by 17%, but it is 48% slower than the mean solution of `acados` with condensing. Note that both the QP solver and the linear algebra routine

in *acados* is highly optimized for NMPC; the results in Table I show that a tailored SQP solver is highly efficient when a good initialization is available. As Step 4) in the proposed scheme is similar to an SQP iteration, the proposed algorithm also shows comparable performance to a tailored SQP solver for warm-started iterations. This computational efficiency aligns with our discussion in Section III. On top of this benefit, the proposed scheme performs much better when a good initialization is unavailable. This justifies the use of mpQP solutions, which improves the detection of the active sets. In summary, this numerical study proves the efficiency and efficacy of the proposed algorithm, and it also suggests a further possibility of performance improvement by a more sophisticated combination of the proposed scheme and a tailored SQP solver; we leave this for a future study.

Besides the observation given in Table I, we also observe that the proposed algorithm is more robust to the choice of initialization strategy. More specifically, if the initialization is only partially warm-started by setting the predictive input sequence to  $\mathbf{0}$  (i.e., cold-start inputs but warm-start all the other variables), both *acados* and ALADIN will return NaN during the simulation for all the Monte-Carlo tests. On the contrary, the proposed scheme will always converge even when all the decision variables are initialized by  $\mathbf{0}$ . This observation aligns with the motivation of the proposed algorithm and justifies the benefit of Step 1) in the proposed algorithm. Meanwhile, this robustness might be beneficial in some applications. For example, set-point change in tracking control makes initialization more challenging.

Last but not least, the property that the proposed algorithm is feasible even with an infeasible initial state is useful in practice, which is typically the case in building control. Due to the uncertain occupant behavior, such as opening the window, the indoor climate can be significantly perturbed, resulting in an infeasible initial state for the MPC problem. Consider a case where the occupant opens the window to bring in the fresh air when he first arrives in room 1 at 10:00 A.M., this move causes a sudden drop in indoor temperature as shown in Fig. 3 (c). Such sudden temperature drop causes infeasibility, which leads to the failures of the *acados* solver. However, the proposed algorithm can still give reasonable control inputs and quickly recovers the indoor temperature to a comfortable level.

### B. Bilinear DC Motor Control with a C2000 Microcontroller

Next, the proposed algorithm is deployed on an embedded system, a Texas Instruments C2000 LaunchPad XL F28379D, to control the speed of a field-controlled DC motor. The dynamics of the field-controlled DC motor are bilinear,

$$\begin{aligned} \frac{dx_1}{dt} &= -\frac{R_a}{L_a}x_1 - \frac{K_m}{L_a}x_2u + \frac{V_s}{L_a}, \\ \frac{dx_2}{dt} &= -\frac{B}{J}x_2 + \frac{K_m}{J}x_1u - \frac{T_e}{J}, \end{aligned}$$

where states  $x_1$ ,  $x_2$  are, respectively, armature current and angular velocity, and the control input  $u$  is field current.  $V_s$  and  $T_e$  define the external torque, respectively, which are chosen as

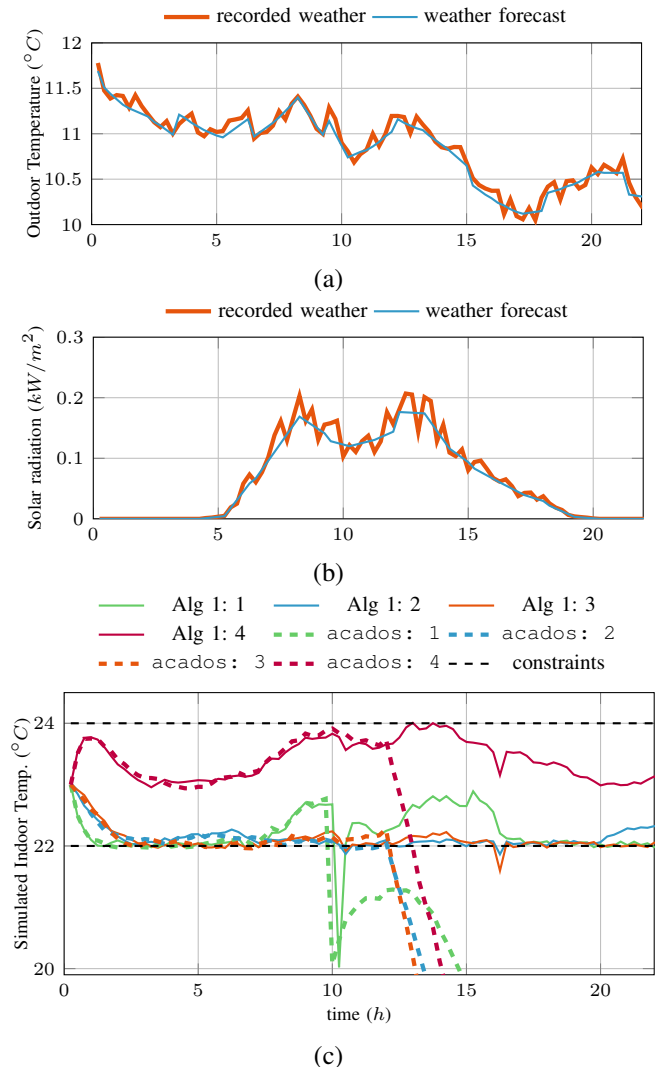


Fig. 3: Case study of sudden indoor temperature drop: (a) - (b): a sample of the recorded and forecast weather condition, (a) outdoor temperature, (b) solar radiation. The forecast is used as the nominal weather in the NMPC problem, while recorded weather is used for the simulation of building dynamics. (c): simulation of indoor temperature of different rooms (room index depicted in Figure 2).

60 V and 0 Nm for this experiment. The remaining parameters are identified on a real field-controlled DC motor (Fig. 5b) as shown in Table II.

We first provide some background on the motor behavior to gain insight into the NMPC solutions. Typically, the armature current dynamics  $x_1$  are much faster than the mechanical dynamics, so it is useful to consider the motor behavior after the current dynamics have decayed. The torque-speed curves of the motor are shown in Fig. 4 (a) for various field currents. This is the electrical torque produced by the motor for the given speed and field current. If the mechanical torques (drag+external) match this torque at a given speed, it is an equilibrium point. For example, we may observe that the no-load speed for this motor (without drag) is 87 rad/s at the full 3 A field current. The typical operating region of this

**TABLE I:** Statistics of the solution time at different tolerance (the entries of the top two performers in each row are stressed by boldface black and boldface blue respectively)

method		Algorithm 1				acados		ALADIN
tol	sol time (ms)	parallel (mpQPs)				condensing		parallel
		yes (yes)	no (yes)	yes (no)	no (no)	yes	no	
$10^{-4}$	max	<b>8.404</b>	<b>6.304</b>	8.926	12.272	19.858	10.711	2553
	mean	0.5280	<b>0.4835</b>	0.6583	0.8751	<b>0.3671</b>	0.7931	897.4
$10^{-5}$	max	<b>8.838</b>	<b>6.629</b>	9.371	12.905	21.755	11.804	2859
	mean	0.6864	<b>0.6286</b>	0.8558	1.137	<b>0.4020</b>	0.8687	953.2
$10^{-6}$	max	<b>8.968</b>	<b>6.753</b>	9.593	13.510	22.532	12.251	3369
	mean	0.7814	<b>0.7156</b>	0.9743	1.295	<b>0.4163</b>	0.8995	1053

**TABLE II:** Parameters of the Field Controlled DC Motor

Parameter	Variable	Value
Armature Resistance	$R_a$	10 [ohm]
Armature Inductance	$L_a$	60 [mH]
Motor Constant	$K_m$	0.2297 [V (A rad/s) $^{-1}$ ]
Damping Ratio	B	0.0024 [Nm(rad/s) $^{-1}$ ]
Inertia	J	0.008949 [kg m $^2$ ]

type of motor is at speeds higher than the full-field line, roughly  $[80, 180]$  rad/s for low torques. Operation below this speed is undesirable because the armature currents exceed the 3 A continuous thermal limits regardless of the field current selected. This is shown in the current-speed curve (Fig. 4 (b)), where armature current is plotted as a function of speed for different field currents (i.e., control input). Hence, the curves below the red dashed line in Figure 4 (b) also show the set of desired operating points that allows long-term operation without overheating.

The continuous dynamics are discretized by the Euler method with a sampling time of 10ms. The prediction horizon is set to  $3^3$  with a convergence tolerance at  $10^{-4}$ . The MPC controller conducts speed control, which tracks a reference speed  $\omega_{\text{ref}}$ . This motor operates around  $[80, 180]$  rad/s, and for most reference torque/speed combinations within this range, there are two possible field current solutions as shown by overlapping lines in Fig. 4 (a). The low field current (i.e., control input) solution results in a higher armature current, usually above the 3 A limit (Fig. 4 (b)). Long-term operation on this equilibrium point will result in armature overheating even though it tracks the reference speed. However, to have an agile motor response, the armature current should be able to operate above 3 A for short intervals. Therefore, we do not enforce a constraint on the armature current, while the field current (i.e., control input) is bounded within  $[1, 3]$  A.

In this control setup, the desired operating point has an

<sup>3</sup>The prediction horizon is set based on some recent results with **commercial** solver from ODYS [70], [71], where they use the same C2000-series hardware to deploy MPC on a synchronous machine. In their setup, the input constraints are neglected, the prediction horizon is two, and the linearized model is used instead of the nonlinear model.

armature current lower than 3 A, which corresponds to the higher of two viable field currents (Fig 4). A proper choice of the loss function can help the solver to converge to the desired operating point. First of all, it is not desirable to use the speed regulating stage cost  $\ell(x, u) = ([x_k]_2 - \omega_{\text{ref}})^2$ , as the solver tends to select the lower field current command which will overheat the armature. This is particularly the case in the presence of noise based on our observation of different hardware in the loop simulations. We suspect that this can be explained by the torque-speed curve (Fig. 4 (a)). When operating with a low field current, the curves are relatively flat, and a slight change in the field current can lead to a rapid change in speed. This implies that the solver can give better local convergence behavior in this region, so the solver tends to converge to this undesired operating point. To avoid this issue and push the solution to the preferred operating point, we offer a reference armature current and field current, whose steady state solution has an explicit form by substituting  $\omega_{\text{ref}}$  into the system dynamics. The stage cost is designed to

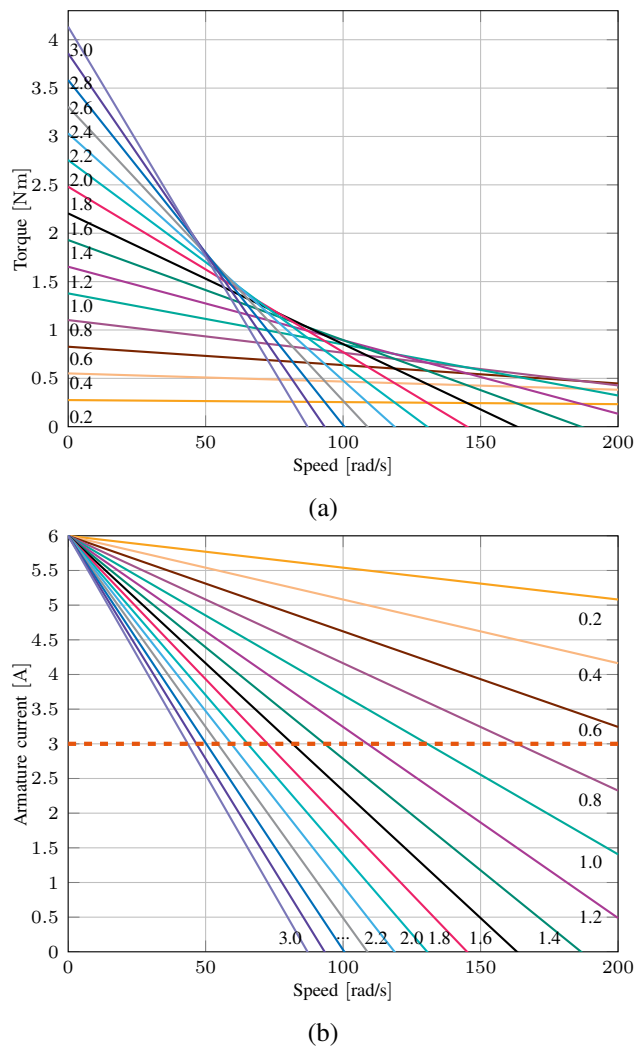
$$\ell(x, u) = 20([x_k]_1 - I_{\text{ref}})^2 + ([x_k]_2 - \omega_{\text{ref}})^2 + 10(u - u_{\text{ref}})^2,$$

where  $I_{\text{ref}}$  and  $u_{\text{ref}}$  are reference armature current and reference field current.

The experimental setup is shown in Fig. 5 with the explanation given in its caption. Two experiments are carried out on velocity tracking control, which both track a triangular reference speed that varies between 100 rad/s and 140 rad/s. In the first case, we only have a speed constraint within  $[80, 180]$  rad/s, while this constraint is tightened to  $[110, 180]$  rad/s in the second case. Thus, the speed lower bound is inactive in the first case but active and satisfied in the second case.

To verify the real-time controller performance, we first simulate these experiments using control hardware in the loop methods. Both the controller and a simulated motor run on the same C2000 microcontroller (Fig. 5c). The results are shown in Fig. 6, with the armature current disturbed by white noise to reproduce the switching noise encountered in real-world experiments. In simulations, the NMPC properly executes the speed control, and the speed constraint is satisfied in the second case as expected.

The hardware in-the-loop result above can already justify the efficiency of the proposed solver, a similar setup for



**Fig. 4:** Torque and current curves of the DC motor. Top: Torque-Speed curve at different field currents (indicated in Amps on the left end of each straight line). Bottom: Armature current as a function of speed for different field currents (indicated in Amps on the right end of each straight line). To avoid armature overheating, the armature current should stay below the 3 A thick red dashed line in long term operation.

performance proof is also used in the **commercial** product in [70], [71]. But for the sake of completeness, we carry out the experiment on a real motor. The hardware experimental results are shown in Fig. 7. The measured signals are post-processed with a low-pass zero-phase filter. In this experiment, the proposed algorithm successfully executes the control in real-time with a 10 ms MPC update rate. In particular, the maximum and average execution time of the proposed algorithm in this embedded system are 2.088 ms and 1.764 ms, respectively. Thus, the solver can run up to 500 Hz, which is sufficiently fast regarding the 10 ms sampling time of the targeted system.

However, in our real-world experiments, the tracking performance is somewhat lacking. From Fig. 7, we can observe that the NMPC tries to track the signal, and periodic triangular speed trajectories are recorded with noticeable tracking errors. It is noteworthy that the NMPC satisfies the lower speed limit

in the second test, which justifies the constraints enforced by the NMPC. The reasons will be investigated in future work but may be due to poor estimates of the unmeasured drag torque, other parameter errors, or inaccurate delivery of the current command. In summary, based on the hardware in the loop experiments and the real-world experiments, the efficiency and real-time capability of the proposed algorithm has been proven. Even though it is not the main focus of this work, we believe that there are still a few improvements to be carried out on our experiment platform to exploit the capability of the NMPC control fully, and we leave this for future investigation.

## VI. CONCLUSION

This paper proposes a novel proximal-point Lagrangian based nonconvex solver for bilinear model predictive control. The proposed algorithm combines the ideas of explicit MPC, horizon splitting, and real-time SQP algorithms and a novel horizon splitting scheme is proposed to enable this integration. The numerical efficiency of the proposed algorithm is validated by a building control simulation and an experiment on a real field-controlled DC motor with a TI C2000 LaunchPad XL F28379D microcontroller. Particularly, the latter experiment proves the real-time capability of the proposed algorithm, which successfully solves the NMPC problem in 1.764 ms on average.

## REFERENCES

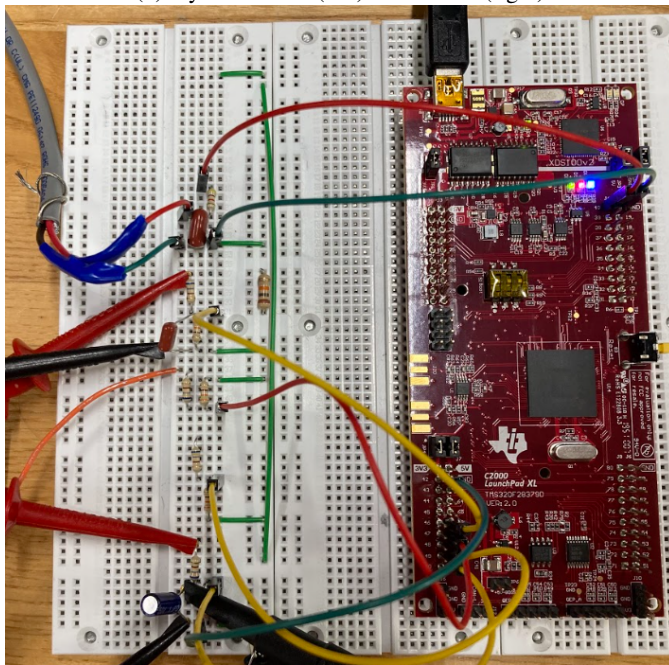
- [1] R. Mohler and R. Rink, "Multivariable bilinear system control," in *Proc. IFAC Symp. Multivariable Control Systems*, vol. 14, 1968.
- [2] R. Rink and R. Mohler, "Completely controllable bilinear systems," *SIAM Journal on Control*, vol. 6, no. 3, pp. 477–486, 1968.
- [3] A. Ruberti and R. Mohler, "Variable structure systems with applications to economics and biology," in *Proceedings of the Second US-Italy Seminar on Variable Structure Systems*, Springer, 1974.
- [4] D. L. Elliott, "Bilinear systems," *Wiley Encyclopedia of Electrical and Electronics Engineering*, 2001.
- [5] W.-H. Steeb, "A note on carleman linearization," *Physics Letters A*, vol. 140, no. 6, pp. 336–338, 1989.
- [6] E. Kaiser, J. N. Kutz, and S. L. Brunton, "Data-driven approximations of dynamical systems operators for control," *The Koopman Operator in Systems and Control*, pp. 197–234, 2020.
- [7] D. Elliott, *Bilinear control systems: matrices in action*, vol. 169. Springer Science & Business Media, 2009.
- [8] W. J. Rugh, *Nonlinear system theory*. Johns Hopkins University Press Baltimore, MD, 1981.
- [9] J. Baillieul, "Geometric methods for nonlinear optimal control problems," *Journal of optimization theory and applications*, vol. 25, no. 4, pp. 519–548, 1978.
- [10] S. M. Rajguru, M. A. Ifediba, and R. D. Rabbitt, "Three-dimensional biomechanical model of benign paroxysmal positional vertigo," *Annals of biomedical engineering*, vol. 32, no. 6, pp. 831–846, 2004.
- [11] D. D'alessandro and M. Dahleh, "Optimal control of two-level quantum systems," *IEEE Transactions on Automatic Control*, vol. 46, no. 6, pp. 866–876, 2001.
- [12] G. Escobar, R. Ortega, H. Sira-Ramirez, J.-P. Vilain, and I. Zein, "An experimental comparison of several nonlinear controllers for power converters," *IEEE Control Systems Magazine*, vol. 19, no. 1, pp. 66–82, 1999.
- [13] H. Sira-Ramirez, "Sliding motions in bilinear switched networks," *IEEE transactions on circuits and systems*, vol. 34, no. 8, pp. 919–933, 1987.
- [14] S. Peitz, "Controlling nonlinear PDEs using low-dimensional bilinear approximations obtained from data," *arXiv preprint arXiv:1801.06419*, 2018.
- [15] J. Haddad and N. Geroliminis, "On the stability of traffic perimeter control in two-region urban cities," *Transportation Research Part B: Methodological*, vol. 46, no. 9, pp. 1159–1176, 2012.



(a) Overall Test setup



(b) Dynamometer (left) and Motor (right)



(c) TI LaunchPad XL and analog I/O filtering

**Fig. 5:** Hardware testing of the proposed MPC algorithm with a real motor controlled by a TI C2000 Delfino 28379D microcontroller. The overall test setup is visible in (a) with the host computer on the right, the power supplies for the armature and field windings in the top left, the motor/dynamometer setup in the bottom left, and the controller in the front left corner of the bench. A detailed view of the motor is shown in (b), with clamp-on current sensors for the armature and field currents. The controller and its associated analog input/output filters are shown in (c).

- [16] M. B. Kane, J. Scruggs, and J. P. Lynch, "Model-predictive control techniques for hydronic systems implemented on wireless sensor and actuator networks," in *2014 American Control Conference*, pp. 3542–3547, IEEE, 2014.
- [17] L. Hetel, M. Defoort, and M. Djemai, "Binary control design for a class of bilinear systems: Application to a multilevel power converter," *IEEE Transactions on Control Systems Technology*, vol. 24, no. 2, pp. 719–726, 2015.
- [18] J. B. Rawlings, D. Q. Mayne, and M. Diehl, *Model predictive control: theory, computation, and design*, vol. 2. Nob Hill Publishing Madison, WI, 2017.
- [19] S. Boyd, N. Parikh, and E. Chu, *Distributed optimization and statistical learning via the alternating direction method of multipliers*. Now Publishers Inc, 2011.
- [20] A. Rantzer, "Dynamic dual decomposition for distributed control," in *2009 American Control Conference*, pp. 884–888, IEEE, 2009.
- [21] I. Necoara and J. Suykens, "Interior-point lagrangian decomposition method for separable convex optimization," *Journal of Optimization Theory and Applications*, vol. 143, no. 3, pp. 567–588, 2009.
- [22] J. V. Frasch, S. Sager, and M. Diehl, "A parallel quadratic programming method for dynamic optimization problems," *Mathematical programming computation*, vol. 7, no. 3, pp. 289–329, 2015.
- [23] S. Richter, M. Morari, and C. N. Jones, "Towards computational complexity certification for constrained mpc based on lagrange relaxation and the fast gradient method," in *2011 50th IEEE Conference on Decision and Control and European Control Conference*, pp. 5223–5229, IEEE, 2011.
- [24] B. Houska, J. Frasch, and M. Diehl, "An augmented lagrangian based algorithm for distributed nonconvex optimization," *SIAM Journal on Optimization*, vol. 26, no. 2, pp. 1101–1127, 2016.
- [25] Y. Jiang, J. Oravec, B. Houska, and M. Kvasnica, "Parallel mpc for linear systems with input constraints," *IEEE Transactions on Automatic Control*, 2020.
- [26] Y. Jiang, C. N. Jones, and B. Houska, "A time splitting based real-time iteration scheme for nonlinear mpc," in *2019 IEEE 58th Conference on Decision and Control (CDC)*, pp. 2350–2355, IEEE, 2019.
- [27] F. Laine and C. Tomlin, "Parallelizing LQR computation through endpoint-explicit riccati recursion," in *2019 IEEE 58th Conference on Decision and Control (CDC)*, pp. 1395–1402, IEEE, 2019.
- [28] H. Deng and T. Ohtsuka, "A parallel newton-type method for nonlinear model predictive control," *Automatica*, vol. 109, p. 108560, 2019.
- [29] S. Shin, T. Faulwasser, M. Zanon, and V. M. Zavala, "A parallel decomposition scheme for solving long-horizon optimal control problems," in *2019 IEEE 58th Conference on Decision and Control (CDC)*, pp. 5264–5271, IEEE, 2019.
- [30] I. Nielsen and D. Axehill, "A parallel structure exploiting factorization algorithm with applications to model predictive control," in *2015 54th IEEE Conference on Decision and Control (CDC)*, pp. 3932–3938, IEEE, 2015.
- [31] Y. Wang and S. Boyd, "Fast model predictive control using online optimization," *IEEE Transactions on control systems technology*, vol. 18, no. 2, pp. 267–278, 2009.
- [32] H. J. Ferreau, H. G. Bock, and M. Diehl, "An online active set strategy to overcome the limitations of explicit mpc," *International Journal of Robust and Nonlinear Control: IFAC-Affiliated Journal*, vol. 18, no. 8, pp. 816–830, 2008.
- [33] M. Diehl, H. G. Bock, and J. P. Schlöder, "A real-time iteration scheme for nonlinear optimization in optimal feedback control," *SIAM Journal on control and optimization*, vol. 43, no. 5, pp. 1714–1736, 2005.
- [34] E. Zafriou, "Robust model predictive control of processes with hard constraints," *Computers & Chemical Engineering*, vol. 14, no. 4-5, pp. 359–371, 1990.
- [35] A. Bemporad, M. Morari, V. Dua, and E. N. Pistikopoulos, "The explicit linear quadratic regulator for constrained systems," *Automatica*, vol. 38, no. 1, pp. 3–20, 2002.
- [36] D. M. Raimondo, S. Riverso, C. N. Jones, and M. Morari, "A robust explicit nonlinear mpc controller with input-to-state stability guarantees," *IFAC Proceedings Volumes*, vol. 44, no. 1, pp. 9284–9289, 2011.
- [37] S. Boyd, S. P. Boyd, and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [38] A. Bemporad, M. Morari, V. Dua, and E. N. Pistikopoulos, "The explicit linear quadratic regulator for constrained systems," *Automatica*, vol. 38, pp. 3–20, Jan. 2002.
- [39] F. Borrelli, A. Bemporad, and M. Morari, *Predictive control for linear and hybrid systems*. Cambridge University Press, 2017.
- [40] S. Wright, J. Nocedal, et al., "Numerical optimization," *Springer Science*, vol. 35, no. 67-68, p. 7, 1999.

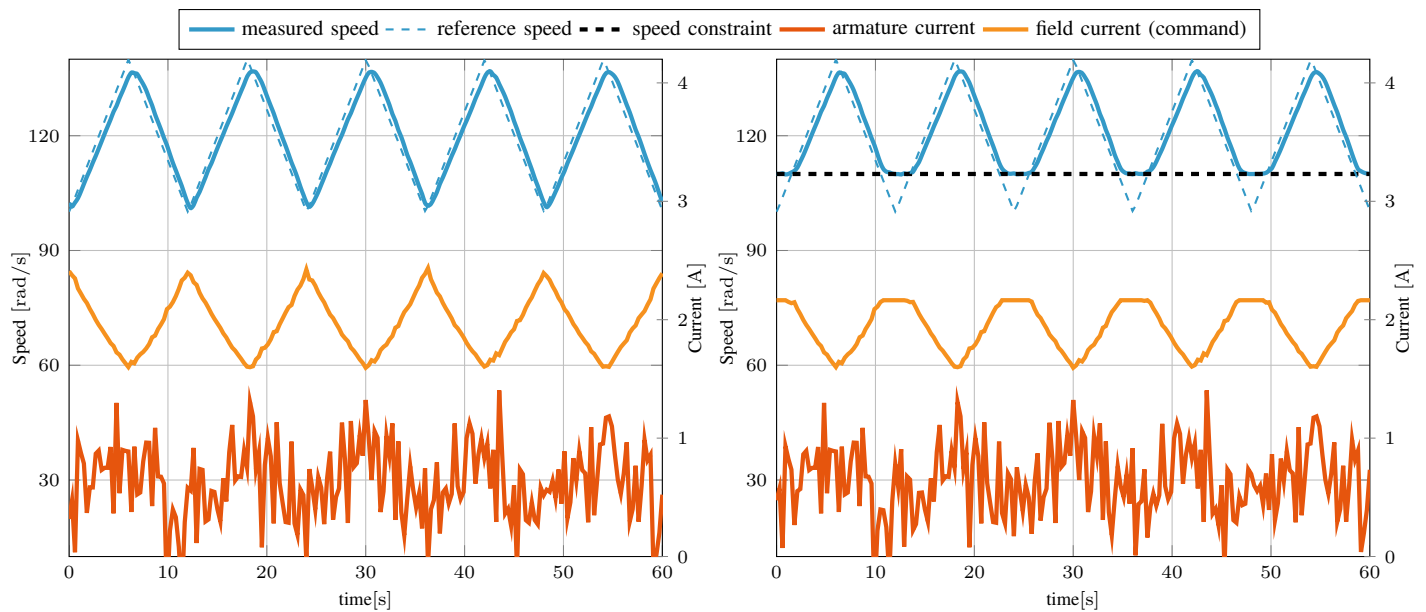


Fig. 6: Hardware in the loop simulation of the motor control with the C2000 microcontroller. Left: speed constraint not active. Right: speed constraint [110, 180] rad/s active (black dashed line). Raw signals are plotted.

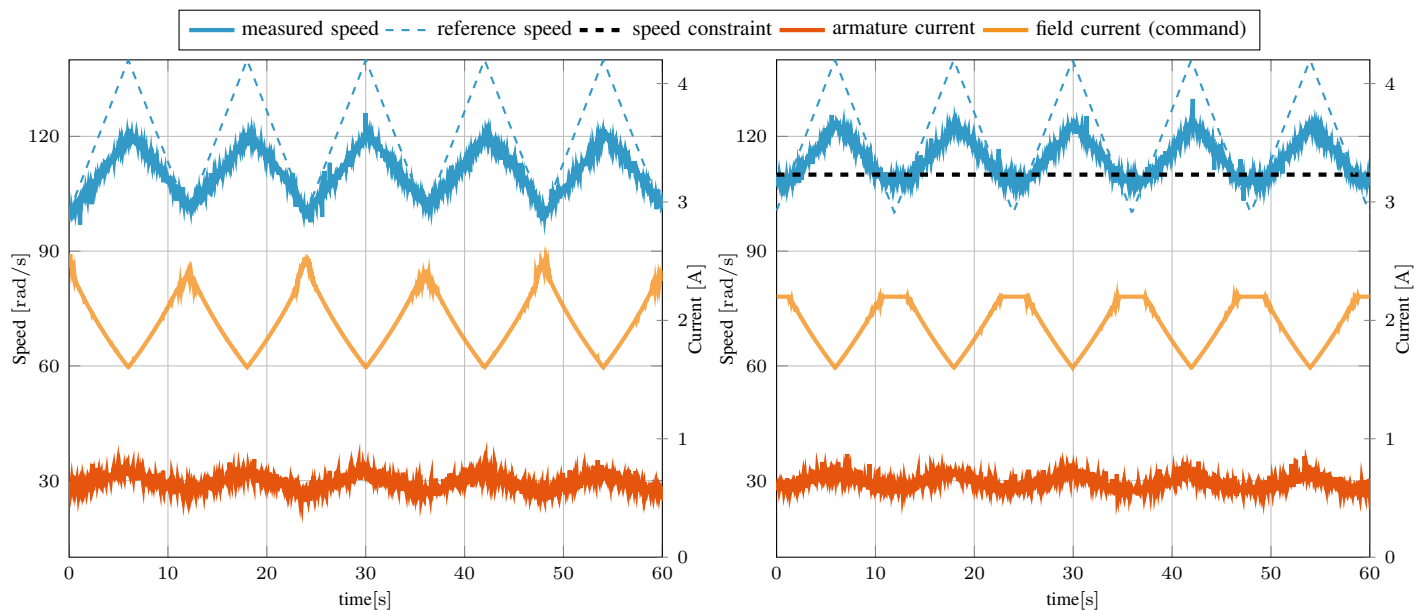
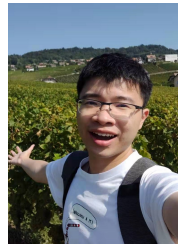


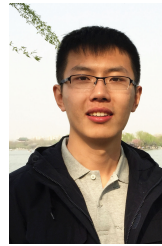
Fig. 7: Hardware testing of the proposed MPC algorithm with a real motor controlled by a C2000 microcontroller. Left: speed constraint not active. Right: speed constraint [110, 180] rad/s active (black dashed line). The measured signals are post-processed with a low-pass zero-phase forward-reverse filter.

- [41] N. G. D. Robinson, "A second derivative sqp method: local convergence," 2009.
- [42] M. R. Hestenes, "Multiplier and gradient methods," *Journal of optimization theory and applications*, vol. 4, no. 5, pp. 303–320, 1969.
- [43] M. J. Powell, "A method for nonlinear constraints in minimization problems," *Optimization*, pp. 283–298, 1969.
- [44] R. T. Rockafellar, "Augmented lagrangians and applications of the proximal point algorithm in convex programming," *Mathematics of operations research*, vol. 1, no. 2, pp. 97–116, 1976.
- [45] R. Verschueren, G. Frison, D. Kouzoupis, J. Frey, N. van Duijkeren, A. Zanelli, B. Novoselnik, T. Albin, R. Quirynen, and M. Diehl, "acados – a modular open-source framework for fast embedded optimal control," *Mathematical Programming Computation*, Oct 2021.
- [46] P. E. Gill, W. Murray, and M. A. Saunders, "Snopt: An sqp algorithm for large-scale constrained optimization," *SIAM review*, vol. 47, no. 1, pp. 99–131, 2005.
- [47] R. T. Rockafellar and R. J.-B. Wets, *Variational analysis*, vol. 317. Springer Science & Business Media, 2009.
- [48] G. Di Pillo and L. Grippo, "A new augmented lagrangian function for inequality constraints in nonlinear programming problems," *Journal of Optimization Theory and Applications*, vol. 36, no. 4, pp. 495–519, 1982.
- [49] S. Alwadani, H. H. Bauschke, J. P. Revalski, and X. Wang, "The difference vectors for convex sets and a resolution of the geometry conjecture," *Open Journal of Mathematical Optimization*, vol. 2, pp. 1–18, 2021.
- [50] S. Alwadani, H. H. Bauschke, J. P. Revalski, and X. Wang, "Resolvents

- and Yosida approximations of displacement mappings of isometries,” *Set-Valued and Variational Analysis*, vol. 29, no. 3, pp. 721–733, 2021.
- [51] B. Houska and Y. Jiang, “Distributed optimization and control with ALADIN,” *Recent Advances in Model Predictive Control: Theory, Algorithms, and Applications*, p. 135–163, 2021.
- [52] R. T. Rockafellar, “Augmented Lagrange multiplier functions and duality in nonconvex programming,” *SIAM Journal on Control*, vol. 12, no. 2, pp. 268–285, 1974.
- [53] J. Bolte, S. Sabach, and M. Teboulle, “Proximal alternating linearized minimization for nonconvex and nonsmooth problems,” *Mathematical Programming*, vol. 146, no. 1, pp. 459–494, 2014.
- [54] D. P. Bertsekas, “Multiplier methods: A survey,” *Automatica*, vol. 12, no. 2, pp. 133–145, 1976.
- [55] B. Houska, D. Kouzoupis, Y. Jiang, and M. Diehl, “Convex optimization with aladin,” *Optimization Online preprint*, 2017.
- [56] B. Contesse, “Une caractérisation complète des minima locaux en programmation quadratique,” 1980.
- [57] A. Forsgren, P. Gill, and W. Murray, “On the identification of local minimizers in inertia-controlling methods for quadratic programming,” *SIAM journal on matrix analysis and applications*, vol. 12, no. 4, pp. 730–746, 1991.
- [58] Y. Wang, W. Yin, and J. Zeng, “Global convergence of admm in nonconvex nonsmooth optimization,” *Journal of Scientific Computing*, vol. 78, no. 1, pp. 29–63, 2019.
- [59] M. Herceg, C. N. Jones, M. Kvasnica, and M. Morari, “Enumeration-based approach to solving parametric linear complementarity problems,” *Automatica*, no. 62, pp. 243–248, 2015.
- [60] P. Tøndel, T. Johansen, and A. Bemporad, “Evaluation of piecewise affine control via binary search tree,” *Automatica*, vol. 39, no. 5, pp. 945–950, 2003.
- [61] R. Verschueren, M. Zanon, R. Quirynen, and M. Diehl, “A sparsity preserving convexification procedure for indefinite quadratic programs arising in direct optimal control,” *SIAM Journal on Optimization*, vol. 27, no. 3, pp. 2085–2109, 2017.
- [62] G. Frison, *Algorithms and methods for high-performance model predictive control*. PhD thesis, 2016.
- [63] M. Kvasnica, P. Grieder, M. Baotić, and M. Morari, “Multi-parametric toolbox (mpt),” in *International Workshop on Hybrid Systems: Computation and Control*, pp. 448–462, Springer, 2004.
- [64] A. Engelmann, Y. Jiang, H. Benner, R. Ou, B. Houska, and T. Faulwasser, “Aladin—an open-source matlab toolbox for distributed non-convex optimization,” *Optimal Control Applications and Methods*, vol. 43, no. 1, pp. 4–22, 2022.
- [65] G. Frison, D. Kouzoupis, T. Sartor, A. Zanelli, and M. Diehl, “Blasfeo: Basic linear algebra subroutines for embedded optimization,” *ACM Transactions on Mathematical Software (TOMS)*, vol. 44, no. 4, pp. 1–30, 2018.
- [66] G. Frison and M. Diehl, “Hpipm: a high-performance quadratic programming framework for model predictive control,” *IFAC-PapersOnLine*, vol. 53, no. 2, pp. 6563–6569, 2020.
- [67] F. Belić, D. Slišković, and Ž. Hocenski, “Detailed thermodynamic modeling of multi-zone buildings with resistive-capacitive method,” *Energies*, vol. 14, no. 21, p. 7051, 2021.
- [68] F. Oldewurtel, C. N. Jones, A. Parisio, and M. Morari, “Stochastic model predictive control for building climate control,” *IEEE Transactions on Control Systems Technology*, vol. 22, no. 3, pp. 1198–1205, 2013.
- [69] “Tomorrow.io.” <https://www.tomorrow.io>. Accessed: 2022-03-01.
- [70] G. Cimini, D. Bernardini, A. Bemporad, and S. Levijoki, “Online model predictive torque control for permanent magnet synchronous motors,” in *2015 IEEE International Conference on Industrial Technology (ICIT)*, pp. 2308–2313, IEEE, 2015.
- [71] G. Cimini, D. Bernardini, S. Levijoki, and A. Bemporad, “Embedded model predictive control with certified real-time optimization for synchronous motors,” *IEEE Transactions on Control Systems Technology*, vol. 29, no. 2, pp. 893–900, 2020.



**Yingzhao Lian** graduated with an MSc in Mechanical Engineering from EPFL in 2018. He specialized in optimization and learning theory. He jointly conducted his master’s thesis with the Automatic Control Lab, EPFL, and ABB Corporate Research, Baden, where he investigated the problem of data-driven model-based optimal control. In August 2018, he joined the Automatic Control Laboratory at EPFL as a Ph.D. student under the supervision of Professor Colin Jones.



**Yuning Jiang** received the B.Sc. degree in electronic engineering from Shandong University, Jinan, China, in 2014, and the Ph.D. degree in information engineering from ShanghaiTech University, Shanghai, China, and the University of Chinese Academy of Sciences, Beijing, China, in 2020. He has ever been a Visiting Scholar with the University of California at Berkeley (UC Berkeley), Berkeley, CA, USA, the University of Freiburg, Freiburg im Breisgau, Germany, and Technische Universität Ilmenau (TU Ilmenau), Ilmenau, Germany, during his Ph.D. study. He is currently a Post-Doctoral Researcher with the Automatic Control Laboratory at the EPFL in Switzerland. His research focuses on learning- and optimization-based policy for operating complex systems such as nonlinear autonomous systems (e.g., autonomous vehicles, robotics and smart buildings), and large-scale multi-agent systems (e.g., power and energy systems, IoT and traffic networks).



**Daniel F. Opila** (SM) received the B.S. and M.S. degrees from the Massachusetts Institute of Technology, Cambridge, MA, USA, in 2002 and 2003, respectively, and the Ph.D. degree from the University of Michigan, Ann Arbor, MI, USA, in 2010. He is currently an Associate Professor of Electrical and Computer Engineering at the United States Naval Academy, Annapolis, MD, USA. He has previously worked in various engineering positions at GE Power Conversion, Ford Motor Company, Orbital Sciences Corporation, and Bose Corporation. He specializes in optimal control of energy systems, including hybrid vehicles, naval power systems, power converters, and renewables.



**Colin N. Jones** has been an Associate Professor in the Automatic Control Laboratory at the EPFL in Switzerland since 2017 and an assistant professor from 2011. He was a Senior Researcher at the Automatic Control Lab at ETH Zurich until 2010 and obtained a Ph.D. in 2005 from the University of Cambridge for his work on polyhedral computational methods for constrained control. Prior to that, he was at the University of British Columbia in Canada, where he took his bachelor and master degrees in Electrical Engineering and Mathematics. He is the author or coauthor of more than 200 publications and was awarded an ERC starting grant to study the optimal control of building networks. His current research interests are in the areas of high-speed predictive control and optimization, as well as the control of green energy generation, distribution and management.